

---

# AWS Neuron

**AWS**

**Jul 02, 2025**



## OVERVIEW

|           |                                  |             |
|-----------|----------------------------------|-------------|
| <b>1</b>  | <b>Overview</b>                  | <b>3</b>    |
| <b>2</b>  | <b>ML Frameworks</b>             | <b>25</b>   |
| <b>3</b>  | <b>NeuronX Distributed (NxD)</b> | <b>445</b>  |
| <b>4</b>  | <b>Additional ML Libraries</b>   | <b>797</b>  |
| <b>5</b>  | <b>Developer Flows</b>           | <b>839</b>  |
| <b>6</b>  | <b>Runtime &amp; Tools</b>       | <b>911</b>  |
| <b>7</b>  | <b>Compiler</b>                  | <b>1123</b> |
| <b>8</b>  | <b>Learning Neuron</b>           | <b>1537</b> |
| <b>9</b>  | <b>Legacy Software</b>           | <b>1637</b> |
| <b>10</b> | <b>About Neuron</b>              | <b>1671</b> |
|           | <b>Python Module Index</b>       | <b>1681</b> |
|           | <b>Index</b>                     | <b>1683</b> |





AWS Neuron is the software development kit (SDK) used to run deep learning and generative AI workloads on [AWS Inferentia](#) and [AWS Trainium](#) powered Amazon EC2 instances and UltraServers ([Inf1](#), [Inf2](#), [Trn1](#), [Trn2](#) and [Trn2 UltraServer](#)). It includes a compiler, runtime, training and inference libraries, and profiling tools. Neuron supports customers in their end-to-end ML development lifecycle including building and deploying deep learning and AI models.

- **ML Frameworks and Libraries** - Neuron integrates with [PyTorch](#) and [JAX](#), and offers [NxD Training](#) and [NxN Inference](#) PyTorch libraries for distributed workflows. It also supports [third party libraries](#) such as Hugging Face Optimum Neuron, PyTorch Lightning, and AXLearn library for JAX model training.
- **Frontier Models Support** - Neuron supports frontier models such as Llama3.3-70b and Llama Llama3.1-405b.
- **Developer Tools** - Neuron provides [health monitoring](#), [observability](#), and [profiling tools](#) for AWS Inferentia and Trainium instances. It tracks hardware utilization, model execution metrics, and device information. The Neuron Profiler identifies performance bottlenecks. Neuron also integrates with [third-party](#) monitoring tools like Datadog and Weights and Biases.
- **Compute Kernels** - [Neuron Kernel Interface \(NKI\)](#) provides direct hardware access on AWS Trainium and Inferentia, enabling customer to write optimized kernel. NKI provides a Python-based environment with Triton-like syntax. Neuron supports custom C++ operators, allowing developers to extend functionality and enhance deep learning models.
- **Workloads Orchestrations and Managed Services** - Neuron enables you to use Trainium and Inferentia-based instances with Amazon services such as SageMaker, EKS, ECS, ParallelCluster, and Batch. and [third-party solutions](#) like Ray (Anyscale) and Domino Data Lab.
- **Architecture** - To understand the architecture of AWS AI Chips, Trn/Inf instances, and NeuronCores visit [Instance and UltraServer Architecture](#), [Amazon EC2 AI Chips Architecture](#) and [AWS NeuronCore Architecture](#).

For more information about the latest AWS Neuron release, see [Neuron 2.24.1 \(06/30/2025\)](#) and check the [Announcements](#) page.

For list of AWS Neuron model samples and tutorials on Amazon EC2 [Inf1](#), [Inf2](#), [Trn1](#), and [Trn2](#) instances, see [Model samples and tutorials](#).

[Get Started with Neuron](#)

[Neuron Quick Links](#)



## OVERVIEW

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 1.1 Neuron Quick Links

#### Overview

- *Get Started with Neuron*
- *Ask Q Developer*
- *Model samples and tutorials*
- *Neuron performance*
- *What's New*
- *Announcements*

#### ML frameworks

- *PyTorch Neuron*
- *JAX Neuron (beta)*
- *TensorFlow Neuron*
- *MXNet Neuron (maintenance)*

#### ML libraries

- *NxD Training*
- *NxD Inference*
- *NxD Core*
- *Transformers NeuronX (transformers-neuronx)*
- *AWS Neuron Reference for NeMo Megatron*

#### User Guides

- *NeuronX Runtime*
- *Neuron Compiler*
- *Neuron Kernel Interface (NKI) (beta)*
- *Neuron Custom C++ Operators (beta)*
- *Monitoring Tools*

- [Profiling Tools](#)
- [Other Tools](#)
- [Setup Guide](#)
- [Neuron DLAMI User Guide](#)
- [Neuron Containers](#)
- [AWS Workload Orchestration](#)

### Learn AWS Neuron

- [Neuron Architecture](#)
- [Neuron Features](#)
- [Neuron Application Notes](#)
- [Neuron FAQ](#)
- [Troubleshooting Guide](#)

### About AWS Neuron

- [Release Details](#)
- [Roadmap](#)
- [Support](#)

This document is relevant for: Inf1, Inf2, Trn1, Trn2

This document is relevant for: Inf1, Inf2, Trn1, Trn2



## 1.2 Ask Q Developer

Use Q Developer as your Neuron Expert for general Neuron technical guidance and to jumpstart your NKI kernel development.

Ask Q through Chat

Ask Q in your IDE

Guidelines for Quality Results

### 1.2.1 Guidelines for Quality Results

1. Be Specific: Clearly state the task, desired output, and any constraints.
2. Provide Context: Mention specific versions, strategies, and any relevant performance requirements.
3. Request Complete Code: Ask for full implementations including imports, decorators, and main functions. Remember to always review and test the generated code before using it in production.
4. Ask for Explanations: Request comments or separate explanations for complex parts of the code.
5. Iterate: If the initial response isn't satisfactory, refine your prompt based on the output. If you encounter issues or inaccuracies, consider rephrasing your prompt or breaking down complex tasks into smaller, more specific questions.

6. Fact check: Use Q as a starting point and supplement its output with official documentation, AWS NKI Samples repository, and your own expertise.

## Example Prompts

**Note:** Amazon Q Developer support for Neuron is currently in Beta. Therefore, Q may not always produce optimal or fully accurate results.

1. “Explain the key features and benefits of AWS Neuron Kernel Interface (NKI).”
2. “How do different parallelism strategies (data, pipeline, tensor) affect training performance on Neuron?”
3. “What are the best practices for optimizing matrix multiplication operations using Neuron Kernel Interface (NKI)?”
4. “Provide complete Neuron Kernel Interface (NKI) code for a matrix multiplication kernel, including imports, decorators, and explanations of key optimizations. Focus on efficient tiling and data movement strategies.”

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 1.3 Get Started with Neuron

This section walks you through the various options to get started with Neuron. You have to install Neuron on Trainium and Inferentia powered instances to enable deep-learning acceleration.

Get started with PyTorch      Get Started with JAX      Get Started with TensorFlow      Get Started with Q Developer      *This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 1.4 Model samples and tutorials

*This document is relevant for:* Trn1

### 1.4.1 Training Samples/Tutorials (Trn1/Trn1n)

#### Table of contents

- *Encoders*
- *Decoders*
- *Encoder-Decoders*
- *Vision Transformers*
- *Stable Diffusion*
- *Multi Modal*

- *Convolutional Neural Networks(CNN)*

## Encoders

| Model   | Frame-works/Libraries | Samples and Tutorials  |
|---|-----------------------|--|
| bert-base-cased                                 | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “bert-base-cased” PyTorch model for Text Classification</li> <li>• How to fine-tune a “bert base cased” PyTorch model with AWS Trainium (Trn1 instances) for Sentiment Analysis</li> </ul>  |
| bert-base-uncased                               | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “bert-base-uncased” PyTorch model</li> <li>• Fine tuning BERT base model from HuggingFace on Amazon SageMaker</li> </ul>  |
| bert-large-cased                                | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “bert-large-cased” PyTorch model</li> </ul>   |
| bert-large-uncased                              | torch-neuronx         | <ul style="list-style-type: none"> <li>• <i>Hugging Face BERT Pretraining Tutorial (Data-Parallel)</i></li> <li>• Launch Bert Large Phase 1 pretraining job on Parallel Cluster</li> <li>• Launch a Multi-Node PyTorch Neuron Training Job on Trainium Using TorchX and EKS</li> <li>• <i>PyTorch Neuron for Trainium Hugging Face BERT MRPC task finetuning using Hugging Face Trainer API</i></li> <li>• Fine-tune a “bert-large-uncased” PyTorch model</li> </ul> |
| roberta-base                                    | tensorflow-neuronx    | <ul style="list-style-type: none"> <li>• Fine-tune a “roberta-base” PyTorch model</li> </ul>   |
| roberta-large                                   | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “roberta-large” PyTorch model</li> </ul>  |
| xlm-roberta-base                                | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “xlm-roberta-base” PyTorch model</li> </ul>   |
| alberta-base-v2                                 | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “alberta-base-v2” PyTorch model</li> </ul>  |
| distilbert-base-uncased                         | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “distilbert-base-uncased” PyTorch model</li> </ul>  |
| camembert-base                                  | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tune a “camembert-base PyTorch model</li> </ul>  |
| cl-tohoku/bert-base-japanese-whole-word-masking | torch-neuronx         | <ul style="list-style-type: none"> <li>• Fine-tuning &amp; Deployment Hugging Face BERT Japanese model</li> </ul>  |



## Decoders

| Model                    | Frame-works/Libraries | Samples and Tutorials  |
|--------------------------|-----------------------|--|
| gpt-2                    | nxd-training          | <ul style="list-style-type: none"> <li>• <a href="#">Megatron GPT Pretraining</a></li> </ul>   |
| gpt-2                    | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">How to run training jobs for “gpt2” PyTorch model with AWS Trainium</a></li> <li>• <a href="#">ZeRO-1 Tutorial</a></li> </ul>   |
| gpt-3                    | neuronx-nemo-megatron | <ul style="list-style-type: none"> <li>• <a href="#">Launch a GPT-3 23B pretraining job using neuronx-nemo-megatron</a></li> <li>• <a href="#">Launch a GPT-3 46B pretraining job using neuronx-nemo-megatron</a></li> <li>• <a href="#">Launch a GPT-3 175B pretraining job using neuronx-nemo-megatron</a></li> </ul>              |
| GPT-NEOX-20B             | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training GPT-NeoX 20B with Tensor Parallelism and ZeRO-1 Optimizer</a></li> <li>• <a href="#">Training GPT-NEOX 20B model using neuronx-distributed</a></li> <li>• <a href="#">Pre-train GPT Neox 20b on Wikicorpus dataset using Neuronx Distributed library</a></li> </ul>    |
| GPT-NEOX-6.9B            | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training GPT-NeoX 6.9B with Tensor Parallelism and ZeRO-1 Optimizer</a></li> <li>• <a href="#">Training GPT-NEOX 6.9B model using neuronx-distributed</a></li> <li>• <a href="#">Pre-train GPT Neox 6.9b on Wikicorpus dataset using Neuronx Distributed library</a></li> </ul> |
| meta-llama/Llama-3.1-70b | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training Llama-3.1-70B, Llama-3-70B or Llama-2-13B/70B with Tensor Parallelism and Pipeline Parallelism</a></li> </ul>  |
| meta-llama/Llama-3.1-8b  | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training Llama3.1-8B, Llama3-8B and Llama2-7B with Tensor Parallelism and ZeRO-1 Optimizer</a></li> </ul>   |
| meta-llama/Llama-3-70b   | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training Llama-3.1-70B, Llama-3-70B or Llama-2-13B/70B with Tensor Parallelism and Pipeline Parallelism</a></li> </ul>  |
| meta-llama/Llama-3-8b    | nxd-training          | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Llama3.1/Llama3-8B Pretraining</a></li> <li>• <a href="#">HuggingFace Llama3.1/Llama3-8B Supervised Fine-tuning</a></li> </ul>  |
| meta-llama/Llama-3-8b    | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training Llama3 8B Model with Tensor Parallelism and ZeRO-1 Optimizer</a></li> <li>• <a href="#">Tutorial for Fine-tuning Llama3 8B with tensor parallelism and LoRA using Neuron PyTorch Lightning with NeuronX Distributed</a></li> </ul>                                     |
| 8                        |                       |  |
| meta-llama/Llama-2-7b    | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">Training Llama3.1-8B, Llama3-8B and Llama2-7B with</a></li> </ul>   |



## Encoder-Decoders

| Model               | Frame-works/Libraries   | Samples and Tutorials  |
|---------------------|---|--|
| t5-small            | <ul style="list-style-type: none"> <li>torch-neuronx</li> <li>optimum-neuron</li> </ul> | <ul style="list-style-type: none"> <li><i>Fine-tune T5 model on Trn1</i></li> </ul>  |
| facebook/bart-large | <ul style="list-style-type: none"> <li>torch-neuronx</li> </ul>                         | <ul style="list-style-type: none"> <li>How to fine-tune a “Bart-Large” PyTorch model with AWS Trainium (trn1 instances)</li> </ul> |

## Vision Transformers

| Model                             | Frame-works/Libraries | Samples and Tutorials   |
|-----------------------------------|-----------------------|---|
| google/vit-base-patch16-224-in21k | torch-neuronx         | <ul style="list-style-type: none"> <li>Fine-tune a pretrained HuggingFace vision transformer PyTorch model</li> </ul>           |
| openai/clip-vit-base-patch32      | torch-neuronx         | <ul style="list-style-type: none"> <li>Fine-tune a pretrained HuggingFace CLIP-base PyTorch model with AWS Trainium</li> </ul>  |
| openai/clip-vit-large-patch14     | torch-neuronx         | <ul style="list-style-type: none"> <li>Fine-tune a pretrained HuggingFace CLIP-large PyTorch model with AWS Trainium</li> </ul> |

## Stable Diffusion

| Model                                 | Frame-works/Libraries | Samples and Tutorials   |
|---------------------------------------|-----------------------|---|
| stabilityai/stable-diffusion-2-1-base | torch-neuronx         | <ul style="list-style-type: none"> <li>[Beta] Train stabilityai/stable-diffusion-2-1-base with AWS Trainium (trn1 instances)</li> </ul> |
| runwayml/stable-diffusion-v1-5        | torch-neuronx         | <ul style="list-style-type: none"> <li>[Beta] Train runwayml/stable-diffusion-v1-5 with AWS Trainium (trn1 instances)</li> </ul>        |

**Multi Modal**

| Model                 | Frame-works/Libraries | Samples and Tutorials  |
|-----------------------|-----------------------|--|
| language-perceiver    | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">How to fine-tune a “language perceiver” PyTorch model with AWS Trainium (trn1 instances)</a></li> </ul> |
| vision-perceiver-conv | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">How to fine-tune a pretrained HuggingFace Vision Perceiver Conv</a></li> </ul>                          |

**Convolutional Neural Networks(CNN)**

| Model                 | Frame-works/Libraries | Samples and Tutorials  |
|-----------------------|-----------------------|--|
| resnet50              | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">How to fine-tune a pretrained ResNet50 Pytorch model with AWS Trainium (trn1 instances) using NeuronSDK</a></li> </ul>                  |
| milesial/Pytorch-UNet | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">This notebook shows how to fine-tune a pretrained UNET PyTorch model with AWS Trainium (trn1 instances) using NeuronSDK.</a></li> </ul> |

*This document is relevant for: Trn1*

*This document is relevant for: Inf2, Trn1*

**1.4.2 Inference Samples/Tutorials (Inf2/Trn1/Trn2)****Table of contents**

- *Encoders*
- *Decoders*
- *Encoder-Decoders*
- *Vision Transformers*
- *Convolutional Neural Networks(CNN)*
- *Stable Diffusion*
- *Diffusion Transformers*
- *Audio*
- *Multi Modal*

## Encoders

| Model                          | Frame-works/Libraries | Samples and Tutorials  |
|--------------------------------|-----------------------|--|
| bert-base-cased-finetuned-mrpc | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">BERT TorchServe tutorial</a></li> <li>• <a href="#">HuggingFace pretrained BERT tutorial [html]</a> [notebook]</li> <li>• <a href="#">LibTorch C++ Tutorial for HuggingFace Pretrained BERT</a></li> <li>• <a href="#">Compiling and Deploying HuggingFace Pretrained BERT on Inf2 on Amazon SageMaker</a></li> </ul> |
| bert-base-cased-finetuned-mrpc | neuronx-distributed   | <ul style="list-style-type: none"> <li>• <a href="#">tp_inference_tutorial</a></li> </ul>  |
| bert-base-uncased              | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Pretrained BERT Inference on Trn1</a></li> </ul>  |
| distilbert-base-uncased        | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Pretrained DistilBERT Inference on Trn1</a></li> </ul>  |
| roberta-base                   | tensorflow-neuronx    | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Roberta-Base [html]</a> [notebook]</li> </ul>   |
| roberta-large                  | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Pretrained RoBERTa Inference on Trn1</a></li> </ul>   |

## Decoders

| Model                                    | Frame-works/Libraries         | Samples and Tutorials  |
|--|-------------------------------|--|
| gpt2                                     | torch-neuronx                 | <ul style="list-style-type: none"> <li>HuggingFace Pretrained GPT2 Feature Extraction on Trn1</li> </ul>   |
| meta-llama/Llama-3.3-70B                 | neuronx-distributed-inference | <ul style="list-style-type: none"> <li><i>Tutorial: Using Speculative Decoding to improve Llama-3.3-70B inference performance on Trn2 instances</i></li> <li><i>Tutorial: Scaling LLM Inference with Data Parallelism on Trn2</i></li> </ul>   |
| meta-llama/Llama-3.2-11B-Vision-Instruct | neuronx-distributed-inference | <ul style="list-style-type: none"> <li><i>Tutorial for deploying Llama3.2 Multimodal Models on Trn1 &amp; Inf2 instances</i></li> </ul>  |
| meta-llama/Llama-3.2-90B-Vision-Instruct | neuronx-distributed-inference | <ul style="list-style-type: none"> <li><i>Tutorial for deploying Llama3.2 Multimodal Models on Trn1 &amp; Inf2 instances</i></li> </ul>  |
| meta-llama/Llama-3.1-8b                  | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face Llama 3.1 8B autoregressive sampling on Inf2 &amp; Trn1 with 32k sequence length</li> <li>Run Hugging Face Llama 3.1 8B autoregressive sampling on Inf2 &amp; Trn1 with 128k sequence length</li> <li>Run meta-llama/Meta-Llama-3.1-8B autoregressive sampling on Inf2 &amp; Trn1</li> </ul> |
| meta-llama/Llama-3.1-70b                 | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face Llama 3.1 70B autoregressive sampling on Trn1 with 64k sequence length</li> <li>Run Hugging Face meta-llama/Meta-Llama-3.1-70B autoregressive sampling on Inf2 &amp; Trn1</li> </ul>   |
| meta-llama/Llama-3.1-70b-Instruct        | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face Llama-3.1-70B-Instruct + Llama-3.2-1B-Instruct Speculative Decoding on Trn1 with transformers-neuronx and vLLM</li> <li>Run Hugging Face Llama-3.1-70B-Instruct EAGLE Speculative Decoding on Trn1 with transformers-neuronx and vLLM</li> </ul>   |
| meta-llama/Llama-3.1-405b                | neuronx-distributed-inference | <ul style="list-style-type: none"> <li><i>Tutorial for deploying Llama-3.1-405B on Trn2</i></li> <li><i>Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances</i></li> </ul>   |
| meta-llama/Llama-3.1-405b                | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face Llama 3.1 405B autoregressive sampling on Trn1/Trn1n with 16k sequence length</li> </ul>   |
| meta-llama/Llama-3-8b                    | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face meta-llama/Llama-3-8b autoregressive sampling on Inf2 &amp; Trn1</li> </ul>  |
| meta-llama/Llama-3-70b                   | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face meta-llama/Llama-3-70b autoregressive sampling on Inf2 &amp; Trn1</li> </ul>   |
| meta-llama/Llama-2-13b                   | transformers-neuronx          | <ul style="list-style-type: none"> <li>Run Hugging Face meta-llama/Llama-2-13b autoregressive sampling on Inf2 &amp; Trn1</li> </ul>   |

## Encoder-Decoders

| Model             | Frame-works/Libraries   | Samples and Tutorials   |
|-------------------|---|---|
| t5-large          | <ul style="list-style-type: none"> <li>torch-neuronx</li> <li>optimum-neuron</li> </ul> | <ul style="list-style-type: none"> <li>T5 inference tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>         |
| t5-3b             | neuronx-distributed   | <ul style="list-style-type: none"> <li>T5 inference tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>         |
| google/flan-t5-xl | neuronx-distributed   | <ul style="list-style-type: none"> <li>flan-t5-xl inference tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul> |

## Vision Transformers

| Model                       | Frame-works/Libraries | Samples and Tutorials   |
|-----------------------------|-----------------------|---|
| google/vit-base-patch16-224 | torch-neuronx         | <ul style="list-style-type: none"> <li>HuggingFace Pretrained ViT Inference on Trn1</li> </ul>        |
| clip-vit-base-patch32       | torch-neuronx         | <ul style="list-style-type: none"> <li>HuggingFace Pretrained CLIP Base Inference on Inf2</li> </ul>  |
| clip-vit-large-patch14      | torch-neuronx         | <ul style="list-style-type: none"> <li>HuggingFace Pretrained CLIP Large Inference on Inf2</li> </ul> |

## Convolutional Neural Networks(CNN)

| Model    | Frame-works/Libraries | Samples and Tutorials   |
|----------|-----------------------|---|
| resnet50 | torch-neuronx         | <ul style="list-style-type: none"> <li>Torchvision Pretrained ResNet50 Inference on Trn1 / Inf2</li> <li>Torchvision ResNet50 tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul> |
| resnet50 | tensorflow-neuronx    | <ul style="list-style-type: none"> <li><i>Using NEURON_RT_VISIBLE_CORES with TensorFlow Serving</i></li> </ul>  |
| unet     | torch-neuronx         | <ul style="list-style-type: none"> <li>Pretrained UNet Inference on Trn1 / Inf2</li> </ul>  |
| vgg      | torch-neuronx         | <ul style="list-style-type: none"> <li>Torchvision Pretrained VGG Inference on Trn1 / Inf2</li> </ul>   |

## Stable Diffusion

| Model                         | Frame-works/Libraries | Samples and Tutorials   |
|-------------------------------|-----------------------|---|
| stable-diffusion-v1-5         | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace Stable Diffusion 1.5 (512x512) Inference on Trn1 / Inf2</a></li> </ul>   |
| stable-diffusion-2-1-base     | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace Stable Diffusion 2.1 (512x512) Inference on Trn1 / Inf2</a></li> </ul>   |
| stable-diffusion-2-1          | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace Stable Diffusion 2.1 (768x768) Inference on Trn1 / Inf2</a></li> <li><a href="#">Deploy &amp; Run Stable Diffusion on SageMaker and Inferentia2</a></li> </ul>                   |
| stable-diffusion-xl-base-1.0  | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace Stable Diffusion XL 1.0 (1024x1024) Inference on Inf2</a></li> <li><a href="#">HuggingFace Stable Diffusion XL 1.0 Base and Refiner (1024x1024) Inference on Inf2</a></li> </ul> |
| stable-diffusion-2-inpainting | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">stable-diffusion-2-inpainting model Inference on Trn1 / Inf2</a></li> </ul>  |

## Diffusion Transformers

| Model        | Frame-works/Libraries | Samples and Tutorials  |
|--------------|-----------------------|--|
| pixart-alpha | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace PixArt Alpha (256x256, 512x512 square resolution) Inference on Trn1 / Inf2</a></li> </ul> |
| pixart-sigma | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">HuggingFace PixArt Sigma (256x256, 512x512 square resolution) Inference on Trn1 / Inf2</a></li> </ul> |

## Audio

| Model              | Frame-works/Libraries | Samples and Tutorials   |
|--------------------|-----------------------|---|
| wav2vec2-conformer | torch-neuronx         | <ul style="list-style-type: none"> <li><a href="#">Run HuggingFace Pretrained Wav2Vec2-Conformer with Rotary Position Embeddings Inference on Inf2</a></li> <li><a href="#">Run HuggingFace Pretrained Wav2Vec2-Conformer with Relative Position Embeddings Inference on Inf2 &amp; Trn1</a></li> </ul> |

## Multi Modal

| Model                 | Frame-works/Libraries | Samples and Tutorials   |
|-----------------------|-----------------------|---|
| multimodal-perceiver  | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HuggingFace Multimodal Perceiver Inference on Trn1 / Inf2</a></li> </ul>             |
| language-perceiver    | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HF Pretrained Perceiver Language Inference on Trn1 / Inf2</a></li> </ul>             |
| vision-perceiver-conv | torch-neuronx         | <ul style="list-style-type: none"> <li>• <a href="#">HF Pretrained Perceiver Image Classification Inference on Trn1 / Inf2</a></li> </ul> |

*This document is relevant for: Inf2, Trn1*

*This document is relevant for: Inf1*

### 1.4.3 Inference Samples/Tutorials (Inf1)

#### Table of contents

- *Encoders*
- *Vision Transformers*
- *Convolutional Neural Networks(CNN)*
- *Vision*

## Encoders

| Model   | Frame-works/Libraries | Samples and Tutorials  |
|---|-----------------------|--|
| bert-base-cased-finetuned-mrpc                  | torch-neuron          | <ul style="list-style-type: none"> <li>• HuggingFace pretrained BERT tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> <li>• BertBaseCased Inference on Inf1 instances</li> <li>• Bert TorchServe tutorial <a href="#">[html]</a></li> <li>• Bring your own HuggingFace pretrained BERT container to Sagemaker Tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul> |
| bert-base-uncased                               | torch-neuron          | <ul style="list-style-type: none"> <li>• NeuronCore Pipeline tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>   |
| bert-large-uncased                              | torch-neuron          | <ul style="list-style-type: none"> <li>• BertLargeUncased Inference on Inf1 instances</li> </ul>   |
| roberta-base                                    | torch-neuron          | <ul style="list-style-type: none"> <li>• Roberta-Base inference on Inf1 instances</li> </ul>   |
| distilbert-base-uncased-finetuned-sst-2-english | tensorflow-neuron     | <ul style="list-style-type: none"> <li>• Tensorflow 2.x - HuggingFace Pipelines distilBERT with Tensorflow2 Neuron <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>  |
| glueon bert                                     | mxnet-neuron          | <ul style="list-style-type: none"> <li>• MXNet 1.8: Using data parallel mode tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>   |

## Vision Transformers

| Model                       | Frame-works/Libraries | Samples and Tutorials  |
|-----------------------------|-----------------------|--|
| ssd                         | torch-neuron          | <ul style="list-style-type: none"> <li>• Inference of SSD model on inf1 instances</li> </ul> |
| TrOCR                       | torch-neuron          | <ul style="list-style-type: none"> <li>• TrOCR inference on Inf1 instances</li> </ul>        |
| vgg                         | torch-neuron          | <ul style="list-style-type: none"> <li>• VGG inference on Inf1 instances</li> </ul>          |
| google/vit-base-patch16-224 | torch-neuron          | <ul style="list-style-type: none"> <li>• ViT model inference on Inf1</li> </ul>              |



## Convolutional Neural Networks(CNN)

| Model             | Frame-works/Libraries | Samples and Tutorials   |
|-------------------|-----------------------|---|
| EfficientNet      | torch-neuron          | <ul style="list-style-type: none"> <li>EfficientNet model inference on Inf1 instances</li> </ul>  |
| GFL (MMDetection) | torch-neuron          | <ul style="list-style-type: none"> <li>GFL (MMDetection) inference on Inf1 instances</li> </ul>   |
| HRNet             | torch-neuron          | <ul style="list-style-type: none"> <li>HRNET - Pose Estimation</li> </ul>   |
| MarianMT          | torch-neuron          | <ul style="list-style-type: none"> <li>HuggingFace MarianMT tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> <li>Inference of Pre-trained MarianMT model on Inf1</li> </ul>  |
| Detectron2 R-CNN  | torch-neuron          | <ul style="list-style-type: none"> <li>R-CNN inference on Inf1</li> </ul>   |
| resnet            | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Pre-trained Resnet model (18,34,50,101,152) on Inf1</li> <li>ResNet-50 tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>  |
| resnet            | tensorflow-neuron     | <ul style="list-style-type: none"> <li>Tensorflow 2.x - Using NEURON_RT_VISIBLE_CORES with TensorFlow Serving <a href="#">[html]</a></li> </ul>   |
| resnet            | mxnet-neuron          | <ul style="list-style-type: none"> <li>ResNet-50 tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> <li>Getting started with Gluon tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> <li>NeuronCore Groups tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul> |
| Resnext           | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Resnext model on Inf1</li> </ul>  |
| Yolov4            | torch-neuron          | <ul style="list-style-type: none"> <li>PyTorch YOLOv4 tutorial <a href="#">[html]</a> <a href="#">[notebook]</a></li> </ul>   |
| Yolov5            | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Yolov5 on Inf1</li> </ul>   |
| Yolov6            | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Yolov6 on Inf1 instances</li> </ul>   |
| Yolov7            | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Yolov7 model on Inf1</li> </ul>   |
| Yolof             | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of Yolof model on Inf1</li> </ul>  |
| fairseq           | torch-neuron          | <ul style="list-style-type: none"> <li>Inference of fairseq model on Inf1</li> </ul>  |
| unet              | tensorflow-neuron     | <ul style="list-style-type: none"> <li>Unet - Tensorflow 2.x tutorial</li> </ul>  |

## Vision

| Model         | Frame-works/Libraries | Samples and Tutorials   |
|---------------|-----------------------|---|
| craft-pytorch | torch-neuron          | <ul style="list-style-type: none"> <li>• <a href="#">CRAFT model inference on Inf1</a></li> </ul> |

*This document is relevant for: Inf1*

This section gives you the consolidated list of code samples and tutorials published by AWS Neuron across documentation and various GitHub repositories.

Training on Trn1      Inference on Inf2, Trn1 and Trn2      Inference on Inf1      For links to individual GitHub sample repositories, see [neuron-github-samples](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1*

## 1.5 Neuron performance

The Neuron performance pages provide a reference to the expected performance of various open-source models for popular deep learning in Natural Language Processing (NLP), Computer Vision (CV) and Recommender model tasks. We have included with each model links to allow you to setup and reconstruct the test with a few steps.

Inference performance numbers for Inf1      Inference performance numbers for Inf2      Inference performance numbers for Trn1      Training performance numbers for Trn1      *This document is relevant for: Inf1, Inf2, Trn1*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 1.6 What's New

### 1.6.1 Neuron 2.24.1 (06/30/2025)

Neuron version 2.24.1 resolves an installation issue that could prevent NeuronX Distributed Training from being installed successfully.

### 1.6.2 Neuron 2.24.0 (06/24/2025)

Neuron version 2.24 introduces new inference capabilities including prefix caching, disaggregated inference (Beta), and context parallelization support (Beta). This release also includes NKI language enhancements and enhanced profiling visualizations for improved debugging and performance analysis. Neuron 2.24 adds support for PyTorch 2.7 and JAX 0.6, updates existing DLAMIs and DLCs, and introduces a new vLLM inference container.

#### Table of contents

- [Inference](#)
- [Training](#)

- *Neuron Kernel Interface (NKI)*
- *Neuron Tools*
- *Neuron Deep Learning Containers (DLCs)*
- *Neuron Deep Learning AMIs (DLAMIs)*

## Inference

NxD Inference (NxDI) includes the following enhancements:

- **Prefix caching:** Improves Time To First Token (TTFT) by up to 3x when processing common shared prompts across requests.
- **Disaggregated inference (Beta):** Uses 1P1D (1 Prefill, 1 Decode) architecture to reduce prefill-decode interference and improve goodput.
- **Context parallelism (Beta):** Improves TTFT for longer sequence lengths by processing context encoding in parallel across multiple NeuronCores.
- **Model support:** Added beta support for Qwen 2.5 text models.
- **NxD Inference Library:** Upgraded to support PyTorch 2.7 and Transformers 4.48.

Hugging Face Optimum Neuron 0.2.0 now supports PyTorch-based NxD Core backend for LLM inference, simplifying the implementation of new PyTorch model architectures. Models including Llama 3.1-8B and Llama-3.3-70B have migrated from Transformers NeuronX to the NxD backend.

## Training

### Library Upgrades

- **NxD Training (NxDT) Library:** Upgraded to support PyTorch 2.7 and Transformers 4.48.
- **JAX Training Support:** Upgraded to JAX 0.6.0.

### Neuron Kernel Interface (NKI)

- **New `nki.language.gather_flattened`:** Provides efficient parallel tensor element gathering.
- **Enhanced accuracy:** Improved valid range of `nki.language.sqrt` and `nki.isa.activation(nl.sqrt)`
- **Advanced indexing:** Improved performance for `nki.isa.nc_match_replace8`.

## Neuron Tools

### Neuron Profiler Enhancements

- **Framework stack traces:** Maps device instructions to model source code.
- **Scratchpad memory usage visualization:** Shows tensor-level memory usage over time with HLO name association.
- **On-device collectives barriers:** Identifies synchronization overhead.
- **HBM throughput visualization:** Tracks data movement involving High Bandwidth Memory (HBM) over time.

### NCCOM-TEST Improvements

- Added `--report-to-json-file` flag: Outputs results in JSON format.
- Added `--show-input-output-size` flag: Explicitly displays input and output sizes based on operations.

### Neuron Deep Learning Containers (DLCs)

- Updated containers with PyTorch 2.7 support for inference and training.
- Added new inference container with NxD Inference and vLLM with FastAPI.
- JAX DLCs now support JAX 0.6.0 training.

### Neuron Deep Learning AMIs (DLAMIs)

- Updated MultiFramework DLAMIs to include PyTorch 2.7 and JAX 0.6.0.
- Added new Single Framework DLAMIs for PyTorch 2.7 and JAX 0.6.0.



## 1.6.3 Neuron 2.24 Feature Release Notes

| What's New  | Details   | Instances                  |
|---|---|----------------------------|
| NxD Core (neuronx-distributed)                            | <ul style="list-style-type: none"> <li><a href="#">NxD Core Release Notes (neuronx-distributed)</a></li> </ul>  | Trn1 / Trn1n, Trn2         |
| NxD Inference (neuronx-distributed-inference)             | <ul style="list-style-type: none"> <li><a href="#">NxD Inference Release Notes (neuronx-distributed-inference)</a></li> </ul>   | Inf2, Trn1 / Trn1n, Trn2   |
| NxD Training (neuronx-distributed-training)               | <ul style="list-style-type: none"> <li><a href="#">NxD Training Release Notes (neuronx-distributed-training)</a></li> </ul>   | Trn1 / Trn1n, Trn2         |
| PyTorch NeuronX (torch-neuronx)                           | <ul style="list-style-type: none"> <li><a href="#">PyTorch Neuron (torch-neuronx) release notes</a></li> </ul>  | Inf2, Trn1 / Trn1n, Trn2   |
| Neuron Compiler (neuronx-cc)                              | <ul style="list-style-type: none"> <li><a href="#">Neuron Compiler (neuronx-cc) release notes</a></li> </ul>  | Inf2, Trn1 / Trn1n, Trn2   |
| Neuron Kernel Interface (NKI)                             | <ul style="list-style-type: none"> <li><a href="#">Neuron Kernel Interface (NKI) release notes</a></li> </ul>   | Inf2, Trn1/ Trn1n          |
| Neuron Tools  | <ul style="list-style-type: none"> <li><a href="#">Neuron System Tools</a></li> </ul>   | Inf1, Inf2, Trn1/ Trn1n    |
| Neuron Runtime  | <ul style="list-style-type: none"> <li><a href="#">Neuron Runtime Release Notes</a></li> </ul>  | Inf1, Inf2, Trn1/ Trn1n    |
| Transformers NeuronX (transformers-neuronx) for Inference | <ul style="list-style-type: none"> <li><a href="#">Transformers Neuron (transformers-neuronx) release notes</a></li> </ul>  | Inf2, Trn1 / Trn1n         |
| Neuron Deep Learning AMIs (DLAMIs)                        | <ul style="list-style-type: none"> <li><a href="#">Neuron DLAMI User Guide</a></li> </ul>   | Inf1, Inf2, Trn1 / Trn1n   |
| Neuron Deep Learning Containers (DLCs)                    | <ul style="list-style-type: none"> <li><a href="#">neuron-dlc-release-notes</a></li> </ul>  | Inf1, Inf2, Trn1 / Trn1n   |
| Release Announcements                                     | <ul style="list-style-type: none"> <li><a href="#">announce-no-longer-support-beta-pytorch-neuroncore-placement-apis</a></li> <li><a href="#">announce-eos-block-dimension-nki</a></li> <li><a href="#">announce-eos-pytorch25</a></li> <li><a href="#">announce-eos-tensorflow-tutorial</a></li> <li><a href="#">announce-eos-tnx</a></li> <li><a href="#">announce-eos-longer-support-xla-bf16-vars</a></li> <li><a href="#">announce-eos-block-dimension-nki</a></li> <li><a href="#">announce-no-longer-support-llama-32-meta-checkpoint</a></li> <li><a href="#">announce-no-longer-support-...</a></li> </ul> | Inf1, Inf2, Trn1/ Trn1n    |
| 22  | <ul style="list-style-type: none"> <li><a href="#">announce-no-longer-support-llama-32-meta-checkpoint</a></li> <li><a href="#">announce-no-longer-support-...</a></li> </ul>   | <b>Chapter 1. Overview</b> |

For detailed release artifacts, see [Release Artifacts](#).

### 1.6.4 Previous Releases

- `prev-rn`
- `pre-release-content`
- `prev-n1-rn`

*This document is relevant for:* `Inf1`, `Inf2`, `Trn1`, `Trn2`

*This document is relevant for:* `Inf1`, `Inf2`, `Trn1`, `Trn2`

## 1.7 Announcements

This page will be replaced by ABlog. It's here to make sure it's in the TOC.

*This document is relevant for:* `Inf1`, `Inf2`, `Trn1`, `Trn2`





## ML FRAMEWORKS

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 2.1 PyTorch Neuron

PyTorch Neuron unlocks high-performance and cost-effective deep learning acceleration on AWS Trainium-based and AWS Inferentia-based Amazon EC2 instances.

The PyTorch Neuron plugin architecture enables native PyTorch models to be accelerated on Neuron devices, so you can use your existing framework application and get started easily with minimal code changes.

For help selecting a framework type for inference, see `torch-neuron_vs_torch-neuronx`

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### 2.1.1 Pytorch Neuron Setup

PyTorch Neuron (`torch-neuronx`) Setup for Inf2, Trn1, and Trn2 Instances      PyTorch Neuron (`torch-neuron`) Setup for Inf1 Instances      *This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

#### 2.1.2 Inference with `torch-neuronx` (Inf2 & Trn1/Trn2)

*This document is relevant for: Inf2, Trn1, Trn2*

#### Tutorials for Inference (`torch-neuronx`)

#### Compiling and Deploying HuggingFace Pretrained BERT on Trn1 or Inf2

##### Introduction

In this tutorial we will compile and deploy a HuggingFace Transformers BERT model for accelerated inference on Neuron. In this tutorial, we will be deploying directly on Trn1/Inf2 instances. If you are looking to deploy this model through SageMaker on Inf2 instance, please visit the [Sagemaker samples repository](#).

This tutorial will use the `bert-base-cased-finetuned-mrpc` model. This model has 12 layers, 768 hidden dimensions, 12 attention heads, and 110M total parameters. The final layer is a binary classification head that has been trained on the Microsoft Research Paraphrase Corpus (`mrpc`). The input to the model is two sentences and the output of the model is whether or not those sentences are a paraphrase of each other.

This tutorial has the following main sections:

1. Install dependencies
2. Compile the BERT model
3. Run inference on Neuron and compare results to CPU
4. Benchmark the model using multicore inference
5. Finding the optimal batch size

This Jupyter notebook should be run on a Trn1 instance (trn1.2xlarge or larger.) or Inf2 instance (inf2.xlarge or larger.)

## Install dependencies

The code in this tutorial is written for Jupyter Notebooks. To use Jupyter Notebook on the Neuron instance, you can use this [guide](#).

This tutorial requires the following pip packages:

- torch-neuronx
- neuronx-cc
- transformers

Most of these packages will be installed when configuring your environment using the Trn1/Inf2 setup guide. The additional dependencies must be installed here:

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to_
↪detect
!pip install --upgrade transformers
```

## Compile the model into an AWS Neuron optimized TorchScript

In the following section, we load the BERT model and tokenizer, get a sample input, run inference on CPU, compile the model for Neuron using `torch_neuronx.trace()`, and save the optimized model as TorchScript.

`torch_neuronx.trace()` expects a tensor or tuple of tensor inputs to use for tracing, so we unpack the tokenizer output using the `encode` function.

The result of the trace stage will be a static executable where the operations to be run upon inference are determined during compilation. This means that when inferring, the resulting Neuron model must be executed with tensors that are the exact same shape as those provided at compilation time. If a model is given a tensor at inference time whose shape does not match the tensor given at compilation time, an error will occur.

For language models, the shape of the tokenizer tensors can vary based on the length of input sentence. We can satisfy the Neuron restriction of using a fixed shape input by padding all varying input tensors to a specified length. In a deployment scenario, the padding size should be chosen based on the maximum token length that is expected to occur for the application.

In the following section we will assume that we will receive a maximum of 128 tokens at inference time. We will pad our example inputs by using `padding='max_length'` and to avoid potential errors caused by creating a tensor that is larger than `max_length=128`, we will always tokenize using `truncation=True`.

```
[ ]: import torch
import torch_neuronx
from transformers import AutoTokenizer, AutoModelForSequenceClassification
import transformers

def encode(tokenizer, *inputs, max_length=128, batch_size=1):
    tokens = tokenizer.encode_plus(
        *inputs,
        max_length=max_length,
        padding='max_length',
        truncation=True,
        return_tensors="pt"
    )
    return (
        torch.repeat_interleave(tokens['input_ids'], batch_size, 0),
        torch.repeat_interleave(tokens['attention_mask'], batch_size, 0),
        torch.repeat_interleave(tokens['token_type_ids'], batch_size, 0),
    )

# Create the tokenizer and model
name = "bert-base-cased-finetuned-mrpc"
tokenizer = AutoTokenizer.from_pretrained(name)
model = AutoModelForSequenceClassification.from_pretrained(name, torchscript=True)

# Set up some example inputs
sequence_0 = "The company HuggingFace is based in New York City"
sequence_1 = "Apples are especially bad for your health"
sequence_2 = "HuggingFace's headquarters are situated in Manhattan"

paraphrase = encode(tokenizer, sequence_0, sequence_2)
not_paraphrase = encode(tokenizer, sequence_0, sequence_1)

# Run the original PyTorch BERT model on CPU
cpu_paraphrase_logits = model(*paraphrase)[0]
cpu_not_paraphrase_logits = model(*not_paraphrase)[0]

# Compile the model for Neuron
model_neuron = torch_neuronx.trace(model, paraphrase)

# Save the TorchScript for inference deployment
filename = 'model.pt'
torch.jit.save(model_neuron, filename)
```

## Run inference and compare results

In this section we load the compiled model, run inference on Neuron, and compare the CPU and Neuron outputs.

NOTE: Although this tutorial section uses one NeuronCore (and the next section uses two NeuronCores), by default each Jupyter notebook Python process will attempt to take ownership of all NeuronCores visible on the instance. For multi-process applications where each process should only use a subset of the NeuronCores on the instance you can use `NEURON_RT_NUM_CORES=N` or `NEURON_RT_VISIBLE_CORES=< list of NeuronCore IDs >` when starting the Jupyter notebook as described in [NeuronCore Allocation and Model Placement for Inference](#).

```
[ ]: # Load the TorchScript compiled model
model_neuron = torch.jit.load(filename)

# Verify the TorchScript works on both example inputs
neuron_paraphrase_logits = model_neuron(*paraphrase)[0]
neuron_not_paraphrase_logits = model_neuron(*not_paraphrase)[0]

# Compare the results
print('CPU paraphrase logits:      ', cpu_paraphrase_logits.detach().numpy())
print('Neuron paraphrase logits:   ', neuron_paraphrase_logits.detach().numpy())
print('CPU not-paraphrase logits:   ', cpu_not_paraphrase_logits.detach().numpy())
print('Neuron not-paraphrase logits: ', neuron_not_paraphrase_logits.detach().numpy())
```

## Benchmarking

In this section we benchmark the performance of the BERT model on Neuron. By default, models compiled with `torch_neuronx` will always execute on a *single* NeuronCore. When loading *multiple* models, the default behavior of the Neuron runtime is to evenly distribute models across all available NeuronCores. The runtime places models on the NeuronCore that has the fewest models loaded to it first. In the following section, we will `torch.jit.load` multiple instances of the model which should each be loaded onto their own NeuronCore. It is not useful to load more copies of a model than the number of NeuronCores on the instance since an individual NeuronCore can only execute one model at a time.

To ensure that we are maximizing hardware utilization, we must run inferences using multiple threads in parallel. It is nearly always recommended to use some form of threading/multiprocessing and some form of model replication since even the smallest Neuron EC2 instance has 2 NeuronCores available. Applications with no form of threading are only capable of  $1 / \text{num\_neuron\_cores}$  hardware utilization which becomes especially problematic on large instances.

One way to view the hardware utilization is by executing the `neuron-top` application in the terminal while the benchmark is executing. If the monitor shows >90% utilization on all NeuronCores, this is a good indication that the hardware is being utilized effectively.

In this example we load two models, which utilizes all NeuronCores (2) on a `trn1.2xlarge` or `inf2.xlarge` instance. Additional models can be loaded and run in parallel on larger Trn1 or Inf2 instance sizes to increase throughput.

We define a benchmarking function that loads two optimized BERT models onto two separate NeuronCores, runs multithreaded inference, and calculates the corresponding latency and throughput.

```
[ ]: import time
import concurrent.futures
import numpy as np

def benchmark(filename, example, n_models=2, n_threads=2, batches_per_thread=1000):
```

(continues on next page)

(continued from previous page)

```

"""
Record performance statistics for a serialized model and its input example.

Arguments:
    filename: The serialized torchscript model to load for benchmarking.
    example: An example model input.
    n_models: The number of models to load.
    n_threads: The number of simultaneous threads to execute inferences on.
    batches_per_thread: The number of example batches to run per thread.

Returns:
    A dictionary of performance statistics.
"""

# Load models
models = [torch.jit.load(filename) for _ in range(n_models)]

# Warmup
for _ in range(8):
    for model in models:
        model(*example)

latencies = []

# Thread task
def task(model):
    for _ in range(batches_per_thread):
        start = time.time()
        model(*example)
        finish = time.time()
        latencies.append((finish - start) * 1000)

# Submit tasks
begin = time.time()
with concurrent.futures.ThreadPoolExecutor(max_workers=n_threads) as pool:
    for i in range(n_threads):
        pool.submit(task, models[i % len(models)])
end = time.time()

# Compute metrics
boundaries = [50, 95, 99]
percentiles = {}

for boundary in boundaries:
    name = f'latency_p{boundary}'
    percentiles[name] = np.percentile(latencies, boundary)
duration = end - begin
batch_size = 0
for tensor in example:
    if batch_size == 0:
        batch_size = tensor.shape[0]
inferences = len(latencies) * batch_size

```

(continues on next page)

```

throughput = inferences / duration

# Metrics
metrics = {
    'filename': str(filename),
    'batch_size': batch_size,
    'batches': len(latencies),
    'inferences': inferences,
    'threads': n_threads,
    'models': n_models,
    'duration': duration,
    'throughput': throughput,
    **percentiles,
}

display(metrics)

def display(metrics):
    """
    Display the metrics produced by `benchmark` function.

    Args:
        metrics: A dictionary of performance statistics.
    """
    pad = max(map(len, metrics)) + 1
    for key, value in metrics.items():

        parts = key.split('_')
        parts = list(map(str.title, parts))
        title = ' '.join(parts) + ":"

        if isinstance(value, float):
            value = f'{value:0.3f}'

        print(f'{title :<{pad}} {value}')

# Benchmark BERT on Neuron
benchmark(filename, paraphrase)

```

## Finding the optimal batch size

Batch size has a direct impact on model performance. The NeuronCore architecture is optimized to maximize throughput with relatively small batch sizes. This means that a Neuron compiled model can outperform a GPU model, even if running single digit batch sizes.

As a general best practice, we recommend optimizing your model's throughput by compiling the model with a small batch size and gradually increasing it to find the peak throughput on Neuron. To minimize latency, using `batch_size = 1` will nearly always be optimal. This batch size configuration is typically used for on-demand inference applications. To maximize throughput, *usually* `1 < batch_size < 10` is optimal. A configuration which uses a larger batch size is generally ideal for batched on-demand inference or offline batch processing.

In the following section, we compile BERT for multiple batch size inputs. We then run inference on each batch size and benchmark the performance. Notice that latency increases consistently as the batch size increases. Throughput increases as well, up until a certain point where the input size becomes too large to be efficient.

```
[ ]: # Compile BERT for different batch sizes
for batch_size in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:
    tokenizer = AutoTokenizer.from_pretrained(name)
    model = AutoModelForSequenceClassification.from_pretrained(name, torchscript=True)
    example = encode(tokenizer, sequence_0, sequence_2, batch_size=batch_size)
    model_neuron = torch_neuronx.trace(model, example)
    filename = f'model_batch_size_{batch_size}.pt'
    torch.jit.save(model_neuron, filename)

[ ]: # Benchmark BERT for different batch sizes
for batch_size in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:
    print('-'*50)
    example = encode(tokenizer, sequence_0, sequence_2, batch_size=batch_size)
    filename = f'model_batch_size_{batch_size}.pt'
    benchmark(filename, example)
    print()
```

*This document is relevant for: Inf2, Trn1, Trn2*

## BERT TorchServe Tutorial

### Table of Contents

- [Overview](#)
- [Run the tutorial](#)
- [Setup TorchServe](#)
- [Run TorchServe](#)
- [Benchmark TorchServe](#)

## Overview

This tutorial demonstrates the use of [TorchServe](#) with Neuron, the SDK for EC2 Inf2 and Trn1 instances. By the end of this tutorial, you will understand how TorchServe can be used to serve a model backed by EC2 Inf2/Trn1 instances. We will use a pretrained BERT-Base model to determine if one sentence is a paraphrase of another.

## Run the tutorial

Open a terminal, log into your remote instance, and activate a Pytorch virtual environment setup (see the:ref:Install PyTorch Neuron <setup-torch-neuronx>). To complete this tutorial, you will also need a compiled BERT model. You can run `trace_bert_neuronx.py` to obtain a traced BERT model.

You should now have a compiled `bert_neuron_b6.pt` file, which is required going forward.

Open a shell on the instance you prepared earlier, create a new directory named `torchserve`. Copy your compiled model from the previous tutorial into this new directory.

```
cd torchserve
python trace_bert_neuronx.py
ls
```

```
bert_neuron_b6.pt
```

Prepare a new Python virtual environment with the necessary Neuron and TorchServe components. Use a virtual environment to keep (most of) the various tutorial components isolated from the rest of the system in a controlled way.

```
pip install transformers==4.20.1 torchserve==0.7.0 torch-model-archiver==0.7.0 captum==0.
↪ 6.0
```

Install the system requirements for TorchServe.

## Amazon Linux 2 DLAMI Base

```
sudo yum -y install jq java-11-amazon-corretto-headless
sudo alternatives --config java
sudo alternatives --config javac
```

## Ubuntu 20 DLAMI Base

```
sudo apt install openjdk-11-jdk -y
```

```
java -version
```

```
openjdk version "11.0.17" 2022-10-18
OpenJDK Runtime Environment (build 11.0.17+8-post-Ubuntu-1ubuntu218.04)
OpenJDK 64-Bit Server VM (build 11.0.17+8-post-Ubuntu-1ubuntu218.04, mixed mode, sharing)
```

```
javac -version
```



```
javac 11.0.17
```

Verify that TorchServe is now available.

```
torchserve --version
```

```
TorchServe Version is 0.7.0
```

## Setup TorchServe

During this tutorial you will need to download a few files onto your instance. The simplest way to accomplish this is to paste the download links provided above each file into a `wget` command. (We don't provide the links directly because they are subject to change.) For example, right-click and copy the download link for `config.json` shown below.

Listing 2.1: `config.json`

```
{
  "model_name": "bert-base-cased-finetuned-mrpc",
  "max_length": 128,
  "batch_size": 6
}
```

Now execute the following in your shell:

```
wget <paste link here>
ls
```

```
bert_neuron_b6.pt  config.json
```

Download the `custom handler script` that will eventually respond to inference requests.

Listing 2.2: `handler_bert_neuronx.py`

```
1 import os
2 import json
3 import sys
4 import logging
5 from abc import ABC
6
7 import torch
8 import torch_neuronx
9
10 from transformers import AutoTokenizer
11 from ts.torch_handler.base_handler import BaseHandler
12
13
14 # one core per worker
15 os.environ['NEURON_RT_NUM_CORES'] = '1'
16
17 logger = logging.getLogger(__name__)
18
19 class BertEmbeddingHandler(BaseHandler, ABC):
```

(continues on next page)

(continued from previous page)

```

20  """
21  Handler class for Bert Embedding computations.
22  """
23  def __init__(self):
24      super(BertEmbeddingHandler, self).__init__()
25      self.initialized = False
26
27  def initialize(self, ctx):
28      self.manifest = ctx.manifest
29      properties = ctx.system_properties
30      self.device = 'cpu'
31      model_dir = properties.get('model_dir')
32      serialized_file = self.manifest['model']['serializedFile']
33      model_pt_path = os.path.join(model_dir, serialized_file)
34
35      # point sys.path to our config file
36      with open('config.json') as fp:
37          config = json.load(fp)
38      self.max_length = config['max_length']
39      self.batch_size = config['batch_size']
40      self.classes = ['not paraphrase', 'paraphrase']
41
42      self.model = torch.jit.load(model_pt_path)
43      logger.debug(f'Model loaded from {model_dir}')
44      self.model.to(self.device)
45      self.model.eval()
46
47      self.tokenizer = AutoTokenizer.from_pretrained(config['model_name'])
48      self.initialized = True
49
50  def preprocess(self, input_data):
51      """
52      Tokenization pre-processing
53      """
54
55      input_ids = []
56      attention_masks = []
57      token_type_ids = []
58      for row in input_data:
59          seq_0 = row['seq_0'].decode('utf-8')
60          seq_1 = row['seq_1'].decode('utf-8')
61          logger.debug(f'Received text: "{seq_0}", "{seq_1}"')
62
63          inputs = self.tokenizer.encode_plus(
64              seq_0,
65              seq_1,
66              max_length=self.max_length,
67              padding='max_length',
68              truncation=True,
69              return_tensors='pt'
70          )
71

```

(continues on next page)

(continued from previous page)

```

72         input_ids.append(inputs['input_ids'])
73         attention_masks.append(inputs['attention_mask'])
74         token_type_ids.append(inputs['token_type_ids'])
75
76         batch = (torch.cat(input_ids, 0),
77                 torch.cat(attention_masks, 0),
78                 torch.cat(token_type_ids, 0))
79
80         return batch
81
82     def inference(self, inputs):
83         """
84         Predict the class of a text using a trained transformer model.
85         """
86
87         # sanity check dimensions
88         assert(len(inputs) == 3)
89         num_inferences = len(inputs[0])
90         assert(num_inferences <= self.batch_size)
91
92         # insert padding if we received a partial batch
93         padding = self.batch_size - num_inferences
94         if padding > 0:
95             pad = torch.nn.ConstantPad1d((0, 0, 0, padding), value=0)
96             inputs = [pad(x) for x in inputs]
97
98         outputs = self.model(*inputs)[0]
99         predictions = []
100         for i in range(num_inferences):
101             prediction = self.classes[outputs[i].argmax(dim=-1).item()]
102             predictions.append([prediction])
103             logger.debug("Model predicted: '%s'", prediction)
104         return predictions
105
106     def postprocess(self, inference_output):
107         return inference_output

```

Next, we need to associate the handler script with the compiled model using `torch-model-archiver`. Run the following commands in your terminal:

```

mkdir model_store
MAX_LENGTH=$(jq '.max_length' config.json)
BATCH_SIZE=$(jq '.batch_size' config.json)
MODEL_NAME=bert-max_length$MAX_LENGTH-batch_size$BATCH_SIZE
torch-model-archiver --model-name "$MODEL_NAME" --version 1.0 --serialized-file ./bert_
↪neuron_b6.pt --handler "./handler_bert_neuronx.py" --extra-files "./config.json" --
↪export-path model_store

```

**Note:** If you modify your model or a dependency, you will need to rerun the archiver command with the `-f` flag appended to update the archive.

The result of the above will be a `mar` file inside the `model_store` directory.

```
ls model_store
```

```
bert-max_length128-batch_size6.mar
```

This file is essentially an archive associated with a fixed version of your model along with its dependencies (e.g. the handler code).

---

**Note:** The version specified in the `torch-model-archiver` command can be appended to REST API requests to access a specific version of your model. For example, if your model was hosted locally on port 8080 and named “bert”, the latest version of your model would be available at `http://localhost:8080/predictions/bert`, while version 1.0 would be accessible at `http://localhost:8080/predictions/bert/1.0`. We will see how to perform inference using this API in Step 6.

---

Create a `custom config` file to set some parameters. This file will be used to configure the server at launch when we run `torchserve --start`.

Listing 2.3: `torchserve.config`

```
# bind inference API to all network interfaces with SSL enabled
inference_address=http://0.0.0.0:8080
default_workers_per_model=1
```

---

**Note:** This will cause TorchServe to bind on all interfaces. For security in real-world applications, you’ll probably want to use port 8443 and `enable SSL`.

---

## Run TorchServe

It’s time to start the server. Typically we’d want to launch this in a separate console, but for this demo we’ll just redirect output to a file.

```
torchserve --start --ncs --model-store model_store --ts-config torchserve.config 2>&1 >
↪ torchserve.log
```

Verify that the server seems to have started okay.

```
curl http://127.0.0.1:8080/ping
```

```
{
  "status": "Healthy"
}
```

---

**Note:** If you get an error when trying to ping the server, you may have tried before the server was fully launched. Check `torchserve.log` for details.

---

Use the Management API to instruct TorchServe to load our model.

First, determine the number of NeuronCores available based on your instance size.

## Inf2

| Instance Size | # of NeuronCores |
|---------------|------------------|
| xlarge        | 2                |
| 8xlarge       | 2                |
| 24xlarge      | 12               |
| 48xlarge      | 24               |

## Trn1

| Instance Size | # of NeuronCores |
|---------------|------------------|
| 2xlarge       | 2                |
| 32xlarge      | 32               |

```
MAX_BATCH_DELAY=5000 # ms timeout before a partial batch is processed
INITIAL_WORKERS=2 # Number from table above
curl -X POST "http://localhost:8081/models?url=$MODEL_NAME.mar&batch_size=$BATCH_SIZE&
↪initial_workers=$INITIAL_WORKERS&max_batch_delay=$MAX_BATCH_DELAY"
```

```
{
  "status": "Model \"bert-max_length128-batch_size6\" Version: 1.0 registered with X_
↪initial workers"
}
```

**Warning:** You shouldn't set INITIAL\_WORKERS above the number of NeuronCores. If you attempt to load more models than NeuronCores available, one of two things will occur. Either the extra models will fit in device memory but performance will suffer, or you will encounter an error on your initial inference. However, you may want to use fewer cores if you are using the *NeuronCore Pipeline* feature.

**Note:** Any additional attempts to configure the model after the initial curl request will cause the server to return a 409 error. You'll need to stop/start/configure the server to realize any changes.

The MAX\_BATCH\_DELAY is a timeout value that determines how long to wait before processing a partial batch. This is why the handler code needs to check the batch dimension and potentially add padding. TorchServe will instantiate the number of model handlers indicated by INITIAL\_WORKERS, so this value controls how many models we will load onto Inferentia in parallel. If you want to control worker scaling more dynamically, [see the docs](#).

It looks like everything is running successfully at this point, so it's time for an inference.

Create the infer\_bert.py file below on your instance.

Listing 2.4: infer\_bert.py

```
1 import json
2 import concurrent.futures
3 import requests
4
5 with open('config.json') as fp:
```

(continues on next page)

(continued from previous page)

```

6     config = json.load(fp)
7 max_length = config['max_length']
8 batch_size = config['batch_size']
9 name = f'bert-max_length{max_length}-batch_size{batch_size}'
10
11 # dispatch requests in parallel
12 url = f'http://localhost:8080/predictions/{name}'
13 paraphrase = {'seq_0': "HuggingFace's headquarters are situated in Manhattan",
14               'seq_1': "The company HuggingFace is based in New York City"}
15 not_paraphrase = {'seq_0': paraphrase['seq_0'], 'seq_1': 'This is total nonsense.'}
16
17 with concurrent.futures.ThreadPoolExecutor(max_workers=batch_size) as executor:
18     def worker_thread(worker_index):
19         # we'll send half the requests as not_paraphrase examples for sanity
20         data = paraphrase if worker_index < batch_size//2 else not_paraphrase
21         try:
22             response = requests.post(url, data=data)
23
24             # Check if the response status code indicates success
25             if response.status_code == 200:
26                 print(worker_index, response.json())
27             else:
28                 # If the response is not successful, raise an exception with the status_
↳ code and error message
29                 error_message = response.json().get('message', 'Unknown Error')
30                 raise Exception(f"Failed request with status code {response.status_code}:
↳ {error_message}")
31         except Exception as e:
32             # Catch all other exceptions that may be raised
33             print(f"An unexpected error occurred: {e}")
34             raise
35
36     for worker_index in range(batch_size):
37         executor.submit(worker_thread, worker_index)

```

This script will send a `batch_size` number of requests to our model. In this example, we are using a model that estimates the probability that one sentence is a paraphrase of another. The script sends positive examples in the first half of the batch and negative examples in the second half.

Execute the script in your terminal.

```
python infer_bert.py
```

```

1 ['paraphrase']
3 ['not paraphrase']
4 ['not paraphrase']
0 ['paraphrase']
5 ['not paraphrase']
2 ['paraphrase']

```

We can see that the first three threads (0, 1, 2) all report `paraphrase`, as expected. If we instead modify the script to send an incomplete batch and then wait for the timeout to expire, the excess padding results will be discarded.

## Benchmark TorchServe

We've seen how to perform a single batched inference, but how many inferences can we process per second? A separate upcoming tutorial will document performance tuning to maximize throughput. In the meantime, we can still perform a simple naïve stress test. The code below will spawn 64 worker threads, with each thread repeatedly sending a full batch of data to process. A separate thread will periodically print throughput and latency measurements.

Listing 2.5: benchmark\_bert.py

```

1  import os
2  import argparse
3  import time
4  import numpy as np
5  import requests
6  import sys
7  from concurrent import futures
8
9  import torch
10
11
12  parser = argparse.ArgumentParser()
13  parser.add_argument('--url', help='Torchserve model URL', type=str, default=f'http://127.
    ↳ 0.0.1:8080/predictions/bert-max_length128-batch_size6')
14  parser.add_argument('--num_thread', type=int, default=64, help='Number of threads,
    ↳ invoking the model URL')
15  parser.add_argument('--batch_size', type=int, default=6)
16  parser.add_argument('--sequence_length', type=int, default=128)
17  parser.add_argument('--latency_window_size', type=int, default=1000)
18  parser.add_argument('--throughput_time', type=int, default=300)
19  parser.add_argument('--throughput_interval', type=int, default=10)
20  args = parser.parse_args()
21
22  data = { 'seq_0': 'A completely made up sentence.',
23          'seq_1': 'Well, I suppose they are all made up.' }
24  live = True
25  num_infer = 0
26  latency_list = []
27
28
29  def one_thread(pred, feed_data):
30      global latency_list
31      global num_infer
32      global live
33      session = requests.Session()
34      while True:
35          start = time.time()
36          result = session.post(pred, data=feed_data)
37          latency = time.time() - start
38          latency_list.append(latency)
39          num_infer += 1
40          if not live:
41              break
42

```

(continues on next page)

(continued from previous page)

```

43
44 def current_performance():
45     last_num_infer = num_infer
46     for _ in range(args.throughput_time // args.throughput_interval):
47         current_num_infer = num_infer
48         throughput = (current_num_infer - last_num_infer) / args.throughput_interval
49         p50 = 0.0
50         p90 = 0.0
51         if latency_list:
52             p50 = np.percentile(latency_list[-args.latency_window_size:], 50)
53             p90 = np.percentile(latency_list[-args.latency_window_size:], 90)
54         print('pid {}: current throughput {}, latency p50={:.3f} p90={:.3f}'.format(os.
↪getpid(), throughput, p50, p90))
55         sys.stdout.flush()
56         last_num_infer = current_num_infer
57         time.sleep(args.throughput_interval)
58     global live
59     live = False
60
61
62 with futures.ThreadPoolExecutor(max_workers=args.num_thread+1) as executor:
63     executor.submit(current_performance)
64     for _ in range(args.num_thread):
65         executor.submit(one_thread, args.url, data)

```

Run the benchmarking script.

```
python benchmark_bert.py
```

```

pid 1214554: current throughput 0.0, latency p50=0.000 p90=0.000
pid 1214554: current throughput 713.9, latency p50=0.071 p90=0.184
pid 1214554: current throughput 737.9, latency p50=0.071 p90=0.184
pid 1214554: current throughput 731.6, latency p50=0.068 p90=0.192
pid 1214554: current throughput 732.2, latency p50=0.070 p90=0.194
pid 1214554: current throughput 733.9, latency p50=0.070 p90=0.187
pid 1214554: current throughput 739.3, latency p50=0.071 p90=0.184
...

```

**Note:** Your throughput numbers may differ from these based on instance type and size.

**Congratulations!** By now you should have successfully served a batched model over TorchServe.

You can now shutdown torchserve.

```
torchserve --stop
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1



## LibTorch C++ Tutorial

### Table of Contents

- [Overview](#)
- [Notes](#)
- [Run the tutorial](#)
- [Benchmark](#)
- [Troubleshooting](#)

### Overview

This tutorial demonstrates the use of **LibTorch** with Neuron, the SDK for Amazon Inf1, Inf2 and Trn1 instances. By the end of this tutorial, you will understand how to write a native C++ application that performs inference on EC2 Inf1, Inf2 and Trn1 instances. We will use an inf1.6xlarge and a pretrained BERT-Base model to determine if one sentence is a paraphrase of another.

Verify that this tutorial is running in a virtual environment that was set up according to the *Torch-Neuronx Installation Guide* <<https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/setup/torch-neuronx.html#setup-torch-neuronx>> or *Torch-Neuron Installation Guide* <<https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/setup/torch-neuron.html#setup-torch-neuron>>

### Notes

The tutorial has been tested on Inf1, Inf2 and Trn1 instances on ubuntu instances.

### Run the tutorial

This tutorial is self contained. It produces similar output to [\[html\]](#) [\[notebook\]](#).

Note: The tutorial will use about 8.5 GB of disk space. Please ensure you have sufficient space before beginning.

Right-click and copy this link address to the tutorial archive.

```
wget <paste archive URL>
tar xvf libtorch_demo.tar.gz
```

Your directory tree should now look like this:

```
libtorch_demo
├── bert_neuronx
│   ├── compile.py
│   └── detect_instance.py
├── clean.sh
├── core_count
│   ├── build.sh
│   └── main.cpp
└── example_app
```

(continues on next page)

(continued from previous page)

```

├── build.sh
├── core_count.hpp
├── example_app.cpp
├── README.txt
├── utils.cpp
├── utils.hpp
├── neuron.patch
├── run_tests.sh
├── setup.sh
├── tokenizer.json
├── tokenizers_binding
│   ├── build_python.sh
│   ├── build.sh
│   ├── remote_rust_tokenizer.h
│   ├── run_python.sh
│   ├── run.sh
│   ├── tokenizer.json
│   ├── tokenizer_test
│   ├── tokenizer_test.cpp
│   └── tokenizer_test.py

```

This tutorial uses the [HuggingFace Tokenizers](#) library implemented in Rust. Install Cargo, the package manager for the Rust programming language.

| Ubuntu                                 | Amazon Linux                           |
|--|--|
| <code>sudo apt install -y cargo</code> | <code>sudo yum install -y cargo</code> |

Run the setup script to download additional dependencies and build the app. (This may take a few minutes to complete.)

```

cd libtorch_demo
chmod +x setup.sh && ./setup.sh

```

```

...
+ PATH_NEURON_LIB=/opt/aws/neuron/lib/
+ g++ utils.cpp example_app.cpp -o ../example-app -O2 -D_GLIBCXX_USE_CXX11_ABI=0 -I../
↳ libtorch/include -L../tokenizers_binding/lib -L/opt/aws/neuron/lib/ -L../libtorch/lib -
↳ Wl,-rpath,libtorch/lib -Wl,-rpath,tokenizers_binding/lib -Wl,-rpath,/opt/aws/neuron/
↳ lib/ -ltokenizers -ltorchneuron -ltorch_cpu -lc10 -lpthread -lnrt
~/libtorch_demo
Successfully completed setup

```

## Benchmark

The setup script should have compiled and saved a PyTorch model compiled for neuron (bert\_neuron\_b6.pt). Run the provided sanity tests to ensure everything is working properly.

```
./run_tests.sh bert_neuron_b6.pt
```

```
Running tokenization sanity checks.
```

```
None of PyTorch, TensorFlow >= 2.0, or Flax have been found. Models won't be available,
and only tokenizers, configuration and file/data utilities can be used.
```

```
Tokenizing: 100%| 10000/10000 [00:00<00:00, 15021.69it/s]
```

```
Python took 0.67 seconds.
```

```
Sanity check passed.
```

```
Begin 10000 timed tests.
```

```
.....
```

```
End timed tests.
```

```
C++ took 0.226 seconds.
```

```
Tokenization sanity checks passed.
```

```
Running end-to-end sanity check.
```

```
The company HuggingFace is based in New York City
```

```
HuggingFace's headquarters are situated in Manhattan
```

```
not paraphrase: 10%
```

```
paraphrase: 90%
```

```
The company HuggingFace is based in New York City
```

```
Apples are especially bad for your health
```

```
not paraphrase: 94%
```

```
paraphrase: 6%
```

```
Sanity check passed.
```

Finally, run the example app directly to benchmark the BERT model.

**Note:** You can safely ignore the warning about None of PyTorch, Tensorflow >= 2.0, .... This occurs because the test runs in a small virtual environment that doesn't require the full frameworks.

```
./example-app bert_neuron_b6.pt
```

```
Getting ready.....
```

```
Benchmarking.....
```

```
Completed 32000 operations in 43 seconds => 4465.12 pairs / second
```

```
=====
```

```
Summary information:
```

```
=====
```

```
Batch size = 6
```

```
Num neuron cores = 16
```

```
Num runs per neuron core = 2000
```

**Congratulations!** By now you should have successfully built and used a native C++ application with LibTorch.

### Troubleshooting

- In the event of SIGBUS errors you may have insufficient disk space for the creation of temporary model files at runtime. Consider clearing space or mounting additional disk storage.
- In the event of a neuron runtime failure, confirm that the Neuron kernel module is loaded using `sudo modprobe neuron`.

*This document is relevant for: Inf1*

### Compiling and Deploying ResNet50 on Trn1 or Inf2

#### Introduction

In this tutorial we will compile and deploy a torchvision ResNet50 model for accelerated inference on Neuron. To get started with Jupyter Notebook on Neuron Instance you launched, please use this [guide](#).

This tutorial will use the [resnet50](#) model, which is primarily used for arbitrary image classification tasks.

This tutorial has the following main sections:

1. Install dependencies
2. Compile the ResNet model
3. Run inference on Neuron and compare results to CPU
4. Benchmark the model using multicore inference
5. Finding the optimal batch size

This Jupyter notebook should be run on a Trn1 instance (trn1.2xlarge or larger.) or Inf2 instance (inf2.xlarge or larger.)

#### Install Dependencies

The code in this tutorial is written for Jupyter Notebooks. To use Jupyter Notebook on the Neuron instance, you can use this [guide](#).

This tutorial requires the following pip packages:

- `torch-neuronx`
- `neuronx-cc`
- `torchvision`
- `Pillow`

Most of these packages will be installed when configuring your environment using the Trn1 setup guide. The additional dependencies must be installed here:

```
[ ]: !pip install Pillow
```

## Compile the model into an AWS Neuron optimized TorchScript

In the following section, we load the model, get a sample input, run inference on CPU, compile the model for Neuron using `torch_neuronx.trace()`, and save the optimized model as TorchScript.

`torch_neuronx.trace()` expects a tensor or tuple of tensor inputs to use for tracing, so we convert the input image into a tensor using the `get_image` function.

The result of the trace stage will be a static executable where the operations to be run upon inference are determined during compilation. This means that when inferring, the resulting Neuron model must be executed with tensors that are the exact same shape as those provided at compilation time. If a model is given a tensor at inference time whose shape does not match the tensor given at compilation time, an error will occur.

In the following section, we assume that we will receive an image shape of `[1, 3, 224, 224]` at inference time.

```
[ ]: import os
import urllib
from PIL import Image

import torch
import torch_neuronx
from torchvision import models
from torchvision.transforms import functional

def get_image(batch_size=1, image_shape=(224, 224)):
    # Get an example input
    filename = "0000000039769.jpg"
    if not os.path.exists(filename):
        url = "http://images.cocodataset.org/val2017/0000000039769.jpg"
        urllib.request.urlretrieve(url, filename)
    image = Image.open(filename).convert('RGB')
    image = functional.resize(image, (image_shape))
    image = functional.to_tensor(image)
    image = torch.unsqueeze(image, 0)
    image = torch.repeat_interleave(image, batch_size, 0)
    return (image, )

# Create the model
model = models.resnet50(pretrained=True)
model.eval()

# Get an example input
image = get_image()

# Run inference on CPU
output_cpu = model(*image)

# Compile the model
model_neuron = torch_neuronx.trace(model, image)

# Save the TorchScript for inference deployment
filename = 'model.pt'
torch.jit.save(model_neuron, filename)
```

## Run inference and compare results

In this section we load the compiled model, run inference on Neuron, and compare the CPU and Neuron outputs using the ImageNet classes.

```
[ ]: import json

# Load the TorchScript compiled model
model_neuron = torch.jit.load(filename)

# Run inference using the Neuron model
output_neuron = model_neuron(*image)

# Compare the results
print(f"CPU tensor:    {output_cpu[0][0:10]}")
print(f"Neuron tensor: {output_neuron[0][0:10]}")

# Download and read the ImageNet classes
url = urllib.request.urlretrieve("https://s3.amazonaws.com/deep-learning-models/image-models/imagenet_class_index.json", "imagenet_class_index.json")
with open("imagenet_class_index.json", "r") as file:
    class_id = json.load(file)
    id2label = [class_id[str(i)][1] for i in range(len(class_id))]

# Lookup and print the top-5 labels
top5_cpu = output_cpu[0].sort()[1][-5:]
top5_neuron = output_neuron[0].sort()[1][-5:]
top5_labels_cpu = [id2label[idx] for idx in top5_cpu]
top5_labels_neuron = [id2label[idx] for idx in top5_neuron]
print(f"CPU top-5 labels:    {top5_labels_cpu}")
print(f"Neuron top-5 labels: {top5_labels_neuron}")
```

## Benchmarking

In this section we benchmark the performance of the ResNet model on Neuron. By default, models compiled with `torch_neuronx` will always execute on a *single* NeuronCore. When loading *multiple* models, the default behavior of the Neuron runtime is to evenly distribute models across all available NeuronCores. The runtime places models on the NeuronCore that has the fewest models loaded to it first. In the following section, we will `torch.jit.load` multiple instances of the model which should each be loaded onto their own NeuronCore. It is not useful to load more copies of a model than the number of NeuronCores on the instance since an individual NeuronCore can only execute one model at a time.

To ensure that we are maximizing hardware utilization, we must run inferences using multiple threads in parallel. It is nearly always recommended to use some form of threading/multiprocessing and some form of model replication since even the smallest Neuron EC2 instance has 2 NeuronCores available. Applications with no form of threading are only capable of  $1 / \text{num\_neuron\_cores}$  hardware utilization which becomes especially problematic on large instances.

One way to view the hardware utilization is by executing the `neuron-top` application in the terminal while the benchmark is executing. If the monitor shows >90% utilization on all NeuronCores, this is a good indication that the hardware is being utilized effectively.

In this example we load two models, which utilizes all NeuronCores (2) on a `trn1.2xlarge` or `inf2.xlarge` instance. Additional models can be loaded and run in parallel on larger Trn1 or Inf2 instance sizes to increase throughput.

We define a benchmarking function that loads two optimized ResNet models onto two separate NeuronCores, runs multithreaded inference, and calculates the corresponding latency and throughput.

```
[ ]: import time
import concurrent.futures
import numpy as np

def benchmark(filename, example, n_models=2, n_threads=2, batches_per_thread=1000):
    """
    Record performance statistics for a serialized model and its input example.

    Arguments:
        filename: The serialized torchscript model to load for benchmarking.
        example: An example model input.
        n_models: The number of models to load.
        n_threads: The number of simultaneous threads to execute inferences on.
        batches_per_thread: The number of example batches to run per thread.

    Returns:
        A dictionary of performance statistics.
    """

    # Load models
    models = [torch.jit.load(filename) for _ in range(n_models)]

    # Warmup
    for _ in range(8):
        for model in models:
            model(*example)

    latencies = []

    # Thread task
    def task(model):
        for _ in range(batches_per_thread):
            start = time.time()
            model(*example)
            finish = time.time()
            latencies.append((finish - start) * 1000)

    # Submit tasks
    begin = time.time()
    with concurrent.futures.ThreadPoolExecutor(max_workers=n_threads) as pool:
        for i in range(n_threads):
            pool.submit(task, models[i % len(models)])
    end = time.time()

    # Compute metrics
    boundaries = [50, 95, 99]
    percentiles = {}

    for boundary in boundaries:
```

(continues on next page)

(continued from previous page)

```

        name = f'latency_p{boundary}'
        percentiles[name] = np.percentile(latencies, boundary)
    duration = end - begin
    batch_size = 0
    for tensor in example:
        if batch_size == 0:
            batch_size = tensor.shape[0]
    inferences = len(latencies) * batch_size
    throughput = inferences / duration

# Metrics
metrics = {
    'filename': str(filename),
    'batch_size': batch_size,
    'batches': len(latencies),
    'inferences': inferences,
    'threads': n_threads,
    'models': n_models,
    'duration': duration,
    'throughput': throughput,
    **percentiles,
}

display(metrics)

def display(metrics):
    """
    Display the metrics produced by `benchmark` function.

    Args:
        metrics: A dictionary of performance statistics.
    """
    pad = max(map(len, metrics)) + 1
    for key, value in metrics.items():

        parts = key.split('_')
        parts = list(map(str.title, parts))
        title = ' '.join(parts) + ":"

        if isinstance(value, float):
            value = f'{value:0.3f}'

        print(f'{title :<{pad}} {value}')

# Benchmark ResNet on Neuron
benchmark(filename, image)

```



## Finding the optimal batch size

Batch size has a direct impact on model performance. The NeuronCore architecture is optimized to maximize throughput with relatively small batch sizes. This means that a Neuron compiled model can outperform a GPU model, even if running single digit batch sizes.

As a general best practice, we recommend optimizing your model's throughput by compiling the model with a small batch size and gradually increasing it to find the peak throughput on Neuron. To minimize latency, using `batch_size = 1` will nearly always be optimal. This batch size configuration is typically used for on-demand inference applications. To maximize throughput, *usually* `1 < batch_size < 10` is optimal. A configuration which uses a larger batch size is generally ideal for batched on-demand inference or offline batch processing.

In the following section, we compile ResNet for multiple batch size inputs. We then run inference on each batch size and benchmark the performance. Notice that latency increases consistently as the batch size increases. Throughput increases as well, up until a certain point where the input size becomes too large to be efficient.

```
[ ]: # Compile ResNet for different batch sizes
for batch_size in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:
    model = models.resnet50(pretrained=True)
    model.eval()
    example = get_image(batch_size=batch_size)
    model_neuron = torch_neuronx.trace(model, example)
    filename = f'model_batch_size_{batch_size}.pt'
    torch.jit.save(model_neuron, filename)
```

```
[ ]: # Benchmark ResNet for different batch sizes
for batch_size in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]:
    print('-'*50)
    example = get_image(batch_size=batch_size)
    filename = f'model_batch_size_{batch_size}.pt'
    benchmark(filename, example)
    print()
```

## T5 model inference on Trn1 or Inf2

### Introduction

In this tutorial we will compile and deploy a pretrained T5 model for accelerated inference on Neuron.

This tutorial will use the [t5-large](#) model. The T5 model can be used for machine translation, document summarization, question answering, and classification tasks.

This tutorial has the following main sections:

1. Install dependencies
2. Compile the T5 model
3. Run inference with greedy decoding on Neuron
4. Run inference with beam search on Neuron

This Jupyter notebook should be run on a Trn1 instance (trn1.2xlarge or larger.) or Inf2 instance (inf2.xlarge or larger.)

## Install dependencies

The code in this tutorial is written for Jupyter Notebooks. To use Jupyter Notebook on the Neuron instance, you can use this [guide](#).

This tutorial requires the following pip packages:

- `torch-neuronx`
- `neuronx-cc`
- `transformers`
- `optimum-neuron`

Most of these packages will be installed when configuring your environment using the Trn1/Inf2 setup guide. The additional dependencies must be installed here:

```
[ ]: !pip install --upgrade transformers==4.31.0 optimum-neuron==0.0.8 sentencepiece
```

Optimum Neuron is the interface between the Transformers library and AWS Accelerators including AWS Trainium and AWS Inferentia. It provides a set of tools enabling easy model loading, training and inference on single- and multi-Accelerator settings for different downstream tasks. In this tutorial we use HuggingFace Optimum Neuron's `generate()` method instead of `transformers's generate()` to perform greedy decoding. Optimum Neuron takes care of padding the inputs which is necessary to infer on Neuron.

## Compile the model into an AWS Neuron optimized TorchScript

In the following section, we load the T5 model, compile the model's encoder and decoder for Neuron using `torch_neuronx.trace()`, and save the optimized encoder and decoder as TorchScript.

Before we trace the model, we need to make a couple of changes.

1. We need to write encoder and decoder wrappers - `torch_neuronx` can only trace functions with positional arguments. But the T5 encoder and decoder both use keyword arguments. So, in order to trace them, we have to write wrappers that convert keyword arguments to positional arguments
2. We modify the t5 code to maximize the computation on the neuron device - Having sections of code running on cpu will reduce the performance. Moreover, we do not want to move data between the neuron device and cpu during inference. The code we trace with `torch_neuronx` is the code that runs on the neuron device, so we refactor the t5 code to run computationally heavy operations within the wrapper.

Let us start with the EncoderWrapper.

In the huggingface t5 implementation, the encoder block takes in the input ids and returns the encoder hidden states. This hidden states are then used to initialize the KV cache in the decoder blocks during the first decoder invocation. We could trace both the encoder and the cache initialization step separately. But there is a better way, we could just compute the initial KV cache state within the encoder wrapper. This way, we remove the overhead of moving the hidden states from neuron device to cpu and back. This also allows neuron's compiler to optimize execution across both the encoder and cache initialization.

*Why don't we just initialize the cache on the first decoder run?*

This is harder to do on Neuron. Similar to `torch.jit.trace()`, `torch_neuronx.trace()` produces a function that has a fixed control flow, i.e. there are no conditional executions. So we cannot choose to conditionally initialize the cache in the first decoder iteration. Instead, we can compute the initial cache state outside the generation flow and pass the cache to it.

```

[ ]: import torch

from transformers.models.t5.modeling_t5 import T5Stack, T5LayerCrossAttention

class EncoderWrapper(torch.nn.Module):
    '''
        We will trace an instance of the EncoderWrapper.
        This wrapper just converts positional args to kwargs.
    '''

    def __init__(self,
                  encoder,
                  decoder,
                  model_config,
                  batch_size,
                  max_length,
                  device,
                  num_beams,
                  tp_degree=None):

        super().__init__()
        self.encoder = encoder
        self.decoder = decoder
        self.batch_size = batch_size
        self.max_length = max_length
        self.model_config = model_config
        self.device = device
        self.num_beams = num_beams
        self.num_attention_heads_per_partition = model_config.num_heads
        self.tp_degree = tp_degree

    def forward(self, input_ids, attention_mask):
        '''
            This is the core functionality we want to trace.
        '''
        encoder_output = self.encoder(input_ids=input_ids,
                                      attention_mask=attention_mask,
                                      output_attentions=False,
                                      output_hidden_states=False)

        last_hidden_state = encoder_output["last_hidden_state"]
        encoder_hidden_states = torch.concat([tensor.unsqueeze(0).repeat(self.num_beams,
↪1, 1) for tensor in last_hidden_state])

        decoder_blocks = self.decoder.block
        present_key_value_states_sa = []
        present_key_value_states_ca = []

        for i, block in enumerate(decoder_blocks):

            # Cross attention has to be initialized with the encoder hidden state
            cross_attention: T5LayerCrossAttention = block.layer[1]
            attention = cross_attention.EncDecAttention

```

(continues on next page)

(continued from previous page)

```

def shape(states):
    """projection"""
    return states.view(self.batch_size, -1, self.num_attention_heads_per_
↪partition, attention.key_value_proj_dim).transpose(1, 2)

    key_states = shape(attention.k(encoder_hidden_states))
    value_states = shape(attention.v(encoder_hidden_states))

    # cross_attn_kv_state
    present_key_value_states_ca.append(key_states)
    present_key_value_states_ca.append(value_states)

    # Self attention kv states are initialized to zeros. This is done to keep
↪the size of the kv cache tensor constant.
    # The kv cache will be an input to the decoder trace. Any traced function
↪will have a fixed control flow. What this means
    # is that the trace performs the exact same computations on inputs of the
↪same shape in each invocation. So the attention
    # kv cache is padded here to keep a fixed shape.
    present_key_value_states_sa.append(torch.zeros((self.batch_size,
↪
                                # key states
                                self.model_config.num_heads,
                                self.max_length-1,
                                self.model_config.d_kv),
↪dtype=torch.float32, device=self.device))
    present_key_value_states_sa.append(torch.zeros((self.batch_size,
↪
                                # value states
                                self.model_config.num_heads,
                                self.max_length-1,
                                self.model_config.d_kv),
↪dtype=torch.float32, device=self.device))

    return present_key_value_states_sa + present_key_value_states_ca

```

In the decoder wrapper, in addition to converting keyword arguments to positional arguments we add support for attention caching. Generating text from the encoder decoder models is an autoregressive process. For each invocation, we have to compute the key and value states of the attention heads repeatedly. To improve the performance, we cache the key and value states. This cache is what HuggingFace transformers code refers to as `past_key_values`.

In HuggingFace transformers, the `past_key_values` are updated outside the decoder. This works for training and evaluation but for inference we want to perform them within a single trace. This way, we can optimize across both the decoder execution and cache update. So, we move the cache update within the decoder wrapper.

```
[3]: class DecoderWrapper(torch.nn.Module):
```

```

    def __init__(self,
                  decoder: T5Stack,
                  lm_head: torch.nn.Linear,
                  model_config,
                  num_beams: int,
                  max_length: int,

```

(continues on next page)

(continued from previous page)

```

        device: str,
        tp_degree=None):
    super().__init__()
    self.decoder = decoder
    self.lm_head = lm_head
    self.model_dim=model_config.d_model
    self.device = device
    self.num_beams = num_beams
    self.batch_size = 1
    self.config = model_config

    num_heads=model_config.num_heads
    num_decoder_layers=model_config.num_decoder_layers

    self.num_attention_heads_per_partition = num_heads

    # (num_beams, n_heads, seq_length, dim_per_head)
    if device == "cpu":
        self.past_key_values_sa = [torch.ones((num_beams,num_heads,max_length-1,
↪model_config.d_kv), dtype=torch.float32) for _ in range(num_decoder_layers * 2)]
        self.past_key_values_ca = [torch.ones((num_beams,num_heads,max_length,model_
↪config.d_kv), dtype=torch.float32) for _ in range(num_decoder_layers * 2)]
    elif device == "xla":
        self.past_key_values_sa = torch.nn.ParameterList([torch.nn.Parameter(torch.
↪ones((num_beams,self.num_attention_heads_per_partition,max_length-1,model_config.d_kv),
↪ dtype=torch.float32), requires_grad=False) for _ in range(num_decoder_layers * 2)])
        self.past_key_values_ca = torch.nn.ParameterList([torch.nn.Parameter(torch.
↪ones((num_beams,self.num_attention_heads_per_partition,max_length,model_config.d_kv),
↪dtype=torch.float32), requires_grad=False) for _ in range(num_decoder_layers * 2)])

    def update_past(self, past_key_values):
        new_past_sa = []
        new_past_ca = []
        for past_layer in past_key_values:
            new_past_layer = list(past_layer)
            for i in range(len(new_past_layer[:2])):
                new_past_layer[i] = past_layer[i][:, :, 1:]
            new_past_sa += [new_past_layer[:2],]
            new_past_ca += [new_past_layer[2:],]
        return new_past_sa, new_past_ca

    def reorder_cache(self, past_key_values, beam_idx):
        for i in range(len(past_key_values)):
            gather_index = beam_idx.view([beam_idx.shape[0],1,1,1]).expand_as(past_key_
↪values[i])
            past_key_values[i] = torch.gather(past_key_values[i], dim = 0, index=gather_
↪index)
        return past_key_values

    def forward(self,
                input_ids,
                decoder_attention_mask,

```

(continues on next page)

(continued from previous page)

```

        encoder_hidden_states,
        encoder_attention_mask,
        beam_idx,
        beam_scores,
        **kwargs):

    if self.num_beams > 1:
        # We reorder the cache based on the beams selected in each iteration.
    ↪ Required step for beam search.
        past_key_values_sa = self.reorder_cache(self.past_key_values_sa, beam_idx)
        past_key_values_ca = self.reorder_cache(self.past_key_values_ca, beam_idx)
    else:
        # We do not need to reorder for greedy sampling
        past_key_values_sa = self.past_key_values_sa
        past_key_values_ca = self.past_key_values_ca

    # The cache is stored in a flatten form. We order the cache per layer before
    ↪ passing it to the decoder.
    # Each layer has 4 tensors, so we group by 4.
    past_key_values = [[*past_key_values_sa[i*2:i*2+2], *past_key_values_ca[i*2:
    ↪ i*2+2]] for i in range(0, int(len(past_key_values_ca)/2))]

    decoder_output = self.decoder(
        input_ids=input_ids,
        attention_mask=decoder_attention_mask,
        past_key_values=past_key_values,
        encoder_hidden_states=encoder_hidden_states,
        encoder_attention_mask=encoder_attention_mask,
        use_cache=True,
        output_attentions=False,
        output_hidden_states=False)

    last_hidden_state = decoder_output['last_hidden_state']
    past_key_values = decoder_output['past_key_values']

    if self.config.tie_word_embeddings:
        # Rescale output before projecting on vocab
        # See https://github.com/tensorflow/mesh/blob/
    ↪ fa19d69eafc9a482aff0b59ddd96b025c0cb207d/mesh_tensorflow/transformer/transformer.py
    ↪ #L586
        last_hidden_state = last_hidden_state * (self.model_dim**-0.5)

    lm_logits = self.lm_head(last_hidden_state)

    past_key_values_sa, past_key_values_ca = self.update_past(past_key_values)

    # We flatten the cache to a single array. This is required for the input output
    ↪ aliasing to work
    past_key_values_sa = [vec for kv_per_layer in past_key_values_sa for vec in kv_
    ↪ per_layer]
    past_key_values_ca = [vec for kv_per_layer in past_key_values_ca for vec in kv_
    ↪ per_layer]

```

(continues on next page)

(continued from previous page)

```

if self.device == "cpu":
    self.past_key_values_sa = past_key_values_sa
    self.past_key_values_ca = past_key_values_ca

# We calculate topk inside the wrapper
next_token_logits = lm_logits[:, -1, :]

if self.num_beams > 1:
    # This section of beam search is run outside the decoder in the huggingface_
    ↪t5 implementation.
    # To maximize the computation within the neuron device, we move this within_
    ↪the wrapper
    logit_max, _ = torch.max(next_token_logits, dim=-1, keepdim=True)
    logsumexp = torch.log(torch.exp(next_token_logits - logit_max).sum(dim=-1, ↪
    ↪keepdim=True))
    next_token_scores = next_token_logits - logit_max - logsumexp
    next_token_scores = next_token_scores + beam_scores[:, None].expand_as(next_
    ↪token_scores)

    # reshape for beam search
    vocab_size = next_token_scores.shape[-1]
    next_token_scores = next_token_scores.view(self.batch_size, self.num_beams * ↪
    ↪vocab_size)
    next_token_scores = next_token_scores * 1

    # Sample 2 next tokens for each beam (so we have some spare tokens and match_
    ↪output of beam search)
    next_token_scores, next_tokens = torch.topk(
        next_token_scores, 2 * self.num_beams, dim=1, largest=True, sorted=True
    )

    next_indices = torch.div(next_tokens, vocab_size, rounding_mode="floor")
    next_tokens = next_tokens % vocab_size

    return [next_token_scores, next_tokens, next_indices] + past_key_values_sa + ↪
    ↪past_key_values_ca
else:
    # Greedy
    next_tokens = torch.argmax(next_token_logits, dim=-1)
    return [next_tokens] + past_key_values_sa + past_key_values_ca

```

Now let's create a T5 model wrapper to make it compatible with our traced encoder and decoder.

There are two reasons for having this wrapper,

1. The encoder and decoder traces can only be invoked with positional arguments. But the HuggingFace transformers code is written with keyword arguments. So we override the functions that invoke encoder and decoder to call with positional arguments.
2. The generate() function in the NeuronGenerationMixin performs cache update within the CPU. As we are handling the cache within the DecoderWrapper, we disable the cache update on CPU.

3. The topK computation to determine the next tokens for beam search was moved into the decoder wrapper. So, we need to override the huggingface's beam search implementation to accept the next tokens and the beam scores from the decoder.

Let's also override the `generate()` function so that it will initialize the cache using the cache initializer before starting the greedy decoding.

```
[4]: import torch
import torch_xla.core.xla_model as xm

from transformers import T5Tokenizer, T5ForConditionalGeneration
from transformers.modeling_outputs import BaseModelOutput, Seq2SeqLMOutput
from transformers.models.t5.modeling_t5 import T5Stack, T5LayerCrossAttention
from transformers.generation.utils import ModelOutput
from typing import Any, Dict, List, Optional, Tuple, Union
from transformers.generation.beam_search import BeamScorer, BeamSearchScorer

from optimum.neuron.generation import NeuronGenerationMixin

from transformers.generation.logits_process import (
    LogitsProcessorList,
)
from transformers.generation.stopping_criteria import (
    MaxLengthCriteria,
    MaxTimeCriteria,
    StoppingCriteriaList,
    validate_stopping_criteria,
)

from transformers.generation.utils import (
    BeamSearchOutput,
    GreedySearchOutput,
)

class T5Wrapper(T5ForConditionalGeneration, NeuronGenerationMixin):

    def _prepare_encoder_decoder_kwargs_for_generation(
        self,
        inputs_tensor: torch.Tensor,
        model_kwargs,
        model_input_name: Optional[str] = None
    ) -> Dict[str, Any]:
        encoder = self.get_encoder()
        model_kwargs["encoder_outputs"]: ModelOutput = encoder(inputs_tensor, model_
        kwargs["attention_mask"])
        return model_kwargs

    # Override to cut the input_ids to just last token
    def prepare_inputs_for_generation(
        self,
        input_ids,
        past_key_values=None,
        attention_mask=None,
        head_mask=None,
```

(continues on next page)



(continued from previous page)

```

        decoder_head_mask=None,
        decoder_attention_mask=None,
        cross_attn_head_mask=None,
        use_cache=None,
        encoder_outputs=None,
        **kwargs,
    ):
        # cut decoder_input_ids as past is cached
        input_ids = input_ids[:, -1:]

        return {
            "decoder_input_ids": input_ids,
            "past_key_values": past_key_values,
            "encoder_outputs": encoder_outputs,
            "attention_mask": attention_mask,
            "head_mask": head_mask,
            "decoder_head_mask": decoder_head_mask,
            "decoder_attention_mask": decoder_attention_mask,
            "cross_attn_head_mask": cross_attn_head_mask,
            "use_cache": use_cache,
        }

    ...

    We update the cache in the decoder trace, so lets override the _update_model_
    ↪kwargs_for_xla_generation in NeuronGenerationMixin
    ...
    def _update_model_kwargs_for_xla_generation(
        self,
        model_kwargs: Dict[str, Any],
        batch_size: int,
        is_encoder_decoder: bool = False,
        standardize_cache_format: bool = False,
        max_length: Optional[int] = None,
        seq_length: Optional[int] = None,
        use_cache: bool = True,
    ) -> Dict[str, Any]:

        def _update_attention(model_kwargs, is_encoder_decoder):
            """Updates the appropriate attention mask -- encoder-decoder models use_
            ↪`decoder_attention_mask`"""

            attention_mask_name = "decoder_attention_mask" if is_encoder_decoder else
            ↪"attention_mask"
            attention_mask = model_kwargs.pop(attention_mask_name)
            attention_mask_update_slice = torch.ones(
                (batch_size, 1), dtype=attention_mask.dtype, device=attention_mask.device
            )
            attention_mask = torch.cat([attention_mask[:, 1:], attention_mask_update_
            ↪slice], dim=-1)
            mask = {attention_mask_name: attention_mask}
            return mask

```

(continues on next page)

(continued from previous page)

```

mask = _update_attention(model_kwargs, is_encoder_decoder)
# sets the updated variables (mask and past_key_values)
model_kwargs.update(mask)

# Set a mock cache tensor
model_kwargs["past_key_values"] = torch.tensor([])

return model_kwargs

def _reorder_cache(self, past_key_values, beam_idx):
    '''
        This is needed for beam search and not greedy sampling
        We reorder the cache within the trace so we can skip it in modelling_t5.py.
    ↪ So we override the _reorder_cache
    '''
    self.beam_idx = beam_idx
    return past_key_values

def generate(self,
            tokenizer: T5Tokenizer,
            prompt: str,
            max_length: int,
            num_beams: int,
            num_return_sequences: int,
            device: str):

    batch_encoding = tokenizer(prompt, max_length=max_length, truncation=True,
    ↪ padding='max_length',
                                return_tensors="pt")

    past_key_values = self.encoder(batch_encoding['input_ids'], batch_encoding[
    ↪ 'attention_mask'])

    decoder_attention_mask = torch.cat([torch.zeros((1, max_length-1), dtype=torch.
    ↪ int32),
                                torch.ones((1, 1), dtype=torch.int32)],
    ↪ axis=1)

    # copy the new cache state to the decoder
    if device == "xla":
        for state, tensor in zip(self.decoder.parameters(), past_key_values):
            state.copy_(tensor)
    else:
        # First half of the cache is self attention and the rest is cross attention
        self.decoder.past_key_values_sa = past_key_values[:len(past_key_values)//2]
        self.decoder.past_key_values_ca = past_key_values[len(past_key_values)//2:]

    output = super().generate(**batch_encoding,
                            max_length=max_length,
                            num_beams=num_beams,
                            num_return_sequences=num_return_sequences,
                            do_sample=False,

```

(continues on next page)

(continued from previous page)

```

        use_cache=True,
        decoder_attention_mask=decoder_attention_mask,
        encoder_outputs={"last_hidden_state": torch.ones((1,128,
→1))}) # Pass fake encoder_outputs so the transformers code will not invoke the encoder
        return output

def forward(
    self,
    attention_mask: Optional[torch.FloatTensor] = None,
    decoder_input_ids: Optional[torch.LongTensor] = None,
    decoder_attention_mask: Optional[torch.BoolTensor] = None,
    encoder_outputs: Optional[Tuple[Tuple[torch.Tensor]]] = None,
    beam_scores = None,
    **kwargs
) -> Union[Tuple[torch.FloatTensor], Seq2SeqLMOutput]:

    hidden_states = encoder_outputs["last_hidden_state"]

    if not hasattr(self, 'beam_idx'):
        # Inferring the number of beams from the attention mask
        num_beams = attention_mask.shape[0]
        self.beam_idx = torch.arange(0, num_beams, dtype=torch.int64)

    decoder_outputs = self.decoder(
        decoder_input_ids,
        decoder_attention_mask,
        hidden_states,
        attention_mask,
        self.beam_idx,
        beam_scores
    )

    # lm_logits = decoder_outputs[0]
    next_token_scores = decoder_outputs[0]
    next_tokens = decoder_outputs[1]
    next_indices = decoder_outputs[2]

    return next_token_scores, next_tokens, next_indices

def beam_search(
    self,
    input_ids: torch.LongTensor,
    beam_scorer: BeamScorer,
    logits_processor: Optional[LogitsProcessorList] = None,
    stopping_criteria: Optional[StoppingCriteriaList] = None,
    max_length: Optional[int] = None,
    pad_token_id: Optional[int] = None,
    eos_token_id: Optional[Union[int, List[int]]] = None,
    output_attentions: Optional[bool] = None,
    output_hidden_states: Optional[bool] = None,
    output_scores: Optional[bool] = None,
    return_dict_in_generate: Optional[bool] = None,

```

(continues on next page)

(continued from previous page)

```

        synced_gpus: Optional[bool] = False,
        seq_length: Optional[int] = None,
        **model_kwargs,
    ) -> Union[BeamSearchOutput, torch.LongTensor]:

        logits_processor = logits_processor if logits_processor is not None else
↳LogitsProcessorList()
        stopping_criteria = stopping_criteria if stopping_criteria is not None else
↳StoppingCriteriaList()
        pad_token_id = pad_token_id if pad_token_id is not None else self.generation_
↳config.pad_token_id
        eos_token_id = eos_token_id if eos_token_id is not None else self.generation_
↳config.eos_token_id
        if isinstance(eos_token_id, int):
            eos_token_id = [eos_token_id]
        output_scores = output_scores if output_scores is not None else self.generation_
↳config.output_scores
        output_attentions = (
            output_attentions if output_attentions is not None else self.generation_
↳config.output_attentions
        )
        output_hidden_states = (
            output_hidden_states if output_hidden_states is not None else self.
↳generation_config.output_hidden_states
        )

        batch_size = len(beam_scorer._beam_hyps)
        num_beams = beam_scorer.num_beams

        batch_beam_size, cur_len = input_ids.shape

        # Overwrite cur_len
        cur_len = seq_length

        if num_beams * batch_size != batch_beam_size:
            raise ValueError(
                f"Batch dimension of `input_ids` should be {num_beams * batch_size}, but
↳is {batch_beam_size}."
            )

        # init attention / hidden states / scores tuples
        scores = () if (return_dict_in_generate and output_scores) else None
        beam_indices = (
            tuple(
                () for _ in range(batch_beam_size)
            ) if (return_dict_in_generate and
↳output_scores) else None
        )

        # initialise score of first beam with 0 and the rest with -1e9. This makes sure
↳that only tokens
        # of the first beam are considered to avoid sampling the exact same tokens
↳across all beams.
        # beam_scores = torch.zeros((batch_size, num_beams), dtype=torch.float,

```

(continues on next page)

(continued from previous page)

```

↪device=input_ids.device)
    beam_scores_device = "cpu"
    beam_scores = torch.zeros((batch_size, num_beams), dtype=torch.float,
↪device=beam_scores_device)
    beam_scores[:, 1:] = -1e9
    beam_scores = beam_scores.view((batch_size * num_beams,))

    while True:
        # prepare model inputs
        # From max_length-sized input_ids, select first
        # cur_len - 1 values.
        update_indices = torch.stack(
            [torch.arange(input_ids.size(0)), torch.tensor(cur_len - 1).repeat(input_
↪ids.size(0))], dim=-1
        )
        input_ids_ = input_ids[update_indices[:, 0], update_indices[:, 1], None]
        model_inputs = self.prepare_inputs_for_generation(input_ids_, **model_kwargs)

        next_token_scores, next_tokens, next_indices = self(
            **model_inputs,
            return_dict=True,
            output_attentions=output_attentions,
            output_hidden_states=output_hidden_states,
            beam_scores=beam_scores
        )

        # stateless
        beam_outputs = beam_scorer.process(
            input_ids.to("cpu")[:, :cur_len],
            next_token_scores.to("cpu"),
            next_tokens.to("cpu"),
            next_indices.to("cpu"),
            pad_token_id=pad_token_id,
            eos_token_id=eos_token_id,
            beam_indices=beam_indices,
        )

        beam_scores = beam_outputs["next_beam_scores"]
        beam_next_tokens = beam_outputs["next_beam_tokens"]
        beam_idx = beam_outputs["next_beam_indices"]

        update_indices = torch.stack(
            [torch.arange(batch_beam_size), torch.tensor(cur_len - 1).repeat(batch_
↪beam_size)], dim=-1
        )
        update_indices_2 = torch.stack(
            [torch.arange(batch_beam_size), torch.tensor(cur_len).repeat(batch_beam_
↪size)], dim=-1
        )
        # First select beam_indices
        device = input_ids.device
        beam_idx_device = beam_idx.to(device=input_ids.device)

```

(continues on next page)

(continued from previous page)

```

input_ids[:, :] = input_ids[beam_idx_device.long(), :]

# Then append new tokens
input_ids[update_indices_2[:, 0], update_indices_2[:, 1], None] = beam_next_
↳tokens.unsqueeze(-1).to(device).to(torch.long)
input_ids = input_ids * 1 # Hack to materialize tensor

# update generated ids, model inputs, and length for next step
model_kwargs = self._update_model_kwargs_for_xla_generation(
    model_kwargs,
    batch_size=batch_beam_size,
    is_encoder_decoder=self.config.is_encoder_decoder,
    max_length=stopping_criteria.max_length,
    seq_length=cur_len,
    use_cache=model_kwargs["use_cache"],
)
if model_kwargs["past_key_values"] is not None:
    model_kwargs["past_key_values"] = self._reorder_cache(model_kwargs["past_
↳key_values"], beam_idx.to(torch.int64))

if return_dict_in_generate and output_scores:
    beam_indices = tuple((beam_indices[beam_idx[i]] + (beam_idx[i],) for i_
↳in range(len(beam_indices))))

# increase cur_len
cur_len = cur_len + 1

# stop when each sentence is finished, or if we exceed the maximum length
stop_criterion_1 = beam_scorer.is_done
if isinstance(stopping_criteria, list):
    if len(stopping_criteria) == 1:
        stopping_criteria = stopping_criteria[0]

# Cases that can be handled in XLA without requiring
# non-padded input_ids
if isinstance(stopping_criteria, MaxLengthCriteria):
    stop_criterion_2 = cur_len >= stopping_criteria.max_length
elif isinstance(stopping_criteria, MaxTimeCriteria):
    stop_criterion_2 = stopping_criteria(input_ids, scores)
else:
    # Other cases will be handled on CPU
    batch_size, _ = input_ids.shape
    input_ids_cpu = input_ids.to("cpu")
    mask = torch.cat(
        [torch.ones(batch_size, cur_len), torch.zeros(batch_size, input_ids.
↳shape[1] - cur_len)], dim=1
    ).bool()
    input_ids_cpu = torch.masked_select(input_ids_cpu, mask).reshape((batch_
↳size, cur_len))
    scores_cpu = scores.to("cpu") if torch.is_tensor(scores) else scores
    stop_criterion_2 = stopping_criteria(input_ids_cpu, scores_cpu)

```

(continues on next page)

(continued from previous page)

```

        if stop_criterion_1 or stop_criterion_2:
            if not synced_gpus:
                break
            else:
                this_peer_finished = True

sequence_outputs = beam_scorer.finalize(
    input_ids.to("cpu"),
    beam_scores.to("cpu"),
    next_tokens.to("cpu"),
    next_indices.to("cpu"),
    pad_token_id=pad_token_id,
    eos_token_id=eos_token_id,
    max_length=stopping_criteria.max_length,
    beam_indices=beam_indices,
)

for k, v in sequence_outputs.items():
    if type(v) == torch.Tensor:
        sequence_outputs[k] = sequence_outputs[k].to(input_ids.device)

return sequence_outputs["sequences"]

def greedy_search(
    self,
    input_ids: torch.LongTensor,
    logits_processor: Optional[LogitsProcessorList] = None,
    stopping_criteria: Optional[StoppingCriteriaList] = None,
    max_length: Optional[int] = None,
    pad_token_id: Optional[int] = None,
    eos_token_id: Optional[Union[int, List[int]]] = None,
    output_attentions: Optional[bool] = None,
    output_hidden_states: Optional[bool] = None,
    output_scores: Optional[bool] = None,
    return_dict_in_generate: Optional[bool] = None,
    seq_length: Optional[int] = int,
    streamer: Optional["BaseStreamer"] = None,
    **model_kwargs,
) -> Union[GreedySearchOutput, torch.LongTensor]:
    """
    Overriding greedy sampling to use next tokens returned from neuron device
    instead of logits.
    """
    # init values
    logits_processor = logits_processor if logits_processor is not None else
    LogitsProcessorList()
    use_cache = model_kwargs["use_cache"] if "use_cache" in model_kwargs else False
    stopping_criteria = stopping_criteria if stopping_criteria is not None else
    StoppingCriteriaList()
    pad_token_id = pad_token_id if pad_token_id is not None else self.generation_
    config.pad_token_id

```

(continues on next page)

(continued from previous page)

```

        eos_token_id = eos_token_id if eos_token_id is not None else self.generation_
↪ config.eos_token_id
        if isinstance(eos_token_id, int):
            eos_token_id = [eos_token_id]
            eos_token_id_tensor = torch.tensor(eos_token_id).to(input_ids.device) if eos_
↪ token_id is not None else None
            output_scores = output_scores if output_scores is not None else self.generation_
↪ config.output_scores
            output_attentions = (
                output_attentions if output_attentions is not None else self.generation_
↪ config.output_attentions
            )
            output_hidden_states = (
                output_hidden_states if output_hidden_states is not None else self.
↪ generation_config.output_hidden_states
            )

        # init attention / hidden states / scores tuples
        scores = () if (return_dict_in_generate and output_scores) else None
        decoder_attentions = () if (return_dict_in_generate and output_attentions) else
↪ None
        cross_attentions = () if (return_dict_in_generate and output_attentions) else
↪ None
        decoder_hidden_states = () if (return_dict_in_generate and output_hidden_states)
↪ else None

        # keep track of which sequences are already finished
        unfinished_sequences = torch.ones(input_ids.shape[0], dtype=torch.long,
↪ device=input_ids.device)

        this_peer_finished = False # used by synced_gpus only
        while True:

            # prepare model inputs
            # From max_length-sized input_ids, select first
            # seq_length - 1 values.

            if model_kwargs.get("past_key_values") is None:
                input_ids_ = input_ids[:, :seq_length]
            else:
                update_indices = torch.stack(
                    [torch.arange(input_ids.size(0)), torch.tensor(seq_length - 1).
↪ repeat(input_ids.size(0))],
                    dim=-1,
                )
                input_ids_ = input_ids[update_indices[:, 0], update_indices[:, 1], None]

            model_inputs = self.prepare_inputs_for_generation(input_ids_, **model_kwargs)

            # forward pass to get next token
            output = self(

```

(continues on next page)



(continued from previous page)

```

        **model_inputs,
        return_dict=True,
        output_attentions=output_attentions,
        output_hidden_states=output_hidden_states,
    )
    next_tokens = output[0]

    # finished sentences should have their next token be a padding token
    if eos_token_id is not None:
        if pad_token_id is None:
            raise ValueError("If `eos_token_id` is defined, make sure that `pad_
↪token_id` is defined.")
        next_tokens = next_tokens * unfinished_sequences + pad_token_id * (1 -
↪unfinished_sequences)

    # update generated ids, model inputs, and length for next step

    batch_size, _ = input_ids.shape
    update_indices = torch.stack(
        [torch.arange(batch_size), torch.tensor(seq_length).repeat(batch_size)],
↪dim=-1
    )
    input_ids[update_indices[:, 0], update_indices[:, 1]] = next_tokens[:]
    model_kwargs = self._update_model_kwargs_for_xla_generation(
        model_kwargs,
        batch_size=batch_size,
        is_encoder_decoder=self.config.is_encoder_decoder,
        max_length=stopping_criteria.max_length,
        seq_length=seq_length,
        use_cache=use_cache,
    )

    seq_length += 1

    # if eos_token was found in one sentence, set sentence to finished
    if eos_token_id_tensor is not None:
        unfinished_sequences = unfinished_sequences.mul(
            next_tokens.tile(eos_token_id_tensor.shape[0], 1).ne(eos_token_id_
↪tensor.unsqueeze(1)).prod(dim=0)
        )

    # stop when each sentence is finished, or if we exceed the maximum length
    stop_criterion_1 = unfinished_sequences.max() == 0

    if isinstance(stopping_criteria, list):
        if len(stopping_criteria) == 1:
            stopping_criteria = stopping_criteria[0]

    # Cases that can be handled in XLA without requiring
    # non-padded input_ids
    if isinstance(stopping_criteria, MaxLengthCriteria):
        stop_criterion_2 = seq_length >= stopping_criteria.max_length

```

(continues on next page)

(continued from previous page)

```

elif isinstance(stopping_criteria, MaxTimeCriteria):
    stop_criterion_2 = stopping_criteria(input_ids, scores)
else:
    # Other cases will be handled on CPU
    batch_size, _ = input_ids.shape
    mask = torch.cat(
        [torch.ones(batch_size, seq_length), torch.zeros(batch_size, input_
↪ ids.shape[1] - seq_length)],
        dim=1,
    ).bool()
    input_ids_cpu = torch.masked_select(input_ids, mask).reshape((batch_size,
↪ seq_length)).to("cpu")
    scores_cpu = scores.to("cpu") if torch.is_tensor(scores) else scores
    stop_criterion_2 = stopping_criteria(input_ids_cpu, scores_cpu)

if stop_criterion_1 or stop_criterion_2:
    this_peer_finished = True

if this_peer_finished:
    break

if streamer is not None:
    streamer.end()

return input_ids

```

Now let's test inference on CPU with all the wrappers before tracing.

[5]: # Let's set some run parameters

```

model_name = "t5-large"
num_beams = 1
num_return_sequences = 1
max_length = 128

```

[6]: from transformers import T5Tokenizer

```

prompt="translate English to German: Lets eat good food."

tokenizer = T5Tokenizer.from_pretrained(model_name, model_max_length=max_length)
model = T5Wrapper.from_pretrained(model_name)

model.encoder = EncoderWrapper(model.encoder, model.decoder, model.config, num_beams,
↪ max_length, "cpu", num_beams)
setattr(model.encoder, 'main_input_name', 'input_ids') # Attribute required by beam_
↪ search

model.decoder = DecoderWrapper(decoder=model.decoder,
                                lm_head=model.lm_head,
                                model_config=model.config,

```

(continues on next page)

(continued from previous page)

```

        num_beams=num_beams,
        max_length=max_length,
        device="cpu")

output = model.generate(tokenizer=tokenizer,
                        prompt=prompt,
                        max_length=max_length,
                        num_beams=num_beams,
                        num_return_sequences=num_return_sequences,
                        device="cpu")

results = [tokenizer.decode(t, skip_special_tokens=True) for t in output]

print('Results:')
for i, summary in enumerate(results):
    print(i + 1, summary)

```

```

Results:
1 Lassen Sie uns gutes Essen essen.

```

Now that the wrappers are running as expected, let's trace the encoder, and decoder. To trace these functions, we pass the function and a sample input to the trace function. The result of the trace stage will be a static executable where the operations to be run upon inference are determined during compilation. This means that when inferring, the resulting Neuron model must be executed with tensors that are the exact same shape as those provided at compilation time. If a model is given a tensor at inference time whose shape does not match the tensor given at compilation time, an error will occur.

The decoder wrapper returns the new state of the cache as an output which is copied back to the CPU. As the cache is a large tensor, copying it to and from the XLA device for each decoder invocation will significantly slow down the inference. Instead, we can use input output aliasing, a feature of `torch_neuronx` to keep these tensors on device rather than copying back to the CPU. To use input output aliasing, we need to map the outputs to input parameters while tracing.

```

[ ]: import torch
import torch_neuronx

from transformers import T5Tokenizer, T5ForConditionalGeneration

def trace_encoder(model: T5ForConditionalGeneration,
                  tokenizer: T5Tokenizer,
                  max_length: int,
                  num_beams: int):

    # Trace encoder
    batch_encoding = tokenizer("translate English to German: Lets go home now",
                              max_length=max_length, truncation=True, padding='max_
↪length', return_tensors="pt")
    input_ids = batch_encoding['input_ids']
    attention_mask = batch_encoding['attention_mask']

    encoder = EncoderWrapper(model.encoder, model.decoder, model.config, num_beams, max_
↪length, "xla", num_beams)

```

(continues on next page)

(continued from previous page)

```

        traced_encoder = torch_neuronx.trace(encoder, (input_ids, attention_mask), compiler_
↪workdir="/tmp/encoder/")
        setattr(traced_encoder, 'main_input_name', 'input_ids') # Attribute required by ↪
↪beam search

        return traced_encoder

def trace_decoder(model: T5ForConditionalGeneration,
                  num_beams: int,
                  max_length: int):

    decoder = DecoderWrapper(decoder=model.decoder,
                             lm_head=model.lm_head,
                             model_config=model.config,
                             num_beams=num_beams,
                             max_length=max_length,
                             device="xla")

    # We create mock inputs so we can trace the decoder
    decoder_input_ids = torch.ones((num_beams, 1), dtype=torch.int64)
    decoder_attention_mask = torch.ones((num_beams, max_length), dtype=torch.int32)
    encoder_attention_mask = torch.ones((num_beams, max_length), dtype=torch.int64)
    encoder_hidden_states = torch.ones((num_beams, max_length, model.config.d_model), ↪
↪dtype=torch.float32)

    beam_idx = torch.arange(0, num_beams, dtype=torch.int64)
    beam_scores = torch.zeros((num_beams,), dtype=torch.float)

    num_outputs_from_trace = 3 if num_beams > 1 else 1

    aliases = {}
    for i in range(len(decoder.past_key_values_sa)):
        aliases[decoder.past_key_values_sa[i]] = i + num_outputs_from_trace
    for i in range(len(decoder.past_key_values_ca)):
        aliases[decoder.past_key_values_ca[i]] = len(decoder.past_key_values_sa) + i + ↪
↪num_outputs_from_trace

    traced_decoder = torch_neuronx.trace(decoder, (
        decoder_input_ids,
        decoder_attention_mask,
        encoder_hidden_states,
        encoder_attention_mask,
        beam_idx,
        beam_scores,
    ), input_output_aliases=aliases, compiler_workdir="/tmp/decoder/")

    return traced_decoder

tokenizer = T5Tokenizer.from_pretrained(model_name, model_max_length=max_length)
model = T5ForConditionalGeneration.from_pretrained(model_name)

```

(continues on next page)

(continued from previous page)

```
# We enable this flag to ensure model uses attention key value caching
model.config.use_cache = True

traced_encoder = trace_encoder(model, tokenizer, max_length, num_beams)
traced_decoder = trace_decoder(model, num_beams, max_length)

torch.jit.save(traced_encoder, "TracedEncoder.pt")
torch.jit.save(traced_decoder, "TracedDecoder.pt")
```

## Run inference with greedy decoding

Now that we have the traced model, let's use it for inference.

```
[8]: runtime = torch.classes.neuron.Runtime()
runtime.initialize()
runtime.set_default_neuron_cores(0, 1)

tokenizer = T5Tokenizer.from_pretrained(model_name)
model = T5Wrapper.from_pretrained(model_name)

model.encoder = torch.jit.load("TracedEncoder.pt")
# Attribute required by beam search
setattr(model.encoder, 'main_input_name', 'input_ids')

model.decoder = torch.jit.load("TracedDecoder.pt")
torch_neuronx.move_trace_to_device(model.decoder, 0)

output = model.generate(tokenizer=tokenizer,
                        prompt="translate English to German: Lets eat good food.",
                        max_length=max_length,
                        num_beams=num_beams,
                        num_return_sequences=num_return_sequences,
                        device="xla")

results = [tokenizer.decode(t, skip_special_tokens=True) for t in output]

print('Results:')
for i, summary in enumerate(results):
    print(i + 1, summary)
```

```
Results:
1 Lassen Sie uns gutes Essen essen.
```

## Run inference with beam search

```
[ ]: # Let's set some run parameters for beam search

model_name = "t5-large"
num_beams = 4
num_return_sequences = 4
max_length = 128

tokenizer = T5Tokenizer.from_pretrained(model_name, model_max_length=max_length)
model = T5ForConditionalGeneration.from_pretrained(model_name)
model.config.use_cache = True

traced_encoder = trace_encoder(model, tokenizer, max_length, num_beams)
traced_decoder = trace_decoder(model, num_beams, max_length)

torch.jit.save(traced_encoder, "TracedEncoder.pt")
torch.jit.save(traced_decoder, "TracedDecoder.pt")
```

```
[10]: tokenizer = T5Tokenizer.from_pretrained(model_name)
model = T5Wrapper.from_pretrained(model_name)

model.encoder = torch.jit.load("TracedEncoder.pt")
# Attribute required by beam search
setattr(model.encoder, 'main_input_name', 'input_ids')

model.decoder = torch.jit.load("TracedDecoder.pt")
torch_neuronx.move_trace_to_device(model.decoder, 0)

output = model.generate(tokenizer=tokenizer,
                        prompt="translate English to German: Lets eat good food.",
                        max_length=max_length,
                        num_beams=num_beams,
                        num_return_sequences=num_return_sequences,
                        device="xla")

results = [tokenizer.decode(t, skip_special_tokens=True) for t in output]

print('Results:')
for i, summary in enumerate(results):
    print(i + 1, summary)
```

Results:

```
1 Lassen Sie uns gutes Essen essen.
2 Lassen Sie uns gutes Essen zu essen.
3 Lassen Sie uns essen gutes Essen.
4 Lassen Sie uns gutes Essen.
```

- HuggingFace pretrained BERT tutorial [\[html\]](#) [\[notebook\]](#)
- TorchServe tutorial [\[html\]](#)
- LibTorch C++ tutorial (for torch-neuron and torch-neuronx) [\[html\]](#)

- Torchvision ResNet50 tutorial [\[html\]](#) [\[notebook\]](#)
- T5 inference tutorial [\[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### Additional Examples (torch-neuronx)

- [AWS Neuron Samples GitHub Repository](#)
- [Transformers Neuron GitHub samples](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### API Reference Guide (torch-neuronx)

*This document is relevant for:* Inf2, Trn1, Trn2

### PyTorch NeuronX Tracing API for Inference

`torch_neuronx.trace(func, example_inputs, *, input_output_aliases={}, compiler_workdir=None, compiler_args=None, partitioner_config=None, inline_weights_to_neff=True, cpu_backend=False)`

Trace and compile operations in the `func` by executing it using `example_inputs`.

This function is similar to a `torch.jit.trace()` since it produces a `ScriptModule` that can be saved with `torch.jit.save()` and reloaded with `torch.jit.load()`. The resulting module is an optimized fused graph representation of the `func` that is *only* compatible with Neuron.

Tracing a module produces a more efficient *inference-only* version of the model. XLA Lazy Tensor execution should be used during training. See: [Comparison of Traced Inference versus XLA Lazy Tensor Inference \(torch-neuronx\)](#)

**Warning:** Currently this only supports NeuronCore-v2 type instances (e.g. trn1, inf2). To compile models compatible with NeuronCore-v1 (e.g. inf1), please see [torch\\_neuron.trace\(\)](#)

#### Parameters

- **func** (*Module, callable*) – The function/module that that will be run using the `example_inputs` arguments in order to record the computation graph.
- **example\_inputs** (*Tensor, tuple[Tensor]*) – A tuple of example inputs that will be passed to the `func` while tracing.

#### Keyword Arguments

- **input\_output\_aliases** (*dict*) – Marks input tensors as state tensors which are device tensors.
- **compiler\_workdir** (*str*) – Work directory used by *neuronx-cc*. This can be useful for debugging and/or inspecting intermediary *neuronx-cc* outputs
- **compiler\_args** (*str*, *list[str]*) – List of strings representing *neuronx-cc* compiler arguments. See *Neuron Compiler CLI Reference Guide (neuronx-cc)* for more information about compiler options.
- **partitioner\_config** (*PartitionerConfig*) – A *PartitionerConfig* object, which can be optionally supplied if there are unsupported ops in the model that need to be partitioned out to CPU.
- **inline\_weights\_to\_neff** (*bool*) – A boolean indicating whether the weights should be inlined to the NEFF. If set to False, weights will be separated from the NEFF. The default is True.
- **cpu\_backend** (*bool*) – A boolean indicating whether CPU should be used for tracing. If set to True, tracing can be done completely on CPU. This keyword needs to be used with the *compiler\_args* option to set the `--target` flag. The default is False.

**Returns**

The traced *ScriptModule* with the embedded compiled Neuron graph. Operations in this module will execute on Neuron.

**Return type**

*ScriptModule*

**Warning:** Behavior Change! The use of using args for kwargs is deprecated starting from release 2.15.0 (*torch-neuronx==1.13.1.1.12.0*). The current behavior is that a warning will be raised, but *torch\_neuronx.trace()* will attempt to infer the keyword arguments. This is likely to become an error in future releases, so to avoid the warning/error, assign kwargs as kwargs and not args.

**Notes**

This function records operations using *torch-xla* to create a *HloModule* representation of the *func*. This fixed graph representation is compiled to the Neuron Executable File Format (NEFF) using the *neuronx-cc* compiler. The NEFF binary executable is embedded into an optimized *ScriptModule* for *torchscript* execution.

In contrast to a regular *torch.jit.trace()* that produces a graph of many separate operations, tracing with Neuron produces a graph with a single fused operator that is executed entirely on device. In *torchscript* this appears as a stateful *neuron::Model* component with an associated *neuron::forward\** operation.

Tracing can be performed on any EC2 machine with sufficient memory and compute resources, but inference can only be executed on a Neuron instance.

Unlike some devices (such as *torch-xla*) that use *to()* to move *Parameter* and *Tensor* data between CPU and device, upon loading a Neuron traced *ScriptModule*, the model binary executable is automatically moved to a *NeuronCore*. When the underlying *neuron::Model* is initialized after tracing or upon *torch.jit.load()*, it is loaded to a Neuron device without specifying a device or *map\_location* argument.

**Warning:** One small exception is models traced with *inline\_weights\_to\_neff=False*. For these models, the NEFF is loaded onto the *NeuronCore* automatically, but the weights are not moved automatically. To move the weights to the *NeuronCore*, call *torch\_neuronx.move\_trace\_to\_device()*. If this is not done,



a performance penalty is incurred per inference, because on every inference call, the weights move from CPU to Neuron.

Furthermore, the Neuron traced `ScriptModule` expects to consume CPU tensors and produces CPU tensors. The underlying operation performs all data transfers to and from the Neuron device without explicit data movement. This is a significant difference from the training XLA device mechanics since XLA operations are no longer required to be recorded after a trace. See: [Developer Guide for Training with PyTorch NeuronX](#)

By *default*, when multiple NeuronCores are available, every Neuron traced model `ScriptModule` within in a process is loaded to each available NeuronCore in round-robin order. This is useful at deployment to fully utilize the Neuron hardware since it means that multiple calls to `torch.jit.load()` will attempt to load to each available NeuronCore in linear order. The default start device is chosen according to the [NeuronX Runtime Configuration](#).

A traced Neuron module has limitations that are not present in regular torch modules:

- **Fixed Control Flow:** Similar to `torch.jit.trace()`, tracing a model with Neuron statically preserves control flow (i.e. `if/for/while` statements) and will not re-evaluate the branch conditions upon inference. If a model result is based on data-dependent control flow, the traced function may produce inaccurate results.
- **Fixed Input Shapes:** After a function has been traced, the resulting `ScriptModule` will always expect to consume tensors of the same shape. If the tensor shapes used at inference differs from the tensor shapes used in the `example_inputs`, this will result in an error. See: [Running inference on variable input shapes with bucketing](#).
- **Fixed Tensor Shapes:** The intermediate tensors within the `func` must always stay the same shape for the same shaped inputs. This means that certain operations which produce data-dependent sized tensors are not supported. For example, `nonzero()` produces a different tensor shape depending on the input data.
- **Fixed Data Types:** After a model has been traced, the input, output, and intermediate data types cannot be changed without recompiling.
- **Device Compatibility:** Due to Neuron using a specialized compiled format (NEFF), a model traced with Neuron can no longer be executed in any non-Neuron environment.
- **Operator Support:** If an operator is unsupported by `torch-xla`, then this will throw an exception.

## Examples

### Function Compilation

```
import torch
import torch_neuronx
def func(x, y):
    return 2 * x + y
example_inputs = torch.rand(3), torch.rand(3)
# Runs `func` with the provided inputs and records the tensor operations
trace = torch_neuronx.trace(func, example_inputs)
# `trace` can now be run with the TorchScript interpreter or saved
# and loaded in a Python-free environment
torch.jit.save(trace, 'func.pt')
# Executes on a NeuronCore
loaded = torch.jit.load('func.pt')
loaded(torch.rand(3), torch.rand(3))
```

### Module Compilation

```

import torch
import torch_neuronx
import torch.nn as nn
class Model(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = nn.Conv2d(1, 1, 3)
    def forward(self, x):
        return self.conv(x) + 1
model = Model()
model.eval()
example_inputs = torch.rand(1, 1, 3, 3)
# Traces the forward method and constructs a `ScriptModule`
trace = torch_neuronx.trace(model, example_inputs)
torch.jit.save(trace, 'model.pt')
# Executes on a NeuronCore
loaded = torch.jit.load('model.pt')
loaded(torch.rand(1, 1, 3, 3))

```

#### Weight Separated Module

```

import torch
import torch_neuronx
import torch.nn as nn

class Model(nn.Module):

    def __init__(self):
        super().__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        return self.conv(x) + 1

model = Model()
model.eval()

example_inputs = torch.rand(1, 1, 3, 3)

# Traces the forward method and constructs a `ScriptModule`
trace = torch_neuronx.trace(model, example_inputs, inline_weights_to_neff=False)

# Model can be saved like a normally traced model
torch.jit.save(trace, 'model.pt')

# Executes on a NeuronCore like a normally traced model
loaded = torch.jit.load('model.pt')
torch_neuronx.move_trace_to_device(loaded, 0) # necessary for performance
loaded(torch.rand(1, 1, 3, 3))

```

#### CPU Compilation

On CPU:

```

import torch
import torch_neuronx
import torch.nn as nn
class Model(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = nn.Conv2d(1, 1, 3)
    def forward(self, x):
        return self.conv(x) + 1
model = Model()
model.eval()
example_inputs = torch.rand(1, 1, 3, 3)
# Traces the forward method on CPU, compiling for Trn1
trace = torch_neuronx.trace(model, example_inputs, compiler_args="--target trn1",
    ↪cpu_backend=True)
torch.jit.save(trace, 'model.pt')
# Move model.pt to a Neuron instance

```

On Neuron:

```

import torch
import torch_neuronx
import torch.nn as nn

loaded = torch.jit.load('model.pt')
loaded(torch.rand(1, 1, 3, 3))

```

---

**Note:** Weight Separated models can have its weights replaced via the `torch_neuronx.replace_weights` API.

---

## Moving a Traced Module to a Neuron Core

**Warning:** This function will be deprecated in a future release, and instead, `torch_neuronx.experimental.set_neuron_cores()` will move out of experimental, and become a stable API.

`torch_neuronx.move_trace_to_device(trace, device_id)`

This function moves a model traced with `torch_neuronx.trace()`, to a Neuron Core. Here are some reasons to use this function:

1. Explicit control of device placement for models

By default, the Neuron Runtime assigns neffs to devices in a Round Robin manner, meaning it will allocate a neff onto Neuron Core 0, then 1, 2, and then loop around.

2. Allocating Weights onto the Neuron Core for Weight Separated models.

This is necessary for performance reasons. If this is not done, the weights would remain on CPU and would need to move to device on every inference call, which is an expensive operation.

### Parameters

- **trace** (*ScriptModule*) – This is the torchscript model returned from `torch_neuronx.trace()`
- **device\_id** (*int*) – The Neuron Core to move the traced model to. This number will need to be between 0 to the max number of NCs on the instance - 1. For example, a trn1.32xlarge has 32 Neuron Cores, so the acceptable values are from 0-31.

**Returns**

Nothing, the movement of the model happens in-place.

**Return type**

None

## Autobucketing

---

**Note:** See `neuronx_distributed.parallel_model_trace()` for the API to use the autobucketing feature along with tensor parallelism.

---

```
class torch_neuronx.BucketModelConfig(bucket_kernel, *, shared_state_buffer=None,
                                       shared_state_buffer_preprocessor=None, func_kwargs=None)
```

This object contains configuration data for how buckets are selected based on input via the `bucket_kernel`.

This also supports the concept of a shared buffer between bucket models. You can use this to define how the shared buffer can be manipulated to be fed as input to a bucket model via the `shared_state_buffer_preprocessor`. Details on how these are defined are found below.

**Parameters**

**bucket\_kernel** (*callable*) – A function that returns a new TorchScript function. The TorchScript function has been adapted to the TorchScript representation using `torch.jit.script()`. This new function takes in a list of input tensors and outputs a list of tensors and an index tensor.

**Keyword Arguments**

- **shared\_state\_buffer** (*Optional[List[torch.Tensor]]*) – A list of tensors that is used as the initial values for a shared state for bucket models via aliasing.
- **shared\_state\_buffer\_preprocessor** (*Optional[Callable]*) – Similar to `bucket_kernel`, this is a function that returns a new TorchScript function that has been adapted to the TorchScript representation using `torch.jit.script()`. This new TorchScript function takes in 3 arguments: an n-dimensional integer list representing a list of tensor shapes, the `state_buffer` list of tensors, and a tensor representing the bucket index. This function outputs a reshaped `state_buffer` to be supplied to the bucket model. If `shared_state_buffer_preprocessor` is not supplied when `shared_state_buffer` is supplied, the preprocessor returns the full `shared_state_buffer`.
- **func\_kwargs** (*Optional[Union[Dict[str, Any], List[Any]]]*) – A single dictionary or a list of dictionaries that can be used to supply custom arguments to the function supplied to the `func` argument in `torch_neuronx.bucket_model_trace()`. If you are using a list of dictionaries, verify that `func_kwargs` equals the bucket degree, or number of buckets. By default `func_kwargs` is None, which means no arguments.

**Returns**

The `torch_neuronx.BucketModelConfig` with the configuration defining bucket selection for inputs and shared buffers.

**Return type**

`BucketModelConfig`

```
torch_neuronx.bucket_model_trace(func, example_inputs, bucket_config, compiler_workdir=None,
                                compiler_args=None)
```

This function traces a single model with multiple `example_inputs` and a `bucket_config` object to produce a single compiled model that can take in multiple input shapes. This trace function is very similar to `torch_neuronx.trace()`, but it has a few key differences:

1. In this case, `func` does not take in a `Model`. Instead, it takes in a function that returns a tuple containing a `Model` and `input_output_aliases`. This is like `neuronx_distributed.parallel_model_trace()`, and is done for the same reason, which is that bucket models are traced in parallel.
2. Instead of taking in one input, the function takes in multiple inputs in the form of a list. For example, `[torch.rand(128, 128), torch.rand(256, 256)]`.
3. The `bucket_config` argument is of type `torch_neuronx.BucketModelConfig()`, which defines how an input is mapped to a bucket. For more details, see the `torch_neuronx.BucketModelConfig()` API Reference. You can use this for a variety of bucketing applications, such as sequence length bucketing for language models or image resolution bucketing for computer vision models.

Apart from the aforementioned differences, the rest of the function behaves similarly to `torch_neuronx.trace()`. You can save the model with `torch.jit.save()` and load it with `torch.jit.load()`.

#### Parameters

- **func** (*Module, callable*) – This is a function that returns a `Model` object and a dictionary of states, or `input_output_aliases`. Similar to `neuronx_distributed.parallel_model_trace()`, this API calls this function inside each worker and runs trace against them. Note: This differs from the `torch_neuronx.trace` where the `torch_neuronx.trace` requires a model object to be passed.
- **example\_inputs** (*List[Union[Tensor, tuple[Tensor]]]*) – A list of possible inputs to the bucket model.
- **bucket\_config** (*BucketModelConfig*) – The config object that defines bucket selection behavior.

#### Keyword Arguments

- **compiler\_workdir** (*str*) – Work directory used by `neuronx-cc`. This can be useful for debugging and inspecting intermediary `neuronx-cc` outputs.
- **compiler\_args** (*str, list[str]*) – List of strings representing `neuronx-cc` compiler arguments. See *Neuron Compiler CLI Reference Guide (neuronx-cc)* for more information about compiler options.

#### Returns

The traced `ScriptModule` with the embedded compiled Neuron graphs for each bucket model. Operations in this module will execute on Neuron.

#### Return type

*ScriptModule*

**Warning:** If you receive the Too Many Open Files error message, increase the ulimit via `ulimit -n 65535`. There is a limitation in `torch_xla`'s `xmp.spawn` function when dealing with large amounts of data.

The developer guide for Autobucketing is located here, which contains an example usage of autobucketing with BERT.

## Dynamic Batching

`torch_neuronx.dynamic_batch(neuron_script)`

Enables a compiled Neuron model to be called with variable sized batches.

When tracing with Neuron, usually a model can only consume tensors that are the same size as the example tensor used in the `torch_neuronx.trace()` call. Enabling dynamic batching allows a model to consume inputs that may be either smaller or larger than the original trace-time tensor size. Internally, dynamic batching splits & pads an input batch into chunks of size equal to the original trace-time tensor size. These chunks are passed to the underlying model(s). Compared to serial inference, the expected runtime scales by  $\text{ceil}(\text{inference\_batch\_size} / \text{trace\_batch\_size}) / \text{neuron\_cores}$ .

This function modifies the `neuron_script` network in-place. The returned result is a reference to the modified input.

Dynamic batching is only supported by chunking inputs along the 0th dimension. A network that uses a non-0 batch dimension is incompatible with dynamic batching. Upon inference, inputs whose shapes differ from the compile-time shape in a non-0 dimension will raise a `ValueError`. For example, take a model was traced with a single example input of size `[2, 3, 5]`. At inference time, when dynamic batching is enabled, a batch of size `[3, 3, 5]` is *valid* while a batch of size `[2, 7, 5]` is *invalid* due to changing a non-0 dimension.

Dynamic batching is only supported when the 0th dimension is the same size for all inputs. For example, this means that dynamic batching would not be applicable to a network which consumed two inputs with shapes `[1, 2]` and `[3, 2]` since the 0th dimension is different. Similarly, at inference time, the 0th dimension batch size for all inputs must be identical otherwise a `ValueError` will be raised.

### Required Arguments

#### Parameters

**neuron\_script** (*ScriptModule*) – The neuron traced *ScriptModule* with the embedded compiled neuron graph. This is the output of `torch_neuronx.trace()`.

#### Returns

The traced *ScriptModule* with the embedded compiled neuron graph. The same type as the input, but with `dynamic_batch` enabled in the neuron graph.

#### Return type

*ScriptModule*

```
import torch
import torch_neuronx
import torch.nn as nn

class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        return self.conv(x) + 1

n = Net()
n.eval()

inputs = torch.rand(1, 1, 3, 3)
inputs_batch_8 = torch.rand(8, 1, 3, 3)
```

(continues on next page)

(continued from previous page)

```
# Trace a neural network with input batch size of 1
neuron_net = torch_neuronx.trace(n, inputs)

# Enable the dynamic batch size feature so the traced network
# can consume variable sized batch inputs
neuron_net_dynamic_batch = torch_neuronx.dynamic_batch(neuron_net)

# Run inference on inputs with batch size of 8
# different than the batch size used in compilation (tracing)
output_batch_8 = neuron_net_dynamic_batch(inputs_batch_8)
```

## Graph Partitioner

```
torch_neuronx.PartitionerConfig(*, trace_kwargs=None, model_support_percentage_threshold=0.5,
                                min_subgraph_size=-1, max_subgraph_count=-1, ops_to_partition=None,
                                analyze_parameters=None)
```

Allows for Neuron to trace a model with unsupported operators and partition these operators to CPU.

This model will contain subgraphs of Neuron and CPU submodules, but it is executed like one model, and can be saved and loaded like one model as well.

The graph partitioner is customized using this class, and is *only* enabled (disabled by default) from the `torch_neuronx.trace` API by setting `partitioner_config` keyword argument to this class. Below are the various configuration options.

### Parameters

- **trace\_kwargs** (*Dict*) – Used if you need to pass trace kwargs to the Neuron subgraphs, such as the `compiler_workdir` and/or `compiler_args`. The default is `None` corresponding to the default trace args.
- **model\_support\_percentage\_threshold** (*float*) – A number between 0 to 1 representing the maximum allowed percentage of operators that must be supported. If the max is breached, the function will throw a `ValueError`. Default is `0.5` (i.e 50% of operators must be supported by Neuron)
- **min\_subgraph\_size** (*int*) – The minimum number of operators in a subgraph. Can be `>= 1` or `== -1`. If `-1`, minimum subgraph size is not checked (i.e no minimum). If `>= 1`, each subgraph must contain at least that many operators. If not, the graph partitioner will throw a `ValueError`.
- **max\_subgraph\_count** (*int*) – The maximum number of subgraphs in the partitioned model. Can be `>= 1` or `== -1`. If `-1`, max subgraph count is not checked (i.e no maximum). If `>= 1`, the partitioned model must contain at most that many subgraphs. If not, the graph partitioner will throw a `ValueError`.
- **ops\_to\_partition** (*Set[str]*) – This is a set of strings of this structure “aten::<operator>”. These are operators that will be partitioned to CPU regardless of Neuron support. The default is `None` (i.e no additional operators will be partitioned).
- **analyze\_parameters** (*Dict*) – This is a dictionary of kwargs used in `torch_neuronx.analyze()`. NOTE: Not all kwargs in `torch_neuronx.analyze()` are supported in the graph partitioner. The following kwargs in `analyze` are supported for use in the graph partitioner.

a) `compiler_workdir`

b) `additional_ignored_ops`

c) `max_workers`

The default is `None`, corresponding to the default `analyze` arguments.

### Returns

The `PartitionerConfig` with the configuration for the graph partitioner.

### Return type

`PartitionerConfig`

## Examples

This example demonstrates using the graph partitioner.

The below model is a simple MLP model with sorted log softmax output. The sort operator, `torch.sort()` or `aten::sort`, is not supported by `neuronx-cc` at this time, so the graph partitioner will partition out the sort operator to CPU.

```
import torch
import torch_neuronx
import torch.nn as nn

import logging

# adjust logger level to see what the partitioner is doing
logger = logging.getLogger("Neuron")

class MLP(nn.Module):
    def __init__(
        self, input_size=28 * 28, output_size=10, layers=[4096, 2048]
    ):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)
        self.relu = nn.ReLU()

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = self.relu(f1)
        f2 = self.fc2(r1)
        r2 = self.relu(f2)
        f3 = self.fc3(r2)
        out = torch.log_softmax(f3, dim=1)
        sort_out, _ = torch.sort(out)
        return sort_out

n = MLP()
n.eval()

inputs = torch.rand(32, 784)

# Configure the graph partitioner with the default values
```

(continues on next page)



(continued from previous page)

```

partitioner_config = torch_neuronx.PartitionerConfig()

# Trace a neural network with graph partitioner enabled
neuron_net = torch_neuronx.trace(n, inputs, partitioner_config=partitioner_config)

# Run inference on the partitioned model
output = neuron_net(inputs)

```

**Note:** Dynamic batching has a case-by-case support with partitioned models, because it is highly dependent on how the final partition scheme looks like.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch Neuron (torch-neuronx) Weight Replacement API for Inference

`torch_neuronx.replace_weights(neuron_model, weights)`

Replaces the weights in a Neuron Model with split weights. This function will emit a warning of the supplied Neuron model does not contain any separated weights.

**Warning:** The below API is only applicable for models traced with the parameter `inline_weights_to_neff=False`, which is True by default. See `torch_neuronx.trace()` for details.

### Parameters

- **neuron\_model** (*RecursiveScriptModule*) – A Neuron model compiled with split weights
- **weights** (*Module, Dict[str, Tensor]*) – Either the original model with the new weights, or the state\_dict of a model.

### Returns

None, this function performs the weight replacement inline.

### Return type

None

## Examples

*Using a model*

```

import torch
import torch_neuronx

class Network(torch.nn.Module):
    def __init__(self, hidden_size=4, layers=3) -> None:
        super().__init__()

```

(continues on next page)

(continued from previous page)

```

        self.layers = torch.nn.Sequential(
            *(torch.nn.Linear(hidden_size, hidden_size) for _ in range(layers)))

    def forward(self, tensor):
        return self.layers(tensor)

# initialize two networks
network = Network()
network2 = Network()
network.eval()
network2.eval()

inp = torch.rand(2,4)

# trace weight separated model with first network
weight_separated_trace = torch_neuronx.trace(network,inp,inline_weights_to_
    ↪neff=False)

# replace with weights from second network
torch_neuronx.replace_weights(weight_separated_trace,network2.state_dict())

# get outputs from neuron and cpu networks
out_network2 = network2(inp)
out_neuron = weight_separated_trace(inp)

# check that they are equal
print(out_network2,out_neuron)

```

*Using safetensors*

The `safetensors` library is useful for storing/loading model tensors safely and quickly.

```

import torch
import torch_neuronx

from safetensors import safe_open
from safetensors.torch import save_model

class Network(torch.nn.Module):
    def __init__(self, hidden_size=4, layers=3) -> None:
        super().__init__()
        self.layers = torch.nn.Sequential(
            *(torch.nn.Linear(hidden_size, hidden_size) for _ in range(layers)))

    def forward(self, tensor):
        return self.layers(tensor)

# initialize two networks
network = Network()
network2 = Network()

```

(continues on next page)

(continued from previous page)

```

network.eval()
network2.eval()

inp = torch.rand(2,4)

# trace weight separated model with first network
weight_separated_trace = torch_neuronx.trace(network,inp,inline_weights_to_
↪neff=False)

# save network2 weights to safetensors
safetensor_path = f"{directory}/network2.safetensors"
save_model(network2,safetensor_path)

#load safetensors from network2 into traced_weight separated model
tensors = {}
with safe_open(safetensor_path,framework="pt") as f:
    for k in f.keys():
        tensors[k] = f.get_tensor(k)

# replace with weights from second network
torch_neuronx.replace_weights(weight_separated_trace,tensors)

# get outputs from neuron and cpu networks
out_network2 = network2(inp)
out_neuron = weight_separated_trace(inp)

# check that they are equal
print(out_network2,out_neuron)

```

---

**Note:** For non-safetensors models, use `torch.load` to load the model, and pass the model's `state_dict` inside like the first example.

---

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch NeuronX NeuronCore Placement APIs [Beta]

**Warning:** The following functionality is beta and **will not be supported** in future releases of the NeuronSDK. This module serves only as a preview for future functionality. In future releases, equivalent functionality may be moved directly to the `torch_neuronx` module and will no longer be available in the `torch_neuronx.experimental` module.

Functions which enable placement of `torch.jit.ScriptModule` to specific NeuronCores. Two sets of functions are provided which can be used interchangeably but have different performance characteristics and advantages:

- The `multicore_context()` & `neuron_cores_context()` functions are context managers that allow a model to be placed on a given NeuronCore *only* at `torch.jit.load()` time. These functions are the most efficient way of loading a model since the model is loaded directly to a NeuronCore. The alternative functions described below require that a model is unloaded from one core and then reloaded to another.

- The `set_multicore()` & `set_neuron_cores()` functions allow a model that has already been loaded to a NeuronCore to be moved to a different NeuronCore. This functionality is less efficient than directly loading a model to a NeuronCore within a context manager but allows device placement to be fully dynamic at runtime. This is analogous to the `torch.nn.Module.to()` function for device placement.

---

**Important:** A prerequisite to enable placement functionality is that the loaded `torch.jit.ScriptModule` has already been compiled with the `torch_neuronx.trace()` API. Attempting to place a regular `torch.nn.Module` onto a NeuronCore prior to compilation will do nothing.

---

`torch_neuronx.experimental.set_neuron_cores(trace: torch.jit.ScriptModule, start_nc: int = -1, nc_count: int = -1)`

Set the NeuronCore start/count for all Neuron subgraphs in a torch Module.

This will unload the model from an existing NeuronCore if it is already loaded.

*Requires Torch 1.8+*

#### Parameters

**trace** (*ScriptModule*) – A torch module which contains one or more Neuron subgraphs.

#### Keyword Arguments

- **start\_nc** (*int*) – The starting NeuronCore index where the Module is placed. The value -1 automatically loads to the optimal NeuronCore (least used). Note that this index is always relative to NeuronCores visible to this process.
- **nc\_count** (*int*) – The number of NeuronCores to use. The value -1 will load a model to exactly one NeuronCore. If `nc_count` is greater than one, the model will be replicated across multiple NeuronCores.

#### Raises

- **[RuntimeError]** – If the Neuron runtime cannot be initialized.
- **[ValueError]** – If the `nc_count` is an invalid number of NeuronCores.

## Examples

*Single Load:* Move a model to the first visible NeuronCore after loading.

```
model = torch.jit.load('example_neuron_model.pt')
torch_neuronx.experimental.set_neuron_cores(model, start_nc=0, nc_count=1)

model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 0
```

*Multiple Core Replication:* Replicate a model to 2 NeuronCores after loading. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores by running round-robin executions.

```
model = torch.jit.load('example_neuron_model.pt')
torch_neuronx.experimental.set_neuron_cores(model, start_nc=2, nc_count=2)

model(example) # Executes on NeuronCore 2
model(example) # Executes on NeuronCore 3
model(example) # Executes on NeuronCore 2
```

*Multiple Model Load:* Move and pin 2 models to separate NeuronCores. This causes each `torch.jit.ScriptModule` to always execute on a specific NeuronCore.

```
model1 = torch.jit.load('example_neuron_model.pt')
torch_neuronx.experimental.set_neuron_cores(model1, start_nc=2)

model2 = torch.jit.load('example_neuron_model.pt')
torch_neuronx.experimental.set_neuron_cores(model2, start_nc=0)

model1(example) # Executes on NeuronCore 2
model1(example) # Executes on NeuronCore 2
model2(example) # Executes on NeuronCore 0
model2(example) # Executes on NeuronCore 0
```

`torch_neuronx.experimental.set_multicore(trace: torch.jit.ScriptModule)`

Loads all Neuron subgraphs in a torch Module to all visible NeuronCores.

This loads each Neuron subgraph within a `torch.jit.ScriptModule` to multiple NeuronCores without requiring multiple calls to `torch.jit.load()`. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores for concurrent threadsafe inferences. Executions use a round-robin strategy to distribute across NeuronCores.

This will unload the model from an existing NeuronCore if it is already loaded.

*Requires Torch 1.8+*

#### Parameters

**trace** (*ScriptModule*) – A torch module which contains one or more Neuron subgraphs.

#### Raises

[[RuntimeError](#)] – If the Neuron runtime cannot be initialized.

## Examples

*Multiple Core Replication:* Move a model across all visible NeuronCores after loading. This allows a single `torch.jit.ScriptModule` to use all NeuronCores by running round-robin executions.

```
model = torch.jit.load('example_neuron_model.pt')
torch_neuronx.experimental.set_multicore(model)

model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 1
model(example) # Executes on NeuronCore 2
```

`torch_neuronx.experimental.neuron_cores_context(start_nc: int = -1, nc_count: int = -1)`

A context which sets the NeuronCore start/count for Neuron models loaded with `torch.jit.load()`.

This context manager may only be used when loading a model with `torch.jit.load()`. A model which has already been loaded into memory will not be affected by this context manager. Furthermore, after loading the model, inferences do not need to occur in this context in order to use the correct NeuronCores.

Note that this context is *not* threadsafe. Using multiple core placement contexts from multiple threads may not correctly place models.

#### Keyword Arguments

- **start\_nc** (*int*) – The starting NeuronCore index where the Module is placed. The value -1 automatically loads to the optimal NeuronCore (least used). Note that this index is always relative to NeuronCores visible to this process.
- **nc\_count** (*int*) – The number of NeuronCores to use. The value -1 will load a model to exactly one NeuronCore. If nc\_count is greater than one, the model will be replicated across multiple NeuronCores.

#### Raises

- **[RuntimeError]** – If the Neuron runtime cannot be initialized.
- **[ValueError]** – If the nc\_count is an invalid number of NeuronCores.

#### Examples

*Single Load:* Directly load a model from disk to the first visible NeuronCore.

```
with torch_neuronx.experimental.neuron_cores_context(start_nc=0, nc_count=1):
    model = torch.jit.load('example_neuron_model.pt') # Load must occur within the
    ↪ context

model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 0
```

*Multiple Core Replication:* Directly load a model from disk to 2 NeuronCores. This allows a single torch.jit.ScriptModule to use multiple NeuronCores by running round-robin executions.

```
with torch_neuronx.experimental.neuron_cores_context(start_nc=2, nc_count=2):
    model = torch.jit.load('example_neuron_model.pt') # Load must occur within the
    ↪ context

model(example) # Executes on NeuronCore 2
model(example) # Executes on NeuronCore 3
model(example) # Executes on NeuronCore 2
```

*Multiple Model Load:* Directly load 2 models from disk and pin them to separate NeuronCores. This causes each torch.jit.ScriptModule to always execute on a specific NeuronCore.

```
with torch_neuronx.experimental.neuron_cores_context(start_nc=2):
    model1 = torch.jit.load('example_neuron_model.pt') # Load must occur within
    ↪ the context

with torch_neuronx.experimental.neuron_cores_context(start_nc=0):
    model2 = torch.jit.load('example_neuron_model.pt') # Load must occur within
    ↪ the context

model1(example) # Executes on NeuronCore 2
model1(example) # Executes on NeuronCore 2
model2(example) # Executes on NeuronCore 0
model2(example) # Executes on NeuronCore 0
```

`torch_neuronx.experimental.multicore_context()`

A context manager which loads models to all visible NeuronCores for Neuron models loaded with torch.jit.load().

This loads each Neuron subgraph within a `torch.jit.ScriptModule` to multiple NeuronCores without requiring multiple calls to `torch.jit.load()`. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores for concurrent threadsafe inferences. Executions use a round-robin strategy to distribute across NeuronCores.

This context manager may only be used when loading a model with `torch.jit.load()`. A model which has already been loaded into memory will not be affected by this context manager. Furthermore, after loading the model, inferences do not need to occur in this context in order to use the correct NeuronCores.

Note that this context is *not* threadsafe. Using multiple core placement contexts from multiple threads may not correctly place models.

#### Raises

`[RuntimeError]` – If the Neuron runtime cannot be initialized.

### Examples

*Multiple Core Replication:* Directly load a model to all visible NeuronCores. This allows a single `torch.jit.ScriptModule` to use all NeuronCores by running round-robin executions.

```
with torch_neuronx.experimental.multicore_context():
    model = torch.jit.load('example_neuron_model.pt') # Load must occur within the_
    ↪ context

model(example) # Executes on NeuronCore 0
model(example) # Executes on NeuronCore 1
model(example) # Executes on NeuronCore 2
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## PyTorch NeuronX Analyze API for Inference

`torch_neuronx.analyze(func, example_inputs, compiler_workdir=None)`

Checks the support of the operations in the `func` by checking each operator against `neuronx-cc`.

#### Parameters

- **func** (*Module, callable*) – The function/module that that will be run using the `example_inputs` arguments in order to record the computation graph.
- **example\_inputs** (*Tensor, tuple[Tensor]*) – A tuple of example inputs that will be passed to the `func` while tracing.

#### Keyword Arguments

- **compiler\_workdir** (*str*) – Work directory used by `neuronx-cc`. This can be useful for debugging and/or inspecting intermediary `neuronx-cc` outputs
- **additional\_ignored\_ops** (*set*) – A set of aten operators to not analyze. Default is an empty set.
- **max\_workers** (*int*) – The max number of workers threads to spawn. The default is 4.
- **is\_hf\_transformers** (*bool*) – If the model is a huggingface transformers model, it is recommended to enable this option to prevent deadlocks. Default is `False`.

- **cleanup** (*bool*) – Specifies whether to delete the compiler artifact directories generated after running analyze. Default is False.

**Returns**

A JSON like Dict with the supported operators and their count, and unsupported operators with the failure mode and location of the operator in the python code.

**Return type**

Dict

**Examples**

*Fully supported model*

```
import json

import torch
import torch.nn as nn
import torch_neuronx

class MLP(nn.Module):
    def __init__(self, input_size=28*28, output_size=10, layers=[120,84]):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(layers[0], layers[1])
    def forward(self, x):
        f1 = self.fc1(x)
        r1 = self.relu(f1)
        f2 = self.fc2(r1)
        r2 = self.relu(f2)
        f3 = self.fc3(r2)
        return torch.log_softmax(f3, dim=1)

model = MLP()
ex_input = torch.rand([32,784])

model_support = torch_neuronx.analyze(model,ex_input)
print(json.dumps(model_support,indent=4))
```

```
{
  "torch_neuronx_version": "1.13.0.1.5.0",
  "neuronx_cc_version": "2.0.0.11796a0+24a26e112",
  "support_percentage": "100.00%",
  "supported_operators": {
    "aten::linear": 3,
    "aten::relu": 2,
    "aten::log_softmax": 1
  },
  "unsupported_operators": []
}
```

*Unsupported Model/Operator*



```
import json
import torch
import torch_neuronx
```

```
def fft(x):
    return torch.fft.fft(x)
```

```
model = fft
ex_input = torch.arange(4)
```

```
model_support = torch_neuronx.analyze(model, ex_input)
print(json.dumps(model_support, indent=4))
```

```
{
  "torch_neuronx_version": "1.13.0.1.5.0",
  "neuronx_cc_version": "2.0.0.11796a0+24a26e112",
  "support_percentage": "0.00%",
  "supported_operators": {},
  "unsupported_operators": [
    {
      "kind": "aten::fft_fft",
      "failureAt": "neuronx-cc",
      "call": "test.py(6): fft\n/home/ubuntu/testdir/venv/lib/python3.8/site-
→packages/torch_neuronx/xla_impl/analyze.py(35): forward\n/home/ubuntu/testdir/
→venv/lib/python3.8/site-packages/torch/nn/modules/module.py(1182): _slow_forward\n
→n/home/ubuntu/testdir/venv/lib/python3.8/site-packages/torch/nn/modules/module.
→py(1194): _call_impl\n/home/ubuntu/testdir/venv/lib/python3.8/site-packages/torch/
→jit/_trace.py(976): trace_module\n/home/ubuntu/testdir/venv/lib/python3.8/site-
→packages/torch/jit/_trace.py(759): trace\n/home/ubuntu/testdir/venv/lib/python3.8/
→site-packages/torch_neuronx/xla_impl/analyze.py(302): analyze\ntest.py(11):
→<module>\n",
      "opGraph": "graph(%x : Long(4, strides=[1], requires_grad=0, device=cpu),\n
→n      %neuron_4 : NoneType,\n      %neuron_5 : int,\n      %neuron_6 : NoneType):\n
→\n      %neuron_7 : ComplexFloat(4, strides=[1], requires_grad=0, device=cpu) = aten::
→fft_fft(%x, %neuron_4, %neuron_5, %neuron_6)\n      return (%neuron_7)\n"
    }
  ]
}
```

**Note:** the failureAt field can either be “neuronx-cc” or “Lowering to HLO”. If the field is “neuronx-cc”, then it indicates that the provided operator configuration failed to be compiled with neuronx-cc. This could either indicate that the operator configuration is unsupported, or there is a bug with that operator configuration.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch NeuronX DataParallel API

The `torch_neuronx.DataParallel()` Python API implements data parallelism on `ScriptModule` models created by *PyTorch NeuronX Tracing API for Inference*. This function is analogous to `DataParallel` in PyTorch. The *Data Parallel Inference on torch\_neuronx* application note provides an overview of how `torch_neuronx.DataParallel()` can be used to improve the performance of inference workloads on Inferentia.

`torch_neuronx.DataParallel(model, device_ids=None, dim=0, set_dynamic_batching=True)`

Applies data parallelism by replicating the model on available NeuronCores and distributing data across the different NeuronCores for parallelized inference.

By default, `DataParallel` will use all available NeuronCores allocated for the current process for parallelism. `DataParallel` will apply parallelism on `dim=0` if `dim` is not specified.

`DataParallel` automatically enables *dynamic batching* on eligible models if `dim=0`. Dynamic batching can be disabled using `torch_neuronx.DataParallel.disable_dynamic_batching()`, or by setting `set_dynamic_batching=False` when initializing the `DataParallel` object. If dynamic batching is not enabled, the batch size at compilation-time must be equal to the batch size at inference-time divided by the number of NeuronCores being used. Specifically, the following must be true when dynamic batching is disabled: `input.shape[dim] / len(device_ids) == compilation_input.shape[dim]`.

`torch.neuron.DataParallel()` requires PyTorch  $\geq 1.8$ .

### Required Arguments

#### Parameters

**model** (`ScriptModule`) – Model created by the *PyTorch NeuronX Tracing API for Inference* to be parallelized.

### Optional Arguments

#### Parameters

- **device\_ids** (`list`) – List of `int` or `'nc: #'` that specify the NeuronCores to use for parallelization (default: all NeuronCores). Refer to the *device\_ids note* for a description of how `device_ids` indexing works.
- **dim** (`int`) – Dimension along which the input tensor is scattered across NeuronCores (default `dim=0`).
- **set\_dynamic\_batching** (`bool`) – Whether to enable dynamic batching.

### Attributes

#### Parameters

- **num\_workers** (`int`) – Number of worker threads used for multithreaded inference (default: `2 * number of NeuronCores`).
- **split\_size** (`int`) – Size of the input chunks (default: `max(1, input.shape[dim] // number of NeuronCores)`).

`torch.neuron.DataParallel.disable_dynamic_batching()`

Disables automatic dynamic batching on the `DataParallel` module. See *Dynamic batching disabled* for example of how `DataParallel` can be used with dynamic batching disabled. Use as follows:

```
>>> model_parallel = torch_neuronx.DataParallel(model_neuron)
>>> model_parallel.disable_dynamic_batching()
```

**Note:** `device_ids` uses per-process NeuronCore granularity and zero-based indexing. Per-process granularity means that each Python process “sees” its own view of the world. Specifically, this means that `device_ids` only “sees” the NeuronCores that are allocated for the current process. Zero-based indexing means that each Python process will index its allocated NeuronCores starting at 0, regardless of the “global” index of the NeuronCores. Zero-based indexing makes it possible to redeploy the exact same code unchanged in different process. This behavior is analogous to the `device_ids` argument in the PyTorch `DataParallel` function.

As an example, assume `DataParallel` is run on an `inf2.48xlarge`, which contains 12 Inferentia chips each of which contains two NeuronCores:

- If `NEURON_RT_VISIBLE_CORES` is not set, a single process can access all 24 NeuronCores. Thus specifying `device_ids=["nc:0"]` will correspond to `chip0:core0` and `device_ids=["nc:13"]` will correspond to `chip6:core1`.
- However, if two processes are launched where: process 1 has `NEURON_RT_VISIBLE_CORES=0-11` and process 2 has `NEURON_RT_VISIBLE_CORES=12-23`, `device_ids=["nc:13"]` cannot be specified in either process. Instead, `chip6:core1` can only be accessed in process 2. Additionally, `chip6:core1` is specified in process 2 with `device_ids=["nc:1"]`. Furthermore, in process 1, `device_ids=["nc:0"]` would correspond to `chip0:core0`; in process 2 `device_ids=["nc:0"]` would correspond to `chip6:core0`.

## Examples

The following sections provide example usages of the `torch_neuronx.DataParallel()` module.

### Default usage

The default `DataParallel` use mode will replicate the model on all available NeuronCores in the current process. The inputs will be split on `dim=0`.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module
model_parallel = torch_neuronx.DataParallel(model_neuron)

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

## Specifying NeuronCores

The following example uses the `device_ids` argument to use the first three NeuronCores for DataParallel inference.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module, run on the first two NeuronCores
# Equivalent to model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1])
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=['nc:0', 'nc:1'])

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

## DataParallel with dim != 0

In this example we run DataParallel inference using two NeuronCores and `dim = 2`. Because `dim != 0`, dynamic batching is not enabled. Consequently, the DataParallel inference-time batch size must be two times the compile-time batch size.

```
import torch
import torch_neuronx

# Create an example model
class Model(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = torch.nn.Conv2d(3, 3, 3)

    def forward(self, x):
        return self.conv(x) + 1

model = Model()
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 8, 8])
model_neuron = torch_neuronx.trace(model, image)
```

(continues on next page)

(continued from previous page)

```

# Create the DataParallel module using 2 NeuronCores and dim = 2
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=[0, 1], dim=2)

# Create a batched input
# Note that image_batched.shape[dim] / len(device_ids) == image.shape[dim]
batch_size = 2 * 8
image_batched = torch.rand([1, 3, batch_size, 8])

# Run inference with a batched input
output = model_parallel(image_batched)

```

## Dynamic batching

In the following example, we use the `torch_neuronx.DataParallel()` module to run inference using several different batch sizes without recompiling the Neuron model.

```

import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module
model_parallel = torch_neuronx.DataParallel(model_neuron)

# Create batched inputs and run inference on the same model
batch_sizes = [2, 3, 4, 5, 6]
for batch_size in batch_sizes:
    image_batched = torch.rand([batch_size, 3, 224, 224])

    # Run inference with a batched input
    output = model_parallel(image_batched)

```

## Dynamic batching disabled

In the following example, we use `torch_neuronx.DataParallel.disable_dynamic_batching()` to disable dynamic batching. We provide an example of a batch size that will not work when dynamic batching is disabled as well as an example of a batch size that does work when dynamic batching is disabled.

```

import torch
import torch_neuronx
from torchvision import models

```

(continues on next page)

(continued from previous page)

```

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module and use 2 NeuronCores
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=[0, 1], dim=0)

# Disable dynamic batching
model_parallel.disable_dynamic_batching()

# Create a batched input (this won't work)
batch_size = 4
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will fail because dynamic batching is disabled and
# image_batched.shape[dim] / len(device_ids) != image.shape[dim]
# output = model_parallel(image_batched)

# Create a batched input (this will work)
batch_size = 2
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will work because
# image_batched.shape[dim] / len(device_ids) == image.shape[dim]
output = model_parallel(image_batched)

```

*This document is relevant for: Inf2, Trn1, Trn2*

## API Reference Guide (torch-neuronx)

- *PyTorch NeuronX Tracing API for Inference*
- *PyTorch Neuron (torch-neuronx) Weight Replacement API for Inference*
- *PyTorch NeuronX NeuronCore Placement APIs [Beta]*
- *PyTorch NeuronX Analyze API for Inference*
- *PyTorch NeuronX DataParallel API*
- `torch_neuronx_lazy_async_load_api`

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer Guide (torch-neuronx)

This document is relevant for: Inf2, Trn1, Trn2

### NeuronCore Allocation and Model Placement for Inference (torch-neuronx)

This programming guide describes the how to allocate NeuronCores to processes and place models onto specific NeuronCores. The models in this guide are expected to have been traced with with `torch_neuronx.trace()`.

**Warning:** This guide is **not** applicable to NeuronCore placement using XLA LazyTensor device execution. See: [Comparison of Traced Inference versus XLA Lazy Tensor Inference \(torch-neuronx\)](#)

In order of precedence, the recommendation is to use the following placement techniques:

1. For nearly all regular models, default core placement should be used to take control of all cores for a single process.
2. For applications using multiple processes, default core placement should be used in conjunction with `NEURON_RT_NUM_CORES` ([Default Core Allocation & Placement](#))
3. For more granular control, then the beta explicit placement APIs may be used ([Explicit Core Placement \[Beta\]](#)).

#### Table of Contents

- [NeuronCore Allocation and Model Placement for Inference \(torch-neuronx\)](#)
  - [Default Core Allocation & Placement](#)
    - \* [Example: Default](#)
    - \* [Example: NEURON\\_RT\\_NUM\\_CORES](#)
    - \* [Example: NEURON\\_RT\\_VISIBLE\\_CORES](#)
    - \* [Example: Multiple Processes](#)
  - [Explicit Core Placement \[Beta\]](#)
    - \* [Example: Manual Core Selection](#)
    - \* [Example: Automatic Multicore](#)

The following guide will assume a machine with 8 NeuronCores:

- NeuronCores will use the notation `nc0`, `nc1`, etc.
- Models will use the notation `m0`, `m1` etc.

NeuronCores and model allocations will be displayed in the following format:

The actual cores that are visible to the process can be adjusted according to the [NeuronX Runtime Configuration](#).

Unlike `torch-neuron` (with `neuron-cc`) instances, `torch-neuronx` (with `neuronx-cc`) does not support [NeuronCore Pipeline](#). This simplifies model core allocations since it means that model pipelines will likely not span across multiple NeuronCores.

## Default Core Allocation & Placement

The most basic requirement of an inference application is to be able to place a single model on a single NeuronCore. More complex applications may use multiple NeuronCores or even multiple processes each executing different models. The important thing to note about designing an inference application is that a single NeuronCore will always be allocated to a single process. *Processes do not share NeuronCores.* Different configurations can be used to ensure that an application process has enough NeuronCores allocated to execute its model(s):

- Default: A process will attempt to take ownership of **all NeuronCores** visible on the instance. This should be used when an instance is only running a single inference process since no other process will be allowed to take ownership of any NeuronCores.
- `NEURON_RT_NUM_CORES`: Specify the **number of NeuronCores** to allocate to the process. This places no restrictions on which NeuronCores will be used, however, the resulting NeuronCores will always be contiguous. This should be used in multi-process applications where each process should only use a subset of NeuronCores.
- `NEURON_RT_VISIBLE_CORES`: Specifies exactly **which NeuronCores** are allocated to the process by index. Similar to `NEURON_RT_NUM_CORES`, this can be used in multi-process applications where each process should only use a subset of NeuronCores. This provides more fined-grained controls over the exact NeuronCores that are allocated to a given process.

See the [NeuronX Runtime Configuration](#) for more environment variable details.

### Example: Default

#### Python Script:

```
import torch
import torch_neuronx

m0 = torch.jit.load('model.pt') # Loads to nc0
m1 = torch.jit.load('model.pt') # Loads to nc1
```

With no environment configuration, the process will take ownership of all NeuronCores. In this example, only two of the NeuronCores are used by the process and the remaining are allocated but left idle.

### Example: `NEURON_RT_NUM_CORES`

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '2'
```

#### Python Script:

```
import torch
import torch_neuronx

m0 = torch.jit.load('model.pt') # Loads to nc0
m1 = torch.jit.load('model.pt') # Loads to nc1
```

Since there is no other process on the instance, only the first 2 NeuronCores will be acquired by the process. Models load in a simple linear order to the least used NeuronCores.



### Example: NEURON\_RT\_VISIBLE\_CORES

#### Environment Setup:

```
export NEURON_RT_VISIBLE_CORES = '4-5'
```

#### Python Script:

```
import torch
import torch_neuronx

m0 = torch.jit.load('model.pt') # Loads to nc4
m1 = torch.jit.load('model.pt') # Loads to nc5
```

Unlike NEURON\_RT\_NUM\_CORES, setting the visible NeuronCores allows the process to take control of a specific contiguous set. This allows an application to have a more fine-grained control of where models will be placed.

### Example: Multiple Processes

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '2'
```

#### Python Script:

```
import torch
import torch_neuronx

m0 = torch.jit.load('model.pt') # Loads to nc0
m1 = torch.jit.load('model.pt') # Loads to nc1
```

In this example, if the script is run **twice**, the following allocations will be made:

Note that each process will take ownership of as many NeuronCores as is specified by the NEURON\_RT\_NUM\_CORES configuration.

### Explicit Core Placement [Beta]

The torch\_neuronx framework allows can be found in the [PyTorch NeuronX NeuronCore Placement APIs \[Beta\]](#) documentation.

### Example: Manual Core Selection

The most direct usage of the placement APIs is to manually select the start NeuronCore that each model is loaded to.

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '4'
```

#### Python Script:

```
import torch
import torch_neuronx

# NOTE: Order of loads does NOT matter
with torch_neuronx.experimental.neuron_cores_context(start_nc=3):
    m0 = torch.jit.load('model.pt') # Loads to nc3

with torch_neuronx.experimental.neuron_cores_context(start_nc=0, nc_count=2):
    m1 = torch.jit.load('model.pt') # Loads replicas to nc0 and nc1

example = torch.rand(1, 3, 224, 224)

m1(example) # Executes on nc3
m1(example) # Executes on nc3

m0(example) # Executes on nc0
m0(example) # Executes on nc1
m0(example) # Executes on nc0
```

### Example: Automatic Multicore

Using explicit core placement it is possible to replicate a model to multiple NeuronCores simultaneously. This means that a single model object within python can utilize all available NeuronCores (or NeuronCores allocated to the process).

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '8'
```

#### Python Script:

```
import torch
import torch_neuronx

with torch_neuronx.experimental.multicore_context():
    m0 = torch.jit.load('model.pt') # Loads replications to nc0-nc7

example = torch.rand(1, 3, 224, 224)

m0(example) # Executes on nc0
m0(example) # Executes on nc1
```

To make full use of a model that has been loaded to multiple NeuronCores, multiple threads should be used to run inferences in parallel.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Comparison of Traced Inference versus XLA Lazy Tensor Inference (torch-neuronx)

### Table of contents

- *Introduction*
- *XLA Lazy Tensor Inference Mechanics*
- *Traced Inference Mechanics*
- *Traced Inference Advantages*
- *Summary*

### Introduction

Using torch-neuronx, there are two ways that a model can be executed for inference:

- **XLA LazyTensor Inference:** A model is executed on Neuron by calling `to()` to move `Parameter` and `Tensor` data using the `xm.xla_device()`. Executing operations uses torch *LazyTensor* to record, compile, and execute the graph. These are the same mechanisms used for *training*.
- **(Recommended) Traced Inference:** A model is traced prior to inference using the `trace()` API. This trace is similar to `torch.jit.trace()` but instead creates a Neuron-specific *TorchScript* artifact. This artifact provides improved performance and portability compared to XLA *LazyTensor* inference.

### XLA Lazy Tensor Inference Mechanics

XLA *LazyTensor* inference uses Just-In-Time (JIT) compilation for Neuron execution.

XLA Device execution uses the built-in torch-xla functionality with torch *LazyTensor* to record torch operations using the `xm.xla_device()`. The graph of operations is sent to the *neuronx-cc* compiler upon calling `xm.mark_step()`. Finally the compiled graph is transferred to a NeuronCore and executed in the Neuron backend.

The initial model inference will be very slow since the model binary file in the Neuron Executable File Format (NEFF) will need to be generated by the compiler. Upon each subsequent call to a model, the application will re-execute the python, rebuild the graph, and check a cache to see if an existing NEFF file is available for the given graph before attempting to recompile.

The process of recording graph operations in python can become a bottleneck for otherwise fast models. This overhead will always have an effect on performance regardless of model size but may be less noticeable on larger models. Note that this XLA *LazyTensor* execution performance may improve significantly with new torch features in the future.

### Example

#### Fixed Shape Example

```
import torch
import torch_neuronx
import torch_xla.core.xla_model as xm

# Create XLA device
```

(continues on next page)

(continued from previous page)

```

device = xm.xla_device()

# Load example model and inputs to Neuron device
model = torch.nn.Sequential(
    torch.nn.Linear(784, 120),
    torch.nn.ReLU(),
    torch.nn.Linear(120, 10),
    torch.nn.Softmax(dim=-1),
)
model.eval()
model.to(device)
example = torch.rand((1, 784), device=device)

# Inference
with torch.no_grad():
    result = model(example)
    xm.mark_step() # Compilation occurs here
    print(result.cpu())

```

## Dynamic Shape Example

The following is an example of a model that dynamically changes the sequence length and batch size of the input token ID tensor to trigger recompilations. This kind of workflow would require padding when using traced inference.

```

import torch
import torch_neuronx
import torch_xla.core.xla_model as xm

# Create XLA device
device = xm.xla_device()

# Load example model and inputs to Neuron device
model = torch.nn.Sequential(
    torch.nn.Embedding(num_embeddings=30522, embedding_dim=512),
    torch.nn.Linear(512, 128),
    torch.nn.ReLU(),
    torch.nn.Linear(128, 2),
    torch.nn.Softmax(dim=-1),
)
model.eval()
model.to(device)

token_ids_1 = torch.tensor([
    [1, 28, 748, 0],
]) # shape: [1, 4]
token_ids_2 = torch.tensor([
    [1, 13087, 10439, 1990, 18912, 0],
    [1, 12009, 7849, 2509, 3500, 0],
]) # shape: [2, 6]

# Inference

```

(continues on next page)

(continued from previous page)

```

with torch.no_grad():

    # First compilation/inference
    result = model(token_ids_1)
    xm.mark_step()
    print(result.cpu()) # shape: [1, 4, 2]

    # Recompile occurs here since token_ids_2 is a different shape. This infer
    # would have failed if the model had been traced with shape [1, 4]
    result = model(token_ids_2)
    xm.mark_step()
    print(result.cpu()) # shape: [2, 6, 2]

```

## Traced Inference Mechanics

Traced inference uses Ahead-Of-Time (AOT) compilation for Neuron execution.

Similar to XLA *Lazy Tensor* inference, `trace()` uses the operation recording mechanisms provided by `torch-xla` to build the graph structure. This graph structure is also sent to the `neuronx-cc` compiler to produce a binary (NEFF) that is executable on Neuron.

The main difference is that the call to `trace()` returns a *new* fully compiled graph as a `TorchScript` Module. Upon calling this new Module, rather than re-executing the python, rebuilding the graph, and checking the cache for a matching model, the new Module simply executes the precompiled graph that was preloaded during tracing. This is a significantly more optimized runtime since it avoids the python operator tracing, graph building, etc.

One disadvantage of this interface is that a model will never dynamically recompile after a trace. This means that dynamic control flow is not supported within a function/module. Tensor input/output shapes are fixed to the shapes passed to the `trace()` API. Dynamic batching and bucketing can be used to avoid the pitfalls of static shapes.

## Example

```

import torch
import torch_neuronx

# Create example model and inputs
model = torch.nn.Sequential(
    torch.nn.Linear(784, 120),
    torch.nn.ReLU(),
    torch.nn.Linear(120, 10),
    torch.nn.Softmax(dim=-1),
)
model.eval()
example = torch.rand((1, 784))

# Create fixed model trace
trace = torch_neuronx.trace(model, example)

# Inference
result = trace(example) # No recompilation. Input shapes must not change
print(result)

```

## Traced Inference Advantages

Traced inference should be used for nearly all deployment purposes since it provides some key advantages over XLA *Lazy Tensor* execution:

- **Reduced Overhead:** There is no overhead associated with graph recording, compilation, and model loading since these steps are performed only once within the call to `trace()`. In contrast, when using XLA *Lazy Tensor* inference, all of these steps are performed just-in-time (with caching to improve performance).
- **Serializable:** The TorchScript Module that is produced from the `trace()` API is serializable using the normal `torch.jit.save()` function. It is able to be reloaded in an inference environment with `torch.jit.load()`. In contrast, XLA device inference does not provide a predetermined serialization format that includes the pre-compiled NEFF artifacts. These must be manually copied to an inference environment to be used.
- **Reduced Dependencies:** When using the traced TorchScript Module in an inference environment, it is no longer required to install the *neuronx-cc* compiler. In contrast, when using the XLA *Lazy Tensor* execution, an execution may require a recompile to successfully infer.
- **Static & Predictable:** The resulting module produced by `trace()` will contain a static model that will consume a predictable amount of Neuron device memory and will never require recompilation based on input changes. In contrast, since XLA device inference performs just-in-time compilation, it can be more difficult to predict memory utilization and the compilations that may be required at inference time.
- **C++ Usability:** If the end application is an inference platform using `libtorch`, it is easy to integrate with `libtorchneuron` to load traced modules. It is not currently possible to set up an environment to use `torch` in C++ in conjunction with Neuron XLA *Lazy Tensor* execution.

## Summary

|               | XLA Device Inference | Traced Inference |
|---------------|----------------------|------------------|
| Compilation   | JIT                  | AOT              |
| Serialization | N/A                  | TorchScript      |
| Performance   | Slower               | Faster           |
| Dynamic       | Yes                  | No               |
| C++ Usage     | No                   | Yes              |

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## Data Parallel Inference on torch\_neuronx

### Table of Contents

- [Introduction](#)
- [Data parallel inference](#)
- [torch\\_neuronx.DataParallel](#)
  - [NeuronCore selection](#)
  - [Batch dim](#)

- *Dynamic batching*
- *Performance optimizations*
- *Examples*
  - *Default usage*
  - *Specifying NeuronCores*
  - *DataParallel with dim != 0*
  - *Dynamic batching*
  - *Dynamic batching disabled*

## Introduction

This guide introduces `torch_neuronx.DataParallel()`, a Python API that implements data parallelism on `ScriptModule` models created by the *PyTorch NeuronX Tracing API for Inference*. The following sections explain how data parallelism can improve the performance of inference workloads on Inferentia, including how `torch_neuronx.DataParallel()` uses dynamic batching to run inference on variable input sizes. It covers an overview of the `torch_neuronx.DataParallel()` module and provides a few *example data parallel applications*.

## Data parallel inference

Data Parallelism is a form of parallelization across multiple devices or cores, referred to as nodes. Each node contains the same model and parameters, but data is distributed across the different nodes. By distributing the data across multiple nodes, data parallelism reduces the total execution time of large batch size inputs compared to sequential execution. Data parallelism works best for smaller models in latency sensitive applications that have large batch size requirements.

### `torch_neuronx.DataParallel`

To fully leverage the Inferentia hardware, we want to use all available NeuronCores. An inf2.xlarge and inf2.8xlarge have two NeuronCores, an inf2.24xlarge has 12 NeuronCores, and an inf2.48xlarge has 24 NeuronCores. For maximum performance on Inferentia hardware, we can use `torch_neuronx.DataParallel()` to utilize all available NeuronCores.

`torch_neuronx.DataParallel()` implements data parallelism at the module level by replicating the Neuron model on all available NeuronCores and distributing data across the different cores for parallelized inference. This function is analogous to `DataParallel` in PyTorch. `torch_neuronx.DataParallel()` requires PyTorch  $\geq 1.8$ .

The following sections provide an overview of some of the features of `torch_neuronx.DataParallel()` that enable maximum performance on Inferentia.

## NeuronCore selection

By default, DataParallel will try to use all NeuronCores allocated to the current process to fully saturate the Inferentia hardware for maximum performance. It is more efficient to make the batch dimension divisible by the number of NeuronCores. This will ensure that NeuronCores are not left idle during parallel inference and the Inferentia hardware is fully utilized.

In some applications, it is advantageous to use a subset of the available NeuronCores for DataParallel inference. DataParallel has a `device_ids` argument that accepts a list of `int` or `'nc:#'` that specify the NeuronCores to use for parallelization. See [Specifying NeuronCores](#) for an example of how to use `device_ids` argument.

## Batch dim

DataParallel accepts a `dim` argument that denotes the batch dimension used to split the input data for distributed inference. By default, DataParallel splits the inputs on `dim = 0` if the `dim` argument is not specified. For applications with a non-zero batch dim, the `dim` argument can be used to specify the inference-time input batch dimension. [DataParallel with `dim != 0`](#) provides an example of data parallel inference on inputs with batch dim = 2.

## Dynamic batching

Batch size has a direct impact on model performance. The Inferentia chip is optimized to run with small batch sizes. This means that a Neuron compiled model can outperform a GPU model, even if running single digit batch sizes.

As a general best practice, we recommend optimizing your model's throughput by compiling the model with a small batch size and gradually increasing it to find the peak throughput on Inferentia.

Dynamic batching is a feature that allows you to use tensor batch sizes that the Neuron model was not originally compiled against. This is necessary because the underlying Inferentia hardware will always execute inferences with the batch size used during compilation. Fixed batch size execution allows tuning the input batch size for optimal performance. For example, batch size 1 may be best suited for an ultra-low latency on-demand inference application, while batch size > 1 can be used to maximize throughput for offline inferencing. Dynamic batching is implemented by slicing large input tensors into chunks that match the batch size used during the `torch_neuronx.trace()` compilation call.

The `torch_neuronx.DataParallel()` class automatically enables dynamic batching on eligible models. This allows us to run inference in applications that have inputs with a variable batch size without needing to recompile the model. See [Dynamic batching](#) for an example of how DataParallel can be used to run inference on inputs with a dynamic batch size without needing to recompile the model.

Dynamic batching using small batch sizes can result in sub-optimal throughput because it involves slicing tensors into chunks and iteratively sending data to the hardware. Using a larger batch size at compilation time can use the Inferentia hardware more efficiently in order to maximize throughput. You can test the tradeoff between individual request latency and total throughput by fine-tuning the input batch size.

Automatic batching in the DataParallel module can be disabled using the `disable_dynamic_batching()` function as follows:

```
>>> model_parallel = torch_neuronx.DataParallel(model_neuron)
>>> model_parallel.disable_dynamic_batching()
```

If dynamic batching is disabled, the compile-time batch size must be equal to the inference-time batch size divided by the number of NeuronCores. [DataParallel with `dim != 0`](#) and [Dynamic batching disabled](#) provide examples of running DataParallel inference with dynamic batching disabled.



## Performance optimizations

The DataParallel module has a `num_workers` attribute that can be used to specify the number of worker threads used for multithreaded inference. By default, `num_workers = 2 * number of NeuronCores`. This value can be fine tuned to optimize DataParallel performance.

DataParallel has a `split_size` attribute that dictates the size of the input chunks that are distributed to each NeuronCore. By default, `split_size = max(1, input.shape[dim] // number of NeuronCores)`. This value can be modified to optimally match the inference input chunk size with the compile-time batch size.

## Examples

The following sections provide example usages of the `torch_neuronx.DataParallel()` module.

### Default usage

The default DataParallel use mode will replicate the model on all available NeuronCores in the current process. The inputs will be split on `dim=0`.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module
model_parallel = torch_neuronx.DataParallel(model_neuron)

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

### Specifying NeuronCores

The following example uses the `device_ids` argument to use the first three NeuronCores for DataParallel inference.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
```

(continues on next page)

(continued from previous page)

```

model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module, run on the first two NeuronCores
# Equivalent to model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1])
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=['nc:0', 'nc:1'])

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)

```

### DataParallel with dim != 0

In this example we run DataParallel inference using two NeuronCores and dim = 2. Because dim != 0, dynamic batching is not enabled. Consequently, the DataParallel inference-time batch size must be two times the compile-time batch size.

```

import torch
import torch_neuronx

# Create an example model
class Model(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = torch.nn.Conv2d(3, 3, 3)

    def forward(self, x):
        return self.conv(x) + 1

model = Model()
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 8, 8])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module using 2 NeuronCores and dim = 2
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=[0, 1], dim=2)

# Create a batched input
# Note that image_batched.shape[dim] / len(device_ids) == image.shape[dim]
batch_size = 2 * 8
image_batched = torch.rand([1, 3, batch_size, 8])

```

(continues on next page)

(continued from previous page)

```
# Run inference with a batched input
output = model_parallel(image_batched)
```

## Dynamic batching

In the following example, we use the `torch_neuronx.DataParallel()` module to run inference using several different batch sizes without recompiling the Neuron model.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)

# Create the DataParallel module
model_parallel = torch_neuronx.DataParallel(model_neuron)

# Create batched inputs and run inference on the same model
batch_sizes = [2, 3, 4, 5, 6]
for batch_size in batch_sizes:
    image_batched = torch.rand([batch_size, 3, 224, 224])

    # Run inference with a batched input
    output = model_parallel(image_batched)
```

## Dynamic batching disabled

In the following example, we use `torch_neuronx.DataParallel.disable_dynamic_batching()` to disable dynamic batching. We provide an example of a batch size that will not work when dynamic batching is disabled as well as an example of a batch size that does work when dynamic batching is disabled.

```
import torch
import torch_neuronx
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch_neuronx.trace(model, image)
```

(continues on next page)

(continued from previous page)

```

# Create the DataParallel module and use 2 NeuronCores
model_parallel = torch_neuronx.DataParallel(model_neuron, device_ids=[0, 1], dim=0)

# Disable dynamic batching
model_parallel.disable_dynamic_batching()

# Create a batched input (this won't work)
batch_size = 4
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will fail because dynamic batching is disabled and
# image_batched.shape[dim] / len(device_ids) != image.shape[dim]
# output = model_parallel(image_batched)

# Create a batched input (this will work)
batch_size = 2
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will work because
# image_batched.shape[dim] / len(device_ids) == image.shape[dim]
output = model_parallel(image_batched)

```

This document is relevant for: Inf2, Trn1, Trn2

## Developer Guide for Inference (torch-neuronx)

- [NeuronCore Allocation and Model Placement for Inference \(torch-neuronx\)](#)
- [Comparison of Traced Inference versus XLA Lazy Tensor Inference \(torch-neuronx\)](#)
- [Data Parallel Inference on torch\\_neuronx](#)
- [torch-neuronx-autobucketing-devguide](#)

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## Misc (torch-neuronx)

This document is relevant for: Inf2, Trn1, Trn2

## PyTorch Neuron (torch-neuronx) release notes

### Table of Contents

- [Release \[2.7.0.2.8.\\*, 2.6.0.2.8.\\*, 2.5.1.2.8.\\*\]](#)
- [Release \[2.6.0.2.7.\\*, 2.5.1.2.7.\\*\]](#)
- [Release \[2.5.1.2.6.0\]](#)

- *Release [2.5.1.2.4.0]*
- *Release [2.1.2.2.4.0]*
- *Release [1.13.1.1.16.0]*
- *Release [2.1.2.2.3.2]*
- *Release [2.1.2.2.3.1]*
- *Release [2.1.2.2.3.0]*
- *Release [1.13.1.1.16.0]*
- *Release [2.1.2.2.2.0]*
- *Release [1.13.1.1.15.0]*
- *Release [2.1.2.2.1.0]*
- *Release [1.13.1.1.14.0]*
- *Release [2.1.1.2.0.0b0] (Beta)*
- *Release [1.13.1.1.13.0]*
- *Release [2.0.0.2.0.0b0] (Beta)*
- *Release [1.13.1.1.12.0]*
- *Release [1.13.1.1.11.0]*
- *Release [1.13.1.1.10.1]*
- *Release [1.13.1.1.10.0]*
- *Release [1.13.1.1.9.0]*
- *Release [1.13.1.1.8.0]*
- *Release [1.13.1.1.7.0]*
- *Release [1.13.0.1.6.1]*
- *Release [1.13.0.1.6.1]*
- *Release [1.13.0.1.6.0]*
- *Release [1.13.0.1.5.0]*
- *Release [1.13.0.1.4.0]*
- *Release [1.12.0.1.4.0]*
- *Release [1.11.0.1.2.0]*
- *Release [1.11.0.1.1.1]*

PyTorch Neuron for Trn1/Inf2 is a software package that enables PyTorch users to train, evaluate, and perform inference on second-generation Neuron hardware (See: NeuronCore-v2).

**Release [2.7.0.2.8.\*, 2.6.0.2.8.\*, 2.5.1.2.8.\*]**

Date: 6/24/2025

**Summary**

- *Introducing PyTorch 2.7 Support*

**Known limitations**

- PyTorch NeuronX currently does not support GSPMD
- PyTorch NeuronX currently does not support torch.compile
- PyTorch NeuronX currently does not support DDP/FSDP

**Resolved issues****[v2.7] Resolved the lower BERT pretraining performance with torch-neuronx 2.6 compared to torch-neuronx 2.5**

With torch-neuronx v2.6, BERT pretraining performance is ~10% lower compared to torch-neuronx 2.5. This issue is fixed with torch-neuronx v2.7. See <https://github.com/pytorch/xla/issues/9037> for more details.

**Known issues**

Please see the *Introducing PyTorch 2.6 Support* for a full list of known issues with v2.6. Please see the *Introducing PyTorch 2.5 Support* for a full list of known issues with v2.5.

**Updating Ubuntu OS kernel version from 5.15 to 6.8 may result in lower performance**

Currently, when switching Ubuntu OS kernel version from 5.15 to 6.8, you may see performance differences due to the new kernel scheduler (CFS vs EEVDF). For example, BERT pretraining performance could be lower by up to 10%. You may try using an older OS kernel (i.e. Amazon Linux 2023) or experiment with the kernel real-time scheduler by running `sudo chrt --fifo 99` before your command (i.e. `sudo chrt --fifo 99 <script>`) to improve the performance. Note that adjusting the real-time scheduler can also result in lower performance. See <https://www.kernel.org/doc/html/latest/scheduler/sched-eevdf.html> for more information.

**Tensor split on second dimension of 2D array not working**

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`. Another work-around is to use `torch.tensor_split`.

### [v2.5] Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)

When using torch 2.5.1+cpu (CPU flavor) on python 3.10, importing torch\_xla crashed with `TypeError: must be called with a dataclass type or instance` due to installed triton version 3.2.0 (<https://github.com/pytorch/xla/issues/8560>). To work-around, please remove the installed triton package or downgrade to triton==3.1.0 or use the regular torch 2.5.1 (GPU flavor).

### [v2.5] Certain sequence of operations with `xm.save()` could corrupt tensors

When using the `xm.save` function to save tensors, please use `xm.mark_step()` before `xm.save` to avoid the error described in <https://github.com/pytorch/xla/issues/8422> where parameter aliasing could corrupt other tensor values. This issue will be fixed in a future release.

(Here `xm` is `torch_xla.core.xla_model` following PyTorch/XLA convention)

### [v2.6] Lower BERT pretraining performance with torch-neuronx 2.6 compared to torch-neuronx 2.5

Currently, BERT pretraining performance is ~10% lower with torch-neuronx 2.6 compared to torch-neuronx 2.5. This is due to a known regression in torch-xla <https://github.com/pytorch/xla/issues/9037> and can affect other models with high graph tracing overhead. To work-around this issue, please build the `r2.6_aws_neuron` branch of torch-xla as follows (see: `ref:pytorch-neuronx-install-cxx11` for C++11 ABI version):

```
# Setup build env (make sure you are in a python virtual env). Replace "apt" with "yum"
↳ on AL2023.
sudo apt install cmake
pip install yapf==0.30.0
wget https://github.com/bazelbuild/bazelisk/releases/download/v1.20.0/bazelisk-linux-
↳ amd64
sudo cp bazelisk-linux-amd64 /usr/local/bin/bazel
# Clone repos
git clone --recursive https://github.com/pytorch/pytorch --branch v2.6.0
cd pytorch/
git clone --recursive https://github.com/pytorch/xla.git --branch r2.6_aws_neuron
_GLIBCXX_USE_CXX11_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist
cd xla/
CXX_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist and can be installed instead of the torch-xla
↳ released in pypi.org
```

### Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`

Currently, BERT pretraining performance is ~11% lower when switching to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16` due to <https://github.com/pytorch/xla/issues/8545>. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 and 2.6 although there will be deprecation warnings (as noted below).

**Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.5 release, please downcast your model directly”**

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see [migration\\_from\\_xla\\_downcast\\_bf16](#))

**[v2.6] AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'xrt\_world\_size'**

This is an error that `torch_xla.core.xla_model.xrt_world_size()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.world_size()` instead.

**[v2.6] AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'get\_ordinal'**

This is an error that `torch_xla.core.xla_model.xla_model.get_ordinal()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.global_ordinal()` instead.

**[v2.5] WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.**

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size()` instead.

**[v2.5] WARNING:torch\_xla.core.xla\_model.xla\_model.get\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.global\_ordinal instead.**

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal()` instead.

**AttributeError: module 'torch\_xla.runtime' has no attribute 'using\_pjrt'**

In Torch-XLA 2.5, `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime. See [commit PR](#).

**Release [2.6.0.2.7.\*, 2.5.1.2.7.\*]**

Date: 5/15/2025



## Summary

- *Introducing PyTorch 2.6 Support*
- Added support for libneuronxla 2.2.\*
- Improved rendezvous performance when a payload is specified (see Resolved Issues below)
- Return non-zero exit code when neuron\_parallel\_compile encounters compilation failure(s)
- Added `torch_neuronx.testing.assert_close`, which provides Neuron allclose in an interface similar to `torch.testing.assert_close`. Neuron allclose is a modified allclose algorithm that multiplies `rto1` by the absolute max, rather than by the absolute value. This means Neuron allclose is less strict to account for our hardware. You should use `torch_neuronx.testing.assert_close` instead of `torch.testing.assert_close` to compare tensors that ran on Neuron.

## Known limitations

- PyTorch NeuronX currently does not support GSPMD
- PyTorch NeuronX currently does not support `torch.compile`
- PyTorch NeuronX currently does not support DDP/FSDP

## Resolved issues

### neuron\_parallel\_compile returns success after compilation failure(s)

Previously, when running `neuron_parallel_compile --command compile` with a Neuron Cache that contains known-bad HLO file, `neuron_parallel_compile` fails to compile the graph and still returns with exit code 0 / success. This issue is now fixed in this release.

### Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with OSError when the Neuron Cache is placed on FSx mount

Previously, the Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with `OSError (Errno 61)` when the Neuron Cache is placed on FSx mount. This issue is fixed in this release.

### Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True

Previously, using `torch.utils.data.DataLoader` with `shuffle=True` would cause the `Check failed: tensor_data` error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)). This issue is fixed in release 2.23 with the updated rendezvous handling when a payload is specified.

### **"EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile**

Previously, HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. These errors are resolved release 2.23 with the updated rendezvous handling when a payload is specified.

### **Known issues**

Please see the *Introducing PyTorch 2.6 Support* for a full list of known issues with v2.6. Please see the *Introducing PyTorch 2.5 Support* for a full list of known issues with v2.5.

### **Updating Ubuntu OS kernel version from 5.15 to 6.8 may result in lower performance**

Currently, when switching Ubuntu OS kernel version from 5.15 to 6.8, you may see performance differences due to the new kernel scheduler (CFS vs EEVDF). For example, BERT pretraining performance could be lower by up to 10%. You may try using an older OS kernel or experiment with the kernel real-time scheduler by running `sudo chrt --fifo 99` before your command (i.e. `sudo chrt --fifo 99 <script>`) to improve the performance. Note that adjusting the real-time scheduler can also result in lower performance.

### **Tensor split on second dimension of 2D array not working**

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`. Another work-around is to use `torch.tensor_split`.

### **[v2.5] Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)**

When using torch 2.5.1+cpu (CPU flavor) on python 3.10, importing torch\_xla crashed with `TypeError: must be called with a dataclass type or instance` due to installed triton version 3.2.0 (<https://github.com/pytorch/xla/issues/8560>). To work-around, please remove the installed triton package or downgrade to triton==3.1.0 or use the regular torch 2.5.1 (GPU flavor).

### **[v2.5] Certain sequence of operations with xm.save() could corrupt tensors**

When using the `xm.save` function to save tensors, please use `xm.mark_step()` before `xm.save` to avoid the error described in <https://github.com/pytorch/xla/issues/8422> where parameter aliasing could corrupt other tensor values. This issue will be fixed in a future release.

(Here `xm` is `torch_xla.core.xla_model` following PyTorch/XLA convention)

## [v2.6] Lower BERT pretraining performance with torch-neuronx 2.6 compared to torch-neuronx 2.5

Currently, BERT pretraining performance is ~10% lower with torch-neuronx 2.6 compared to torch-neuronx 2.5. This is due to a known regression in torch-xla <https://github.com/pytorch/xla/issues/9037> and can affect other models with high graph tracing overhead. To work-around this issue, please build the `r2.6_aws_neuron` branch of torch-xla as follows (see: `ref:pytorch-neuronx-install-cxx11` for C++11 ABI version):

```
# Setup build env (make sure you are in a python virtual env). Replace "apt" with "yum"
↪ on AL2023.
sudo apt install cmake
pip install yapf==0.30.0
wget https://github.com/bazelbuild/bazelisk/releases/download/v1.20.0/bazelisk-linux-
↪ amd64
sudo cp bazelisk-linux-amd64 /usr/local/bin/bazel
# Clone repos
git clone --recursive https://github.com/pytorch/pytorch --branch v2.6.0
cd pytorch/
git clone --recursive https://github.com/pytorch/xla.git --branch r2.6_aws_neuron
_GLIBCXX_USE_CXX11_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist
cd xla/
CXX_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist and can be installed instead of the torch-xla
↪ released in pypi.org
```

## Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`

Currently, BERT pretraining performance is ~11% lower when switching to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16` due to <https://github.com/pytorch/xla/issues/8545>. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 and 2.6 although there will be deprecation warnings (as noted below).

## Warning “`XLA_DOWNCAST_BF16` will be deprecated after the 2.5 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

## [v2.6] `AttributeError: <module 'torch_xla.core.xla_model' ... does not have the attribute 'xrt_world_size'`

This is an error that `torch_xla.core.xla_model.xrt_world_size()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.world_size()` instead.

**[v2.6] AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'get\_ordinal'**

This is an error that `torch_xla.core.xla_model.xla_model.get_ordinal()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.global_ordinal()` instead.

**[v2.5] WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.**

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size()` instead.

**[v2.5] WARNING:torch\_xla.core.xla\_model.xla\_model.get\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.global\_ordinal instead.**

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal()` instead.

**AttributeError: module 'torch\_xla.runtime' has no attribute 'using\_pjrt'**

In Torch-XLA 2.5, `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime. See [commit PR](#).

## Release [2.5.1.2.6.0]

Date: 4/3/2025

### Summary

Minor bug fixes and enhancements.

### Known limitations

- PyTorch NeuronX currently does not support GSPMD
- PyTorch NeuronX currently does not support `torch.compile`
- PyTorch NeuronX currently does not support DDP/FSDP

## Known issues

Please see the *Introducing PyTorch 2.5 Support* for a full list of known issues.

### Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with OSError when the Neuron Cache is placed on FSx mount

Currently, the Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with OSError (Errno 61) when the Neuron Cache is placed on FSx mount:

```
[rank197]: RuntimeError: Bad StatusOr access: INVALID_ARGUMENT: RunNeuronCCImpl: error_
→condition !(error != 400): <class 'OSError': [Errno 61] No data available: '/fsxl/
→neuron_cache/neuronxcc-2.16.372.0+4a9b2326/MODULE_3540044791706521849+4eb52b03/model.
→neff' -> '/tmp/tmpx7bvfpmm/model.neff'
```

We found that the error is due to FSx failing during file copy when there are multiple readers (13 workers fail to copy out of 256). This issue doesn't affect simpler models like BERT.

To work-around the issue, please use the shared NFS mount (/home directory on a Parallel Cluster) instead of FSx to store Neuron Cache. This will be fixed in an upcoming release.

### Running in-place update operations (e.g. all\_reduce) on 0-dimensional tensors result in buffer aliasing errors in torch 2.5 and earlier

Torch's lazy tensor core has a feature where 0-dimensional tensors are stored in a device cache, so scalar constant values can be transferred once and then reused. The values in the device cache are supposed to be marked read-only and never participate in parameter aliasing. However, due to a bug in torch-xla 2.5 (#8499), sometimes the read-only flag can be dropped, allowing these tensors to be donated, resulting in aliasing errors later when the cached value is used again.

A work-around is to avoid using 0-dimensional tensors by changing them to be 1d tensor of length 1 ([example](#)). If modifying library code is not possible, disable XLA parameter aliasing by setting environment variable `XLA_ENABLE_PARAM_ALIASING=0`

### Tensor split on second dimension of 2D array not working

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`.

### Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)

When using torch 2.5.1+cpu (CPU flavor) on python 3.10, importing torch\_xla crashed with `TypeError: must be called with a dataclass type or instance` due to installed triton version 3.2.0 (<https://github.com/pytorch/xla/issues/8560>). To work-around, please remove the installed triton package or downgrade to triton==3.1.0 or use the regular torch 2.5.1 (GPU flavor).

### Certain sequence of operations with `xm.save()` could corrupt tensors

When using the `xm.save` function to save tensors, please use `xm.mark_step()` before `xm.save` to avoid the error described in <https://github.com/pytorch/xla/issues/8422> where parameter aliasing could corrupt other tensor values. This issue will be fixed in a future release.

(Here `xm` is `torch_xla.core.xla_model` following PyTorch/XLA convention)

### Lower BERT pretraining performance with torch-neuronx 2.5 compared to torch-neuronx 2.1

Currently, BERT pretraining performance is ~11% lower with torch-neuronx 2.5 compared to torch-neuronx 2.1. This is due to the switch to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16`. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 although there will be deprecation warnings (as noted below).

### Warning “`XLA_DOWNCAST_BF16` will be deprecated after the 2.5 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

### WARNING:`torch_xla.core.xla_model.xrt_world_size()` will be removed in release 2.7. is deprecated. Use `torch_xla.runtime.world_size` instead.

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size()` instead.

### WARNING:`torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in release 2.7. is deprecated. Use `torch_xla.runtime.global_ordinal` instead.

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal` instead.

### AttributeError: module ‘`torch_xla.runtime`’ has no attribute ‘`using_pjrt`’

In Torch-XLA 2.5, `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime. See [commit PR](#).

## "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile

With PyTorch 2.5 (torch-neuronx), HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. To work-around this issue, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

## Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True

With PyTorch 2.5 (torch-neuronx), using `torch.utils.data.DataLoader` with `shuffle=True` would cause the following error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)):

```
RuntimeError: torch_xla/csrc/xla_graph_executor.cpp:562 : Check failed: tensor_data
```

This is due to `synchronize_rng_states` using `xm.mesh_reduce` to synchronize RNG states. `xm.mesh_reduce` in turn uses `xm.rendezvous()` with payload, which as noted in 2.x migration guide, would result in extra graphs that could lead to lower performance due to change in `xm.rendezvous()` in torch-xla 2.x. In the case of [ZeRO1 tutorial](#), using `xm.rendezvous()` with payload also lead to the error above. This limitation will be fixed in an upcoming release. For now, to work around the issue, please disable shuffle in `DataLoader` when `NEURON_EXTRACT_GRAPHS_ONLY` environment is set automatically by Neuron Parallel Compile:

```
train_dataloader = DataLoader(
    train_dataset, shuffle=(os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None) == None),
    ↪collate_fn=default_data_collator, batch_size=args.per_device_train_batch_size
)
```

Additionally, as in the previous section, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
```

(continues on next page)

(continued from previous page)

```
xdata = xm.all_gather(xdatain, pin_layout=False)
cpu_xdata = xdata.detach().to("cpu")
cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
xldata = [x for x in cpu_xdata_split]
return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

## Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.5 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.21, please disable gradient accumulation in torch-neuronx 2.5.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.8/concurrent/futures/
↳ process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]
```

## Release [2.5.1.2.4.0]

Date: 12/20/2024

### Summary

- *Introducing PyTorch 2.5 Support*
- Added support for Trainium2
- Added support for C++11 ABI
- Added support for Neuron Profiler 2.0
- Added support for libneuronxla 2.1.\*
- Supported Python versions: 3.9, 3.10, 3.11

### Known limitations

- PyTorch NeuronX currently does not support GSPMD
- PyTorch NeuronX currently does not support torch.compile
- PyTorch NeuronX currently does not support DDP/FSDP



## Known issues

Please see the *Introducing PyTorch 2.5 Support* for a full list of known issues.

### Tensor split on second dimension of 2D array not working

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`.

### Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)

When using torch 2.5.1+cpu (CPU flavor) on python 3.10, importing torch\_xla crashed with `TypeError: must be called with a dataclass type or instance` due to installed triton version 3.2.0 (<https://github.com/pytorch/xla/issues/8560>). To work-around, please remove the installed triton package or downgrade to triton==3.1.0 or use the regular torch 2.5.1 (GPU flavor).

### Certain sequence of operations with xm.save() could corrupt tensors

When using the `xm.save` function to save tensors, please use `xm.mark_step()` before `xm.save` to avoid the error described in <https://github.com/pytorch/xla/issues/8422> where parameter aliasing could corrupt other tensor values. This issue will be fixed in a future release.

(Here `xm` is `torch_xla.core.xla_model` following PyTorch/XLA convention)

### Lower BERT pretraining performance with torch-neuronx 2.5 compared to torch-neuronx 2.1

Currently, BERT pretraining performance is ~11% lower with torch-neuronx 2.5 compared to torch-neuronx 2.1. This is due to the switch to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16`. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 although there will be deprecation warnings (as noted below).

### Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.5 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

**WARNING:**`torch_xla.core.xla_model.xrt_world_size()` will be removed in release 2.7. is deprecated. Use `torch_xla.runtime.world_size` instead.

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size` instead.

**WARNING:**`torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in release 2.7. is deprecated. Use `torch_xla.runtime.global_ordinal` instead.

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal()` instead.

**AttributeError:** module `'torch_xla.runtime'` has no attribute `'using_pjrt'`

In Torch-XLA 2.5, `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime. See [commit PR](#).

**"EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile**

With PyTorch 2.5 (torch-neuronx), HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. To work-around this issue, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

**Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True**

With PyTorch 2.5 (torch-neuronx), using `torch.utils.data.DataLoader` with `shuffle=True` would cause the following error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)):

```
RuntimeError: torch_xla/csrc/xla_graph_executor.cpp:562 : Check failed: tensor_data
```

This is due to `synchronize_rng_states` using `xm.mesh_reduce` to synchronize RNG states. `xm.mesh_reduce` in turn uses `xm.rendezvous()` with payload, which as noted in 2.x migration guide, would result in extra graphs that could lead to lower performance due to change in `xm.rendezvous()` in torch-xla 2.x. In the case of [ZeRO1 tutorial](#),

using `xm.rendezvous()` with payload also lead to the error above. This limitation will be fixed in an upcoming release. For now, to work around the issue, please disable shuffle in DataLoader when `NEURON_EXTRACT_GRAPHS_ONLY` environment is set automatically by Neuron Parallel Compile:

```
train_dataloader = DataLoader(
    train_dataset, shuffle=(os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None) == None),
    ↪collate_fn=default_data_collator, batch_size=args.per_device_train_batch_size
)
```

Additionally, as in the previous section, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

## Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.5 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.21/2.22, please disable gradient accumulation in torch-neuronx 2.5.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.8/concurrent/futures/
↪process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]
```

## Release [2.1.2.2.4.0]

Date: 12/xx/2024

## Summary

- Added support for Trainium2
- Added support for C++11 ABI
- Added support for Neuron Profiler 2.0
- Added support for libneuronxla 2.1.\*

---

**Note:** The CVEs [CVE-2024-31583](#) and [CVE-2024-31580](#) affect PyTorch versions 2.1 and earlier. Based on Amazon’s analysis, executing models on Trainium and Inferentia is not exposed to either of these vulnerabilities. We recommend upgrading to the new version of Torch-NeuronX by following the Neuron setup instruction.

---

### Release [1.13.1.1.16.0]

Date: 12/xx/2024

#### Summary

Minor updates

---

**Note:** Torch NeuronX 1.13 currently does not support Trainium2.

---

---

**Note:** The CVEs [CVE-2024-31583](#) and [CVE-2024-31580](#) affect PyTorch versions 2.1 and earlier. Based on Amazon’s analysis, executing models on Trainium and Inferentia is not exposed to either of these vulnerabilities. We recommend upgrading to the new version of Torch-NeuronX by following the Neuron setup instruction.

---

### Release [2.1.2.2.3.2]

Date: 11/20/2024

#### Summary

This patch narrows the range of dependent libneuronxla versions to support minor version bumps and fixes the “list index out of range” error when using the Zero Redundancy Optimizer (ZeRO1) checkpoint loading.

### Release [2.1.2.2.3.1]

Date: 10/25/2024

#### Summary

This patch release removes the excessive lock wait time during neuron\_parallel\_compile graph extraction for large cluster training.

## Release [2.1.2.2.3.0]

Date: 09/16/2024

### Summary

This release adds support for Neuron Kernel Interface (NKI), Python 3.11, and protobuf versions 3.20+, as well as improved BERT performance.

### What's new in this release

- Added support for Neuron Kernel Interface (NKI). Please see [NKI documentation](#) for more information.
- Added support for Python 3.11.
- Added support for protobuf versions 3.20+.
- (Training) Increased performance for BERT-Large pretraining by changing NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS default.
- (Training) Improved Neuron Cache locking mechanism for better Neuron Cache performance during multi-node training
- (Inference) Added support for weight separated models for DataParallel class.

### Known limitations

The following features are not yet supported in this version of Torch-Neuronx 2.1: \* (Training) GSPMD \* (Training/Inference) TorchDynamo (torch.compile) \* (Training) DDP/FSDP

### Resolved Issues

#### Better performance for BERT-Large pretraining

Currently we see about 20% better trn1.32xlarge performance for BERT-Large BF16 pre-training with PyTorch 2.1 (torch-neuronx) when NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS="Embedding" (the new default) instead of the previous default "Embedding,LayerNorm,Linear,Conv2d,BatchNorm2d". No action is needed from users when using release 2.20's torch-neuronx which includes the new default. See [list of environment variables](#) regarding information about NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS.

### Known issues

Please see the [Introducing PyTorch 2.1 Support](#) for a full list of known issues.

### Error cannot import name 'builder' from 'google.protobuf.internal' after installing compiler from earlier releases (2.19 or earlier)

When using torch-neuronx from Neuron SDK release 2.20 and installing the compiler from an earlier release (Neuron SDK release 2.19 or earlier), you may encounter the error `ImportError: cannot import name 'builder' from 'google.protobuf.internal'`. This issue is caused by the compiler's dependency on protobuf version 3.19 in the Neuron SDK release 2.19 or earlier.

To work-around this issue, please install protobuf 3.20.3:

```
pip install protobuf==3.20.3
```

Ignore the pip dependency check error that may occur due to the earlier compiler's dependency on protobuf version 3.19.

### Lower accuracy when fine-tuning Roberta

In the current Neuron SDK release 2.20, we have observed lower accuracy (68% vs expected 89%) when fine-tuning the RoBERTa-large model on the MRPC dataset. This issue will be addressed in a future release.

To work around this problem, you can use the compiler from Neuron SDK release 2.19, while also installing the correct version of the protobuf library. Run the following command:

```
python3 -m pip install neuronx-cc==2.14.227.0+2d4f85be protobuf==3.20.3
```

Please note the protobuf version requirement mentioned in the previous section, as it is necessary to address the compatibility issue between the Neuron SDK 2.19 compiler and the protobuf library.

### Slower loss convergence for NxL LLaMA-3 70B pretraining using ZeRO1 tutorial

Currently, with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the [LLaMA-3 70B tutorial for neuronx-distributed](#) when using the recommended flags (`NEURON_CC_FLAGS="--distribution-strategy llm-training --model-type transformer"`). To work-around this issue, please only use `--model-type transformer` flag (`NEURON_CC_FLAGS="--model-type transformer"`).

### Glibc error on Amazon Linux 2

If using PyTorch 2.1 (torch-neuronx) on Amazon Linux 2, you will see a Glibc error below. Please switch to a newer supported OS such as Ubuntu 20, Ubuntu 22, or Amazon Linux 2023.

```
ImportError: /lib64/libc.so.6: version `GLIBC_2.27' not found (required by /tmp/debug/_  
↳XLAC.cpython-38-x86_64-linux-gnu.so)
```

## "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile

With PyTorch 2.1 (torch-neuronx), HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. To work-around this issue, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

## Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True

With PyTorch 2.1 (torch-neuronx), using `torch.utils.data.DataLoader` with `shuffle=True` would cause the following error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)):

```
RuntimeError: torch_xla/src/xla_graph_executor.cpp:562 : Check failed: tensor_data
```

This is due to `synchronize_rng_states` using `xm.mesh_reduce` to synchronize RNG states. `xm.mesh_reduce` in turn uses `xm.rendezvous()` with payload, which as noted in 2.x migration guide, would result in extra graphs that could lead to lower performance due to change in `xm.rendezvous()` in torch-xla 2.x. In the case of [ZeRO1 tutorial](#), using `xm.rendezvous()` with payload also lead to the error above. This limitation will be fixed in an upcoming release. For now, to work around the issue, please disable shuffle in `DataLoader` when `NEURON_EXTRACT_GRAPHS_ONLY` environment is set automatically by Neuron Parallel Compile:

```
train_dataloader = DataLoader(
    train_dataset, shuffle=(os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None) == None),
    ↪collate_fn=default_data_collator, batch_size=args.per_device_train_batch_size
)
```

Additionally, as in the previous section, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
```

(continues on next page)

(continued from previous page)

```

xdata = xm.all_gather(xdatain, pin_layout=False)
cpu_xdata = xdata.detach().to("cpu")
cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
xldata = [x for x in cpu_xdata_split]
return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce

```

### Compiler error when `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)`

Currently with PyTorch 2.1 (torch-neuronx), using the `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)` (as previously done in the [ZeRO1 tutorial](#)) to unload graphs during execution would cause a compilation error `Expecting value: line 1 column 1 (char 0)`. You can remove this line as it is not recommended for use. Please see the updated [ZeRO1 tutorial](#) in release 2.18.

### Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.1 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.18, please use `torch-neuronx==1.13.*` or disable gradient accumulation in torch-neuronx 2.1.

```

ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.8/concurrent/futures/
↳ process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]

```

### Release [1.13.1.1.16.0]

Date: 09/16/2024

#### Summary

This release adds support for Neuron Kernel Interface (NKI), Python 3.11, and protobuf versions 3.20+.

#### What's new in this release

- Added support for Neuron Kernel Interface (NKI). Please see [NKI documentation](#) for more information.
- Added support for Python 3.11.
- Added support for protobuf versions 3.20+.
- (Inference) Added support for weight separated models for DataParallel class.



## Known Issues and Limitations

### Error cannot import name 'builder' from 'google.protobuf.internal' after installing compiler from earlier releases (2.19 or earlier)

When using torch-neuronx from Neuron SDK release 2.20 and installing the compiler from an earlier release (Neuron SDK release 2.19 or earlier), you may encounter the error `ImportError: cannot import name 'builder' from 'google.protobuf.internal'`. This issue is caused by the compiler's dependency on protobuf version 3.19 in the Neuron SDK release 2.19 or earlier.

To work-around this issue, please install protobuf 3.20.3:

```
pip install protobuf==3.20.3
```

Ignore the pip dependency check error that may occur due to the earlier compiler's dependency on protobuf version 3.19.

### Hang while training Stable Diffusion v1.5 with PyTorch 1.13 (torch-neuronx)

In this release, training Stable Diffusion v1.5 at 512x512 resolution using PyTorch 1.13 (torch-neuronx) currently results in a hang. The fix will be available in an upcoming release. To work-around, you can install compiler from release 2.19 (noting the protobuf issue mentioned above).

```
python3 -m pip install neuronx-cc==2.14.227.0+2d4f85be protobuf==3.20.3
```

Stable Diffusion v2.1 training is unaffected.

### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX` or using jemalloc library (see <https://github.com/aws-neuron/aws-neuron-sdk/issues/728>).

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known Issues and Limitations (Inference)

### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts which exceed 4GB cannot be serialized. Serializing the torchscript artifact will trigger a segfault. This issue is resolved in torch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

### Release [2.1.2.2.0]

Date: 07/03/2024

## Summary

### What's new in this release

- Improvements in ZeRO1 to have FP32 master weights support and BF16 all-gather
- Added custom SILU enabled via NEURON\_CUSTOM\_SILU environment variable
- Neuron Parallel Compile now handle non utf-8 characters in trial-run log and reports compilation time results when enabled with NEURON\_PARALLEL\_COMPILE\_DUMP\_RESULTS
- Support for using DummyStore during PJRT process group initialization by setting TORCH\_DIST\_INIT\_BARRIER=0 and XLA\_USE\_DUMMY\_STORE=1

### Known limitations

The following features are not yet supported in this version of Torch-Neuronx 2.1: \* (Training) GSPMD \* (Training/Inference) TorchDynamo (torch.compile) \* (Training) DDP/FSDP

## Resolved Issues

### Resolved an issue with slower loss convergence for GPT-2 pretraining using ZeRO1 tutorial

Previously with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the *ZeRO1 tutorial*. This issue is now resolved. Customer can now run the tutorial with the recommended flags (NEURON\_CC\_FLAGS="--distribution-strategy llm-training --model-type transformer").

### Resolved an issue with slower loss convergence for NxD LLaMA-2 70B pretraining using ZeRO1 tutorial

Previously with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the *LLaMA-2 70B tutorial for neuronx-distributed*. This issue is now resolved. Customer can now run the tutorial with the recommended flags (NEURON\_CC\_FLAGS="--distribution-strategy llm-training --model-type transformer") and turning on functionalization (XLA\_DISABLE\_FUNCTIONALIZATION=0). Turning on functionalization results in slightly higher device memory usage and ~11% lower in performance due to a known issue with torch-xla 2.1 (<https://github.com/pytorch/xla/issues/7174>). The higher device memory usage also limits LLaMA-2 70B tutorial to run on 16

trn1.32xlarge nodes at the minimum, and running on 8 nodes would result in out-of-memory error. See the list of environment variables for more information about XLA\_DISABLE\_FUNCTIONALIZATION.

### Resolved an issue where upon a compiler error during XLA JIT execution, the framework process exits with a stack dump followed by a core dump

Previously, when there's a compiler error during XLA JIT execution, the framework process exits with a stack dump following by a core dump:

```
2024-06-10 04:31:49.733004: F ./torch_xla/csrc/runtime/debug_macros.h:20] Non-OK-status:
↳ status.status() status: INTERNAL: RunNeuronCCImpl: error condition error != 0: <class
↳ 'subprocess.CallledProcessError': Command '' died with <Signals.SIGHUP: 1>.
*** Begin stack trace ***
    tsl::CurrentStackTrace()
    std::unique_ptr<xla::PjRtLoadedExecutable, std::default_delete<xla::
↳ PjRtLoadedExecutable> > ConsumeValue<std::unique_ptr<xla::PjRtLoadedExecutable, std::
↳ default_delete<xla::PjRtLoadedExecutable> > >(absl::lts_20230125::StatusOr<std::unique_
↳ ptr<xla::PjRtLoadedExecutable, std::default_delete<xla::PjRtLoadedExecutable> > >&&)
    torch_xla::runtime::PjRtComputationClient::Compile(std::vector<torch_xla::
↳ runtime::ComputationClient::CompileInstance, std::allocator<torch_xla::runtime::
↳ ComputationClient::CompileInstance> >)
    ...
    Py_RunMain
    Py_BytesMain
    _start
*** End stack trace ***
Aborted (core dumped)
```

This is now fixed so that the above error is more succinct:

```
RuntimeError: Bad StatusOr access: INTERNAL: RunNeuronCCImpl: error condition error != 0:
↳ <class 'subprocess.CallledProcessError': Command '' died with <Signals.SIGHUP: 1>.
```

### Resolved an issue where S3 caching during distributed training can lead to S3 throttling error

When using S3 location as Neuron Cache path (specified via NEURON\_COMPILE\_CACHE\_URL or `--cache_dir` option in NEURON\_CC\_FLAGS), you may get the error An error occurred (SlowDown) when calling the PutObject operation as in:

```
2024-04-18 01:51:38.231524: F ./torch_xla/csrc/runtime/debug_macros.h:20] Non-OK-status:
↳ status.status() status: INVALID_ARGUMENT: RunNeuronCCImpl: error condition !(error !=
↳ 400): <class 'boto3.exceptions.S3UploadFailedError': Failed to upload /tmp/
↳ tmp4d8d4r2d/model.hlo to bucket/llama-compile-cache/neuronxcc-2.13.68.0+6dfecc895/
↳ MODULE_9048582265414220701+5d2d81ce/model.hlo_module.pb: An error occurred (SlowDown)
↳ when calling the PutObject operation (reached max retries: 4): Please reduce your
↳ request rate.
```

This issue is now resolved in release 2.19.

### Resolved error “ImportError: cannot import name ‘packaging’ from ‘pkg\_resources’” when using latest setuptools version 70

As reported in <https://github.com/aws-neuron/aws-neuron-sdk/issues/893>, When running examples in environment where the latest setuptools version 70 is installed, you may get the following error:

```
ImportError: cannot import name 'packaging' from 'pkg_resources' (/home/ubuntu/aws_
↳neuron_venv_pytorch/lib/python3.8/site-packages/pkg_resources/__init__.py)
```

In release 2.19 torch-neuronx now depends on setuptools version <= 69.5.1.

### Resolved compiler assertion error when training using Hugging Face deepmind/language-perceiver model

The follow assertion error when training with Hugging Face deepmind/language-perceiver model is now resolved in release 2.19 compiler:

```
ERROR 176659 [NeuronAssert]: Assertion failure in usr/lib/python3.8/multiprocessing/
↳process.py at line 108 with exception:
Unsupported batch-norm-training op: tensor_op_name: _batch-norm-training.852 | hlo_id:↳
↳852| file_name: | Line: 0 | Column: 0 | .
```

### Resolved lower accuracy for BERT-base finetuning using HF Trainer API

With release 2.19 compiler, the MRPC dataset accuracy for BERT-base finetuning after 5 epochs is now 87% as expected.

### Resolved the issue with increased in Neuron Parallel Compile time

PyTorch 2.1 (torch-neuronx), the time to run Neuron Parallel Compile for some model configuration has decreased.

### Known issues

Please see the [Introducing PyTorch 2.1 Support](#) for a full list of known issues.

### Slower loss convergence for NxD LLaMA-3 70B pretraining using ZeRO1 tutorial

Currently, with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the [LLaMA-3 70B tutorial for neuronx-distributed](#) when using the recommended flags (`NEURON_CC_FLAGS="--distribution-strategy llm-training --model-type transformer"`). To work-around this issue, please only use `--model-type transformer` flag (`NEURON_CC_FLAGS="--model-type transformer"`).

## Gradient accumulation is not yet supported for Stable Diffusion due to a compiler error

Currently, with PyTorch 2.1 (torch-neuronx), we are seeing a compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. To train Stable Diffusion with gradient accumulation, please use PyTorch 1.13 (torch-neuronx) instead of PyTorch 2.1 (torch-neuronx).

## Enabling functionalization (XLA\_DISABLE\_FUNCTIONALIZATION=0) results in 15% lower performance and non-convergence for the BERT pretraining tutorial

Currently, with PyTorch 2.1 (torch-neuronx), enabling functionalization (XLA\_DISABLE\_FUNCTIONALIZATION=0) would result in 15% lower performance and non-convergence for the BERT pretraining tutorial. The lower performance is due to missing aliasing for gradient accumulation and is a known issue with torch-xla 2.1 (<https://github.com/pytorch/xla/issues/7174>). The non-convergence is due to an issue in marking weights as static (buffer address not changing), which can be worked around by setting NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS to empty string (NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS=""). See the list of environment variables for more information about XLA\_DISABLE\_FUNCTIONALIZATION, and NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS.

```
export NEURON_TRANSFER_WITH_STATIC_RING_OPS=""
```

## Glibc error on Amazon Linux 2

If using PyTorch 2.1 (torch-neuronx) on Amazon Linux 2, you will see a Glibc error below. Please switch to a newer supported OS such as Ubuntu 20, Ubuntu 22, or Amazon Linux 2023.

```
ImportError: /lib64/libc.so.6: version `GLIBC_2.27' not found (required by /tmp/debug/_
↳XLAC.cpython-38-x86_64-linux-gnu.so)
```

## "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile

With PyTorch 2.1 (torch-neuronx), HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. To work-around this issue, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

**Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True**

With PyTorch 2.1 (torch-neuronx), using `torch.utils.data.DataLoader` with `shuffle=True` would cause the following error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)):

```
RuntimeError: torch_xla/csrc/xla_graph_executor.cpp:562 : Check failed: tensor_data
```

This is due to `synchronize_rng_states` using `xm.mesh_reduce` to synchronize RNG states. `xm.mesh_reduce` in turn uses `xm.rendezvous()` with payload, which results in extra graphs that could lead to lower performance due to change in `xm.rendezvous()` in torch-xla 2.x. In the case of [ZeRO1 tutorial](#), using `xm.rendezvous()` with payload also lead to the error above. This limitation will be fixed in an upcoming release. For now, to work around the issue, please disable shuffle in `DataLoader` when `NEURON_EXTRACT_GRAPHS_ONLY` environment is set automatically by Neuron Parallel Compile:

```
train_dataloader = DataLoader(
    train_dataset, shuffle=(os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None) == None),
    collate_fn=default_data_collator, batch_size=args.per_device_train_batch_size
)
```

Additionally, as in the previous section, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

**Compiler error when `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)`**

Currently with PyTorch 2.1 (torch-neuronx), using the `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)` (as previously done in the [ZeRO1 tutorial](#)) to unload graphs during execution would cause a compilation error `Expecting value: line 1 column 1 (char 0)`. You can remove this line as it is not recommended for use. Please see the updated [ZeRO1 tutorial](#) in release 2.18.

## Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.1 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.18, please use `torch-neuronx==1.13.*` or disable gradient accumulation in torch-neuronx 2.1.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.8/concurrent/futures/
↳ process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]
```

## Lower performance for BERT-Large

Currently we see 8% less performance when running the BERT-Large pre-training tutorial with PyTorch 2.1 (torch-neuronx) as compared to PyTorch 1.13 (torch-neuronx).

## Release [1.13.1.15.0]

Date: 07/03/2024

### Summary

#### What's new in this release

Improvements in ZeRO1 to have FP32 master weights support and BF16 all-gather Added custom SILU enabled via NEURON\_CUSTOM\_SILU environment variable Neuron Parallel Compile now handle non utf-8 characters in trial-run log and reports compilation time results when enabled with NEURON\_PARALLEL\_COMPILE\_DUMP\_RESULTS

#### Resolved Issues

#### Known Issues and Limitations

#### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX` or using jemalloc library (see <https://github.com/aws-neuron/aws-neuron-sdk/issues/728>).

#### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

## Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known Issues and Limitations (Inference)

### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts which exceed 4GB cannot be serialized. Serializing the torchscript artifact will trigger a segfault. This issue is resolved in torch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [2.1.2.2.1.0]

Date: 04/01/2024

## Summary

This release of 2.1 includes support for Neuron Profiler, multi-instance distributed training, Nemo Megatron, and HuggingFace Trainer API.

## What's new in this release

In addition to previously supported features (Transformers-NeuronX, Torch-NeuronX Trace API, Torch-NeuronX training, NeuronX Distributed training), PyTorch 2.1 (torch-neuronx) now includes support for:

- (Inference) NeuronX Distributed inference
- (Training/Inference) Neuron Profiler
- (Training) Multi-instance distributed training
- (Training) Nemo Megatron
- (Training) *analyze* feature in *neuron\_parallel\_compile*
- (Training) HuggingFace Trainer API

Additionally, auto-bucketing is a new feature for torch-neuronx and Neuronx-Distributed allowing users to define bucket models that can be serialized into a single model for multi-shape inference.



## Known limitations

The following features are not yet supported in this version of PyTorch 2.1 (torch-neuronx):

- (Training) GSPMD
- (Training) TorchDynamo (torch.compile)
- (Training) DDP/FSDP
- (Training) S3 caching during distributed training can lead to throttling issues

## Resolved issues

### “Attempted to access the data pointer on an invalid python storage”

When using Hugging Face Trainer API with transformers version  $\geq 4.35$  and  $< 4.37.3$ , user would see the error "Attempted to access the data pointer on an invalid python storage" during model checkpoint saving. This issue is fixed in transformers version  $\geq 4.37.3$ . See <https://github.com/huggingface/transformers/issues/27578> for more information.

### Too many graph compilations when using HF Trainer API

When using Hugging Face transformers version  $\geq 4.35$  and  $< 4.37.3$ , user would see many graph compilations (see <https://github.com/aws-neuron/aws-neuron-sdk/issues/813> for more information). To work around this issue, in transformers version  $\geq 4.37.3$ , user can add the option `--save_safetensors False` to Trainer API function call and modify the installed `trainer.py` as follows (don't move model to CPU before saving checkpoint):

```
# Workaround https://github.com/aws-neuron/aws-neuron-sdk/issues/813
sed -i "s/model\.to(\"cpu\")//" `python -c "import site; print(site.getsitepackages()[0])"
↪`/trainer.py
```

### Divergence (non-convergence) of loss for BERT/LLaMA when using release 2.16 compiler

With release 2.18, the divergence (non-convergence) of BERT/LLaMA loss is resolved. No compiler flag change is required.

## Known Issues

Please see the [Introducing PyTorch 2.1 Support](#) for a full list of known issues.

## Glibc error on Amazon Linux 2

If using PyTorch 2.1 (torch-neuronx) on Amazon Linux 2, you will see a Glibc error below. Please switch to a newer supported OS such as Ubuntu 20, Ubuntu 22, or Amazon Linux 2023.

```
ImportError: /lib64/libc.so.6: version `GLIBC_2.27' not found (required by /tmp/debug/_
↳ XLAC.cpython-38-x86_64-linux-gnu.so)
```

## "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile

With PyTorch 2.1 (torch-neuronx), HF Trainer API's use of XLA function `.mesh_reduce` causes "EOFError: Ran out of input" or "\_pickle.UnpicklingError: invalid load key, '!'" errors during Neuron Parallel Compile. This is an issue with the trial execution of empty NEFFs and should not affect the normal execution of the training script.

## Check failed: tensor\_data error during when using torch.utils.data.DataLoader with shuffle=True

With PyTorch 2.1 (torch-neuronx), using `torch.utils.data.DataLoader` with `shuffle=True` would cause the following error in `synchronize_rng_states` (i.e. [ZeRO1 tutorial](#)):

```
RuntimeError: torch_xla/csrc/xla_graph_executor.cpp:562 : Check failed: tensor_data
```

This is due to `synchronize_rng_states` using `xm.mesh_reduce` to synchronize RNG states. `xm.mesh_reduce` in turn uses `xm.rendezvous()` with payload, which as noted in 2.x migration guide, would result in extra graphs that could lead to lower performance due to change in `xm.rendezvous()` in torch-xla 2.x. In the case of [ZeRO1 tutorial](#), using `xm.rendezvous()` with payload also lead to the error above. This limitation will be fixed in an upcoming release. For now, to work around the issue, please disable shuffle in `DataLoader` when `NEURON_EXTRACT_GRAPHS_ONLY` environment is set automatically by Neuron Parallel Compile:

```
train_dataloader = DataLoader(
    train_dataset, shuffle=(os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None) == None),
    ↳ collate_fn=default_data_collator, batch_size=args.per_device_train_batch_size
)
```

Additionally, you can add the following code snippet (after python imports) to replace `xm.mesh_reduce` with a form that uses `xm.all_gather` instead of `xm.rendezvous()` with payload. This will add additional small on-device graphs (as opposed to the original `xm.mesh_reduce` which runs on CPU).

```
import copy
import torch_xla.core.xla_model as xm
def mesh_reduce(tag, data, reduce_fn):
    xm.rendezvous(tag)
    xdatain = copy.deepcopy(data)
    xdatain = xdatain.to("xla")
    xdata = xm.all_gather(xdatain, pin_layout=False)
    cpu_xdata = xdata.detach().to("cpu")
    cpu_xdata_split = torch.split(cpu_xdata, xdatain.shape[0])
    xldata = [x for x in cpu_xdata_split]
    return reduce_fn(xldata)
xm.mesh_reduce = mesh_reduce
```

### Compiler error when `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)`

Currently with PyTorch 2.1 (torch-neuronx), using the `torch_neuronx.xla_impl.ops.set_unload_prior_neuron_models_mode(True)` (as previously done in the [ZeRO1 tutorial](#)) to unload graphs during execution would cause a compilation error `Expecting value: line 1 column 1 (char 0)`. You can remove this line as it is not recommended for use. Please see the updated [ZeRO1 tutorial](#) in release 2.18.

### Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.1 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.18, please use `torch-neuronx==1.13.*`.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.8/concurrent/futures/
↳ process.py at line 239 with exception:
too many partition dims! {{0,+,960}}[10],+,10560}}[10]
```

### Compiler assertion error when training using Hugging Face `deepmind/language-perceiver` model

Currently, with PyTorch 2.1 (torch-neuronx), we are seeing the following compiler assertion error when training with Hugging Face `deepmind/language-perceiver` model. This will be fixed in an upcoming release. For now, if you would like to train Hugging Face `deepmind/language-perceiver` model with Neuron SDK release 2.18, please use `torch-neuronx==1.13.*`.

```
ERROR 176659 [NeuronAssert]: Assertion failure in usr/lib/python3.8/multiprocessing/
↳ process.py at line 108 with exception:
Unsupported batch-norm-training op: tensor_op_name: _batch-norm-training.852 | hlo_id:↳
↳ 852 | file_name: | Line: 0 | Column: 0 | .
```

### Lower performance for BERT-Large

Currently we see 8% less performance when running the BERT-Large pre-training tutorial with PyTorch 2.1 (torch-neuronx) as compared to PyTorch 1.13 (torch-neuronx).

### Slower loss convergence for GPT-2 pretraining using ZeRO1 tutorial when using recommended compiler flags

Currently with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the [ZeRO1 tutorial](#) when using recommended compiler flags. To work-around this issue and restore faster convergence, please replace the `NEURON_CC_FLAGS` as below:

```
# export NEURON_CC_FLAGS="--retry_failed_compilation --distribution-strategy llm-
↳ training --model-type transformer"
export NEURON_CC_FLAGS="--retry_failed_compilation -O1"
```

## Slower loss convergence for NxD LLaMA 70B pretraining using ZeRO1 tutorial when using recommended compiler flags

Currently with PyTorch 2.1 (torch-neuronx), we see slower loss convergence in the [LLaMA-2 70B tutorial for neuronx-distributed](#) when using recommended compiler flags. To work-around this issue and restore faster convergence, please replace the NEURON\_CC\_FLAGS as below:

```
# export NEURON_CC_FLAGS="--retry_failed_compilation --distribution-strategy llm-  
↪training --model-type transformer"  
export NEURON_CC_FLAGS="--retry_failed_compilation"
```

## Lower accuracy for BERT-base finetuning using HF Trainer API

Currently, with PyTorch 2.1 (torch-neuronx), MRPC dataset accuracy for BERT-base finetuning after 5 epochs is 83% instead of 87%. A work-around is to remove the option `--model-type=transformer` from NEURON\_CC\_FLAGS. This will be fixed in an upcoming release.

## Increased in Neuron Parallel Compile time

Currently, with PyTorch 2.1 (torch-neuronx), the time to run Neuron Parallel Compile for some model configuration is increased. In one example, the Neuron Parallel Compile time for NeuronX Nemo-Megatron LLaMA 13B is 2x compared to when using PyTorch 1.13 (torch-neuronx). This will be fixed in an upcoming release.

## Release [1.13.1.1.14.0]

Date: 04/01/2024

### Summary

Auto-bucketing is a new feature for torch-neuronx and Neuronx-Distributed allowing users to define bucket models that can be serialized into a single model for multi-shape inference.

### Resolved issues

- (Inference) Fixed an issue where transformers-neuronx inference errors could crash the application and cause it to hang. Inference errors should now correctly throw a runtime exception.
- (Inference/Training) Fixed an issue where `torch.argmax()` produced incorrect results.
- (Training) `neuron_parallel_compile` tool now use `traceback.print_exc` instead of `format` to support Python 3.10.
- (Training) Fixed an issue in ZeRO1 when sharded params are initialized with `torch.double`.

## Known issues and limitations

### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX` or using jemalloc library (see <https://github.com/aws-neuron/aws-neuron-sdk/issues/728>).

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [2.1.1.2.0.0b0] (Beta)

Date: 12/21/2023

### Summary

Introducing the beta release of Torch-NeuronX with PyTorch 2.1 support.

### What's new in this release

This version of PyTorch 2.1 (torch-neuronx) supports:

- (Inference) Transformers-NeuronX
- (Inference) Torch-NeuronX Trace API
- (Training) NeuronX Distributed training
- (Training) Torch-NeuronX training
- (Training) New snapshotting capability enabled via the XLA\_FLAGS environment variable (see [debug guide](#))

### Known limitations

The following features are not yet supported in this version of PyTorch 2.1 (torch-neuronx):

- (Training/Inference) Neuron Profiler
- (Inference) NeuronX Distributed inference
- (Training) Nemo Megatron
- (Training) GSPMD
- (Training) TorchDynamo (torch.compile)
- (Training) *analyze* feature in *neuron\_parallel\_compile*
- (Training) HuggingFace Trainer API (see *Known Issues* below)

Additional limitations are noted in the *Known Issues* section below.

### Known Issues

Please see the [Introducing PyTorch 2.1 Support \(Beta\)](#) for a full list of known issues.

### Lower performance for BERT-Large

Currently we see 8% less performance when running the BERT-Large pre-training tutorial with PyTorch 2.1 (torch-neuronx) as compared to PyTorch 1.13 (torch-neuronx).

### Divergence (non-convergence) of loss for BERT/LLaMA when using release 2.16 compiler

Currently, when using release 2.16 compiler version 2.12.54.0+f631c2365, you may see divergence (non-convergence) of loss curve. To workaround this issue, please use release 2.15 compiler version 2.11.0.35+4f5279863.

## Error “Attempted to access the data pointer on an invalid python storage” when using HF Trainer API

Currently, if using HuggingFace Transformers Trainer API to train (i.e. [HuggingFace Trainer API fine-tuning tutorial](#)), you may see the error “Attempted to access the data pointer on an invalid python storage”. This is a known issue <https://github.com/huggingface/transformers/issues/27578> and will be fixed in a future release.

### Release [1.13.1.1.13.0]

Date: 12/21/2023

### Summary

#### What’s new in this release

- Added Weight Replacement API For Inference)

#### Resolved issues

- Add bucketting logic to control the size of tensors for all-gather and reduce-scatter
- Fixed ZeRO-1 bug for inferring local ranks in 2-D configuration (<https://github.com/pytorch/xla/pull/5936>)

#### Known issues and limitations

##### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX` or using jemalloc library (see <https://github.com/aws-neuron/aws-neuron-sdk/issues/728>).

##### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don’t use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

##### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [2.0.0.2.0.0b0] (Beta)

Date: 10/26/2023

## Summary

Introducing the beta release of Torch-NeuronX with PyTorch 2.0 and PJRT support.

## What's new in this release

- Updating from XRT to PJRT runtime. For more info see: <link to intro pjrt doc>
- (Inference) Added the ability to partition unsupported ops to CPU during traced inference (See `torch_neuronx.trace` API guide)

## Known issues and limitations

- Snapshotting is not supported
- `NEURON_FRAMEWORK_DEBUG=1` is not supported
- `Analyze` in `neuron_parallel_compile` is not supported
- Neuron Profiler is not supported
- VGG11 with input sizes 300x300 may show accuracy issues
- Possible issues with NeMo Megatron checkpointing
- S3 caching with `neuron_parallel_compile` may show compilation errors
- Compiling without `neuron_parallel_compile` on multiple nodes may show compilation errors
- GPT2 inference may show errors with `torch_neuronx.trace`



**Release [1.13.1.1.12.0]**

Date: 10/26/2023

**Summary****What's new in this release**

- (Training) Added coalescing of all-gather and reduce-scatter inside ZeRO1, which should help in improving performance at high cluster sizes.
- (Inference) Added the ability to partition unsupported ops to CPU during traced inference. (See `torch_neuronx.trace` API guide)
- (Inference) Previously undocumented arguments `trace API args state` and `options` are now unsupported (have no effect) and will result in a deprecation warning if used.

**Resolved issues**

- Fixed an issue where `torch.topk` would fail on specific dimensions
- (Inference) Fixed an issue where NaNs could be produced when using `torch_neuronx.dynamic_batch`
- (Inference) Updated `torch_neuronx.dynamic_batch` to better support Modules (traced, scripted, and normal modules) with multiple Neuron subgraphs
- (Inference) Isolate frontend calls to the Neuron compiler to working directories, so concurrent compilations do not conflict by being run from the same directory.

**Known issues and limitations (Training)****Memory leaking in glibc**

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX`.

**DDP shows slow convergence**

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

**Runtime crash when we use too many workers per node with DDP**

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

### Known issues and limitations (Inference)

#### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

#### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

### Release [1.13.1.1.11.0]

Date: 9/15/2023

#### Summary

#### Resolved issues

- Fixed an issue in `torch_neuronx.analyze()` which could cause failures with scalar inputs.
- Improved performance of `torch_neuronx.analyze()`.

### Release [1.13.1.1.10.1]

Date: 9/01/2023

#### Summary

Minor bug fixes and enhancements.

### Release [1.13.1.1.10.0]

Date: 8/28/2023

#### Summary

#### What's new in this release

- Removed support for Python 3.7
- (Training) Added a `neuron_parallel_compile` command to clear file locks left behind when a `neuron_parallel_compile` execution was interrupted (`neuron_parallel_compile --command clear-locks`)
- (Training) Seedable dropout now enabled by default

## Resolved issues

- (Training) Convolution is now supported
- Fixed segmentation fault when using torch-neuronx to compile models on U22
- Fixed XLA tensor stride information in torch-xla package, which blocked lowering of log\_softmax and similar functions and showed errors like:

```
File "/home/ubuntu/waldronn/asr/test_env/lib/python3.7/site-packages/torch/nn/functional.
↪py", line 1930, in log_softmax
    ret = input.log_softmax(dim)
RuntimeError: dimensionality of sizes (3) must match dimensionality of strides (1)
```

## Known issues and limitations (Training)

### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX`.

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

## No automatic partitioning

Currently, when Neuron encounters an operation that it does not support during `torch_neuronx.trace()`, it may exit with the following compiler error: “Import of the HLO graph into the Neuron Compiler has failed. This may be caused by unsupported operators or an internal compiler error.” The intended behavior when tracing is to automatically partition the model into separate subgraphs that run on NeuronCores and subgraphs that run on CPU. This will be supported in a future release. See *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

## Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [1.13.1.1.9.0]

Date: 7/19/2023

### Summary

#### What’s new in this release

Training support:

- Uses jemalloc as the primary malloc lib to avoid memory leak at checkpointing
- Added support for ZeRO-1 along with *tutorial*

Inference support:

- Add async load and lazy model load options to accelerate model loading
- Optimize DataParallel API to load onto multiple cores simultaneously when device IDs specified in `device_ids` are consecutive

#### Resolved issues (Training)

- Remove extra graph creation in `torch_neuronx.optim.adamw` when the beta/lr parameters values become 0 or 1.
- Stability improvements and faster failure on hitting a fault in XRT server used by XLA.

## Known issues and limitations (Training)

### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX`.

### Convolution is not supported

Convolution is not supported during training.

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

### No automatic partitioning

Currently, when Neuron encounters an operation that it does not support during `torch_neuronx.trace()`, it may exit with the following compiler error: "Import of the HLO graph into the Neuron Compiler has failed. This may be caused by unsupported operators or an internal compiler error." The intended behavior when tracing is to automatically partition the model into separate subgraphs that run on NeuronCores and subgraphs that run on CPU. This will be supported in a future release. See *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

### Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

### Release [1.13.1.1.8.0]

Date: 6/14/2023

### Summary

- Added s3 caching to NeuronCache.
- Added extract/compile/analyze phases to neuron\_parallel\_compile.

### What's new in this release

Training support:

- Added S3 caching support to NeuronCache. Removed NeuronCache options `--cache_size/cache_ttl` (please delete cache directories as needed).
- Added separate extract and compile phases Neuron Parallel Compile.
- Added model analyze API to Neuron Parallel Compile.

### Known issues and limitations (Training)

#### Memory leaking in glibc

glibc malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX`.

#### Convolution is not supported

Convolution is not supported during training.

#### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

## Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

### No automatic partitioning

Currently, when Neuron encounters an operation that it does not support during `torch_neuronx.trace()`, this will cause an error. The intended behavior when tracing is to automatically partition the model into separate subgraphs that run on NeuronCores and subgraphs that run on CPU. See *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

## Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [1.13.1.1.7.0]

Date: 05/01/2023

## Summary

### What's new in this release

Training support:

- Added an improved Neuron-optimized AdamW optimizer implementation.
- Added an improved Neuron-optimized `torch.nn.Dropout` implementation.
- Added an assertion when the `torch.nn.Dropout` argument `inplace=True` during training. This is currently not supported on Neuron.
- Added XLA lowering for `aten::count_nonzero`

Inference support:

- Added profiling support for models compiled with `torch_neuronx.trace()`
- Added `torch_neuronx.DataParallel` for models compiled with `torch_neuronx.trace()`

### Resolved issues (Training)

#### Unexpected behavior with `torch.autocast`

Fixed an issue where `torch.autocast` did not correctly autocast when using `torch.bfloat16`

#### Resolved slower BERT bf16 Phase 1 Single Node Performance

As of the Neuron 2.9.0 release, *BERT phase 1 pretraining* performance has regressed by approximately 8-9% when executed on a *single node* only (i.e. just one `trn1.32xlarge` instance). This is resolved in 2.10 release.

#### Resolved lower throughput for BERT-large training on AL2 instances

Starting in release 2.7, we see a performance drop of roughly 5-10% for BERT model training on AL2 instances. This is resolved in release 2.10.

### Resolved issues (Inference)

#### Error when using the original model after `torch_neuronx.trace`

Fixed an issue where model parameters would be moved to the Neuron 'xla' device during `torch_neuronx.trace()` and would no longer be available to execute on the original device. This made it more difficult to compare Neuron models against CPU since previously this would require manually moving parameters back to CPU.

#### Error when using the `xm.xla_device()` object followed by using `torch_neuronx.trace`

Fixed an issue where XLA device execution and `torch_neuronx.trace()` could not be performed in the same python process.

#### Error when executing `torch_neuronx.trace` with `torch.bfloat16` input/output tensors

Fixed an issue where `torch_neuronx.trace()` could not compile models which consumed or produced `torch.bfloat16` values.

### Known issues and limitations (Training)

#### Memory leaking in `glibc`

`glibc` malloc memory leaks affect Neuron and may be temporarily limited by setting `MALLOC_ARENA_MAX`.



## Convolution is not supported

Convolution is not supported during training.

## DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

## Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this:

```
bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files``.
```

Hence, it is recommended to use fewer workers per node with DDP.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` produces incorrect results.

### No automatic partitioning

Currently, when Neuron encounters an operation that it does not support during `torch_neuronx.trace()`, this will cause an error. The intended behavior when tracing is to automatically partition the model into separate subgraphs that run on NeuronCores and subgraphs that run on CPU. See *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

## Torchscript serialization error with compiled artifacts larger than 4GB

When using `torch_neuronx.trace()`, compiled artifacts that exceed 4GB cannot be serialized. Serializing the TorchScript artifact triggers a segmentation fault. This issue is resolved in PyTorch but is not yet released: <https://github.com/pytorch/pytorch/pull/99104>

## Release [1.13.0.1.6.1]

Date: 04/19/2023

### Summary

#### What's new in this release

Training support:

- No changes

Inference support:

- Enable deserialized TorchScript modules to be compiled with `torch_neuronx.trace()`

#### Release [1.13.0.1.6.1]

Date: 04/19/2023

### Summary

#### What's new in this release

Training support:

- No changes

Inference support:

- Enable deserialized TorchScript modules to be compiled with `torch_neuronx.trace()`

#### Release [1.13.0.1.6.0]

Date: 03/28/2023

### Summary

#### What's new in this release

Training support:

- Added pipeline parallelism support in AWS Samples for Megatron-LM

Inference support:

- Added model analysis API: `torch_neuronx.analyze`
- Added HLO opcode support for:
  - `kAtan2`
  - `kAfterAll`
  - `kMap`
- Added XLA lowering support for:
  - `aten::glu`
  - `aten::scatter_reduce`

- Updated torch.nn.MSELoss to promote input data types to a compatible type

## Resolved issues (Training)

### GRPC timeout errors when running Megatron-LM GPT 6.7B tutorial on multiple instances

When running AWS Samples for Megatron-LM GPT 6.7B tutorial over multiple instances, you may encounter GRPC timeout errors like below:

```
E0302 01:10:20.511231294 138645 http2_transport.cc:1098] Received a GOAWAY with
↳ error code ENHANCE_YOUR_CALM and debug data equal to "too_many_pings"
2023-03-02 01:10:20.511500: W tensorflow/core/distributed_runtime/rpc/grpc_remote_master.
↳ cc:157] RPC failed with status = "UNAVAILABLE: Too many pings" and grpc_error_string =
↳ "{"created": "@1677719420.511317309", "description": "Error received from peer ipv4:10.1.
↳ 35.105:54729", "file": "external/com_github_grpc_grpc/src/core/lib/surface/call.cc",
↳ "file_line": 1056, "grpc_message": "Too many pings", "grpc_status": 14}", maybe retrying
↳ the RPC
```

or:

```
2023-03-08 21:18:27.040863: F tensorflow/compiler/xla/xla_client/xrt_computation_client.
↳ cc:476] Non-OK-status: session->session()->Run(session_work->feed_inputs, session_work-
↳ >outputs_handles, &outputs) status: UNKNOWN: Stream removed
```

This is due to excessive DNS lookups during execution, and is fixed in this release.

### NaNs seen with transformers version $\geq 4.21.0$ when running HF GPT fine-tuning or pretraining with XLA\_USE\_BF16=1 or XLA\_DOWNCAST\_BF16=1

Using Hugging Face transformers version  $\geq 4.21.0$  can produce NaN outputs for GPT models when using full BF16 (XLA\_USE\_BF16=1 or XLA\_DOWNCAST\_BF16=1) plus stochastic rounding. This issue occurs due to large negative constants used to implement attention masking (<https://github.com/huggingface/transformers/pull/17306>). To workaround this issue, please use transformers version  $\leq 4.20.0$ .

## Resolved issues (Inference)

### torch.argmax() now supports single argument call variant

Previously only the 3 argument variant of torch.argmax() was supported. Now the single argument call variant is supported.

## Known issues and limitations (Training)

### Slower BERT bf16 Phase 1 Single Node Performance

In the Neuron 2.9.0 release, *BERT phase 1 pretraining* performance has regressed by approximately 8-9% when executed on a *single node* only (i.e. just one `trn1.32xlarge` instance).

### Convolution is not supported

In this release, convolution is not supported.

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this: `bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files.`

Hence, it is recommended to use fewer workers per node with DDP.

### Lower throughput for BERT-large training on AL2 instances

We see a performance drop of roughly 5-10% for BERT model training on AL2 instances. This is because of the increase in time required for tracing the model.

## Known issues and limitations (Inference)

### `torch.argmax()` produces incorrect results

`torch.argmax()` now supports both the single argument call variant and the 3 argument variant. However, `torch.argmax()` currently produces incorrect results.

### Error when using the `xm.xla_device()` object followed by using `torch_neuronx.trace`

Executing a model using the `xm.xla_device()` object followed by using `torch_neuronx.trace` in the same process can produce errors in specific situations due to torch-xla caching behavior. It is recommended that only one type of execution is used per process.

## Error when executing `torch_neuronx.trace` with `torch.bfloat16` input/output tensors

Executing `torch_neuronx.trace` with `torch.bfloat16` input/output tensors can cause an error. It is currently recommended to use an alternative torch data type in combination with compiler casting flags instead.

## No automatic partitioning

Currently, there's no automatic partitioning of a model into subgraphs that run on NeuronCores and subgraphs that run on CPU Operations in the model that are not supported by Neuron would result in compilation error. Please see *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

## Release [1.13.0.1.5.0]

Date: 02/24/2023

## Summary

### What's new in this release

Training support:

- Added SPMD flag for XLA backend to generate global collective-compute replica groups

Inference support:

- Expanded inference support to `inf2`
- Added Dynamic Batching

## Resolved issues

### Known issues and limitations (Training)

#### Convolution is not supported

In this release, convolution is not supported.

#### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

## Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this: `bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files.`

Hence, it is recommended to use fewer workers per node with DDP.

## Lower throughput for BERT-large training on AL2 instances

We see a performance drop of roughly 5-10% for BERT model training on AL2 instances. This is because of the increase in time required for tracing the model.

## Known issues and limitations (Inference)

### `torch.argmax()` and `torch.argmin()` do not support the single argument call variant

`torch.argmax()` and `torch.argmin()` do not support the single argument call variant. Only the 3 argument variant of these functions is supported. The `dim` argument *must be* specified or this function will fail at the call-site. Secondly, `torch.argmin()` may produce incorrect results.

## No automatic partitioning

Currently, there's no automatic partitioning of a model into subgraphs that run on NeuronCores and subgraphs that run on CPU Operations in the model that are not supported by Neuron would result in compilation error. Please see *PyTorch Neuron (torch-neuronx) - Supported Operators* for a list of supported operators.

## Release [1.13.0.1.4.0]

Date: 02/08/2023

## Summary

### What's new in this release

Training support:

- Added support for PyTorch 1.13
- Added support for Python version 3.9
- Added support for `torch.nn.parallel.DistributedDataParallel` (DDP) along with a tutorial
- Added optimized lowering for Softmax activation
- Added support for LAMB optimizer in BF16 mode

Added initial support for inference on Trn1, including the following features:

- Trace API (`torch_neuronx.trace`)
- Core placement API (Beta)

- Python 3.7, 3.8 and 3.9 support
- Support for tracing models larger than 2 GB

The following inference features are not included in this release:

- Automatic partitioning of a model into subgraphs that run on NeuronCores and subgraphs that run on CPU
- cxx11 ABI wheels

## Resolved issues

## Known issues and limitations

### Convolution is not supported

In this release, convolution is not supported.

### DDP shows slow convergence

Currently we see that the models converge slowly with DDP when compared to the scripts that don't use DDP. We also see a throughput drop with DDP. This is a known issue with torch-xla: <https://pytorch.org/xla/release/1.13/index.html#mnist-with-real-data>

### Runtime crash when we use too many workers per node with DDP

Currently, if we use 32 workers with DDP, we see that each worker generates its own graph. This causes an error in the runtime, and you may see errors that look like this: `bootstrap.cc:86 CCOM WARN Call to accept failed : Too many open files.`

Hence, it is recommended to use fewer workers per node with DDP.

### Lower throughput for BERT-large training on AL2 instances

We see a performance drop of roughly 5-10% for BERT model training on AL2 instances. This is because of the increase in time required for tracing the model.

## Release [1.12.0.1.4.0]

Date: 12/12/2022

## Summary

### What's new in this release

- Added support for PyTorch 1.12.
- Setting `XLA_DOWNCAST_BF16=1` now also enables stochastic rounding by default (as done with `XLA_USE_BF16=1`).
- Added support for *capturing snapshots* of inputs, outputs and graph HLO for debug.
- Fixed issue with parallel compile error when both train and evaluation are enabled in HuggingFace fine-tuning tutorial.
- Added support for LAMB optimizer in FP32 mode.

### Resolved issues

#### NaNs seen with transformers version $\geq 4.21.0$ when running HF BERT fine-tuning or pretraining with `XLA_USE_BF16=1` or `XLA_DOWNCAST_BF16=1`

When running HuggingFace BERT (any size) fine-tuning tutorial or pretraining tutorial with transformers version  $\geq 4.21.0$  and using `XLA_USE_BF16=1` or `XLA_DOWNCAST_BF16=1`, you will see NaNs in the loss immediately at the first step. More details on the issue can be found at [pytorch/xla#4152](#). The workaround is to use 4.20.0 or earlier (the tutorials currently recommend version 4.15.0) or add the line `transformers.modeling_utils.get_parameter_dtype = lambda x: torch.bfloat16` to your Python training script (as now done in latest tutorials). A [permanent fix](#) will become part of an upcoming HuggingFace transformers release.

### Known issues and limitations

#### Convolution is not supported

In this release, convolution is not supported.

#### Number of data parallel training workers on one Trn1 instance

The number of workers used in single-instance data parallel training can be one of the following values: 1 or 2 for trn1.2xlarge and 1, 2, 8 or 32 for trn1.32xlarge.

#### Release [1.11.0.1.2.0]

Date: 10/27/2022



## Summary

### What's new in this release

- Added support for argmax.
- Clarified error messages for runtime errors `NRT_UNINITIALIZED` and `NRT_CLOSED`.
- When multi-worker training is launched using `torchrun` on one instance, framework now handles runtime state cleanup at end of training.

### Resolved issues

#### Drop-out rate ignored in dropout operation

A known issue in the compiler's implementation of dropout caused drop-rate to be ignored in the last release. It is fixed in the current release.

#### Runtime error “invalid offset in Coalesced\_memloc\_...” followed by “Failed to process dma block: 1703”

Previously, when running MRPC fine-tuning tutorial with `bert-base-*` model, you would encounter runtime error “invalid offset in Coalesced\_memloc\_...” followed by “Failed to process dma block: 1703”. This is fixed in the current release.

#### Compilation error: “TongaSBTensor[0x7fb2a46e0830]:TongaSB partitions[0] uint8 %138392[128, 512]”

Previously, when compiling MRPC fine-tuning tutorial with `bert-large-*` and FP32 (no `XLA_USE_BF16=1`) for two workers or more, you would encounter compiler error that looks like `Error message: TongaSBTensor[0x7fb2a46e0830]:TongaSB partitions[0] uint8 %138392[128, 512] followed by Error class: KeyError`. Single worker fine-tuning is not affected. This is fixed in the current release.

### Known issues and limitations

#### Convolution is not supported

In this release, convolution is not supported.

#### Number of data parallel training workers on one Trn1 instance

The number of workers used in single-instance data parallel training can be one of the following values: 1 or 2 for `trn1.2xlarge` and 1, 2, 8 or 32 for `trn1.32xlarge`.

### Release [1.11.0.1.1.1]

Date: 10/10/2022

### Summary

This is the initial release of PyTorch Neuron that supports Trainium for users to train their models on the new EC2 Trn1 instances.

### What's new in this release

Announcing the first PyTorch Neuron release for training.

- XLA device support for Trainium
- PyTorch 1.11 with XLA backend support in `torch.distributed`
- `torch-xla` distributed support
- Single-instance and multi-instance distributed training using `torchrun`
- Support for `ParallelCluster` and `SLURM` with node-level scheduling granularity
- Persistent cache for compiled graph
- `neuron_parallel_compile` utility to help speed up compilation
- Optimizer support: SGD, AdamW
- Loss functions supported: `NLLLoss`
- Python versions supported: 3.7, 3.8
- Multi-instance training support with EFA
- Support PyTorch's BF16 automatic mixed precision

### Known issues and limitations

#### Convolution is not supported

In this release, convolution is not supported.

#### Number of data parallel training workers on one Trn1 instance

The number of workers used in single-instance data parallel training can be one of the following values: 1 or 2 for `trn1.2xlarge` and 1, 2, 8 or 32 for `trn1.32xlarge`.

## Drop-out rate ignored in dropout operation

A known issue in the compiler's implementation of dropout caused drop-rate to be ignored. Will be fixed in a follow-on release.

## Runtime error “invalid offset in Coalesced\_memloc\_...” followed by “Failed to process dma block: 1703”

Currently, when running MRPC fine-tuning tutorial with `bert-base-*` model, you will encounter runtime error “invalid offset in Coalesced\_memloc\_...” followed by “Failed to process dma block: 1703”. This issue will be fixed in an upcoming release.

## Compilation error: “TongaSBTensor[0x7fb2a46e0830]:TongaSB partitions[0] uint8 %138392[128, 512]”

When compiling MRPC fine-tuning tutorial with `bert-large-*` and FP32 (no `XLA_USE_BF16=1`) for two workers or more, you will encounter compiler error that looks like `Error message: TongaSBTensor[0x7fb2a46e0830]: TongaSB partitions[0] uint8 %138392[128, 512]` followed by `Error class: KeyError`. Single worker fine-tuning is not affected. This issue will be fixed in an upcoming release.

*This document is relevant for:* `Inf2`, `Trn1`, `Trn2`

- [PyTorch Neuron \(torch-neuronx\) release notes](#)

*This document is relevant for:* `Inf2`, `Trn1`, `Trn2`

Setup (`torch-neuronx`)

## Tutorials (`torch-neuronx`)

- HuggingFace pretrained BERT tutorial [\[html\]](#) [\[notebook\]](#)
- TorchServe tutorial [\[html\]](#)
- LibTorch C++ tutorial (for `torch-neuron` and `torch-neuronx`) [\[html\]](#)
- Torchvision ResNet50 tutorial [\[html\]](#) [\[notebook\]](#)
- T5 inference tutorial [\[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
-

### Additional Examples (torch-neuronx)

- [AWS Neuron Samples GitHub Repository](#)
- [Transformers Neuron GitHub samples](#)

### API Reference Guide (torch-neuronx)

- [PyTorch NeuronX Tracing API for Inference](#)
- [PyTorch Neuron \(torch-neuronx\) Weight Replacement API for Inference](#)
- [PyTorch NeuronX NeuronCore Placement APIs \[Beta\]](#)
- [PyTorch NeuronX Analyze API for Inference](#)
- [PyTorch NeuronX DataParallel API](#)
- [torch\\_neuronx\\_lazy\\_async\\_load\\_api](#)

### Developer Guide (torch-neuronx)

- [NeuronCore Allocation and Model Placement for Inference \(torch-neuronx\)](#)
- [Comparison of Traced Inference versus XLA Lazy Tensor Inference \(torch-neuronx\)](#)
- [Data Parallel Inference on torch\\_neuronx](#)
- [torch-neuronx-autobucketing-devguide](#)

### Misc (torch-neuronx)

- [PyTorch Neuron \(torch-neuronx\) release notes](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1

### 2.1.3 Inference with torch-neuron (Inf1)

*This document is relevant for:* Inf1

### Tutorials for Inference with torch-neuron (Inf1)

*This document is relevant for:* Inf1

## Computer Vision Tutorials (torch-neuron)

- ResNet-50 tutorial [html] [notebook]
- PyTorch YOLOv4 tutorial [html] [notebook]

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Natural Language Processing (NLP) Tutorials (torch-neuron)

- HuggingFace pretrained BERT tutorial [html] [notebook]
- HuggingFace pretrained BERT tutorial with shared weights [html] [notebook]
- Bring your own HuggingFace pretrained BERT container to Sagemaker Tutorial [html] [notebook]
- LibTorch C++ tutorial [html]
- TorchServe tutorial [html]
- HuggingFace MarianMT tutorial [html] [notebook]

## Compiling and Deploying HuggingFace Pretrained BERT

### Introduction

In this tutorial we will compile and deploy BERT-base version of HuggingFace Transformers BERT for Inferentia. The full list of HuggingFace's pretrained BERT models can be found in the BERT section on this page [https://huggingface.co/transformers/pretrained\\_models.html](https://huggingface.co/transformers/pretrained_models.html).

This Jupyter notebook should be run on an instance which is inf1.6xlarge or larger. The compile part of this tutorial requires inf1.6xlarge and not the inference itself. For simplicity we will run this tutorial on inf1.6xlarge but in real life scenario the compilation should be done on a compute instance and the deployment on inf1 instance to save costs.

Verify that this Jupyter notebook is running the Python kernel environment that was set up according to the [PyTorch Installation Guide](#). You can select the kernel from the “Kernel -> Change Kernel” option on the top of this Jupyter notebook page.

### Install Dependencies:

This tutorial requires the following pip packages:

- torch-neuron
- neuron-cc[tensorflow]
- transformers

Most of these packages will be installed when configuring your environment using the Neuron PyTorch setup guide. The additional dependencies must be installed here.

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
↪ detect
!pip install --upgrade "transformers==4.6.0"
```

## Compile the model into an AWS Neuron optimized TorchScript

```
[ ]: import tensorflow # to workaround a protobuf version conflict issue
import torch
import torch.neuron
from transformers import AutoTokenizer, AutoModelForSequenceClassification, AutoConfig
import transformers
import os
import warnings

# Setting up NeuronCore groups for inf1.6xlarge with 16 cores
num_cores = 16 # This value should be 4 on inf1.xlarge and inf1.2xlarge
os.environ['NEURON_RT_NUM_CORES'] = str(num_cores)

# Build tokenizer and model
tokenizer = AutoTokenizer.from_pretrained("bert-base-cased-finetuned-mrpc")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-cased-finetuned-
↳mrpc", return_dict=False)

# Setup some example inputs
sequence_0 = "The company HuggingFace is based in New York City"
sequence_1 = "Apples are especially bad for your health"
sequence_2 = "HuggingFace's headquarters are situated in Manhattan"

max_length=128
paraphrase = tokenizer.encode_plus(sequence_0, sequence_2, max_length=max_length,
↳padding='max_length', truncation=True, return_tensors="pt")
not_paraphrase = tokenizer.encode_plus(sequence_0, sequence_1, max_length=max_length,
↳padding='max_length', truncation=True, return_tensors="pt")

# Run the original PyTorch model on compilation exaple
paraphrase_classification_logits = model(**paraphrase)[0]

# Convert example inputs to a format that is compatible with TorchScript tracing
example_inputs_paraphrase = paraphrase['input_ids'], paraphrase['attention_mask'],
↳paraphrase['token_type_ids']
example_inputs_not_paraphrase = not_paraphrase['input_ids'], not_paraphrase['attention_
↳mask'], not_paraphrase['token_type_ids']

# Run torch.neuron.trace to generate a TorchScript that is optimized by AWS Neuron
model_neuron = torch.neuron.trace(model, example_inputs_paraphrase)

# Verify the TorchScript works on both example inputs
paraphrase_classification_logits_neuron = model_neuron(*example_inputs_paraphrase)
not_paraphrase_classification_logits_neuron = model_neuron(*example_inputs_not_
↳paraphrase)

# Save the TorchScript for later use
model_neuron.save('bert_neuron.pt')
```

You may inspect `model_neuron.graph` to see which part is running on CPU versus running on the accelerator. All native aten operators in the graph will be running on CPU.

```
[ ]: print(model_neuron.graph)
```

## Deploy the AWS Neuron optimized TorchScript

To deploy the AWS Neuron optimized TorchScript, you may choose to load the saved TorchScript from disk and skip the slow compilation.

```
[ ]: # Load TorchScript back
model_neuron = torch.jit.load('bert_neuron.pt')
# Verify the TorchScript works on both example inputs
paraphrase_classification_logits_neuron = model_neuron(*example_inputs_paraphrase)
not_paraphrase_classification_logits_neuron = model_neuron(*example_inputs_not_
↳paraphrase)
classes = ['not paraphrase', 'paraphrase']
paraphrase_prediction = paraphrase_classification_logits_neuron[0][0].argmax().item()
not_paraphrase_prediction = not_paraphrase_classification_logits_neuron[0][0].argmax().
↳item()
print('BERT says that "{}" and "{}" are {}'.format(sequence_0, sequence_2,
↳classes[paraphrase_prediction]))
print('BERT says that "{}" and "{}" are {}'.format(sequence_0, sequence_1, classes[not_
↳paraphrase_prediction]))
```

Now let's run the model in parallel on four cores

```
[ ]: def get_input_with_padding(batch, batch_size, max_length):
    ## Reformulate the batch into three batch tensors - default batch size batches the
    ↳outer dimension
    encoded = batch['encoded']
    inputs = torch.squeeze(encoded['input_ids'], 1)
    attention = torch.squeeze(encoded['attention_mask'], 1)
    token_type = torch.squeeze(encoded['token_type_ids'], 1)
    quality = list(map(int, batch['quality']))

    if inputs.size()[0] != batch_size:
        print("Input size = {} - padding".format(inputs.size()))
        remainder = batch_size - inputs.size()[0]
        zeros = torch.zeros( [remainder, max_length], dtype=torch.long )
        inputs = torch.cat( [inputs, zeros] )
        attention = torch.cat( [attention, zeros] )
        token_type = torch.cat( [token_type, zeros] )

    assert(inputs.size()[0] == batch_size and inputs.size()[1] == max_length)
    assert(attention.size()[0] == batch_size and attention.size()[1] == max_length)
    assert(token_type.size()[0] == batch_size and token_type.size()[1] == max_length)

    return (inputs, attention, token_type), quality

def count(output, quality):
    assert output.size(0) >= len(quality)
    correct_count = 0
    count = len(quality)
```

(continues on next page)

(continued from previous page)

```

batch_predictions = [ row.argmax().item() for row in output ]

for a, b in zip(batch_predictions, quality):
    if int(a)==int(b):
        correct_count += 1

return correct_count, count

```

## Data parallel inference

In the below cell, we use the data parallel approach for inference. In this approach, we load multiple models, all of them running in parallel. Each model is loaded onto a single NeuronCore. In the below implementation, we launch 16 models, thereby utilizing all the 16 cores on an inf1.6xlarge.

Note: Now if you try to decrease the num\_cores in the above cells, please restart the notebook and run `!sudo rmmmod neuron; sudo modprobe neuron` step in cell 2 to clear the Neuron cores.

Since, we can run more than 1 model concurrently, the throughput for the system goes up. To achieve maximum gain in throughput, we need to efficiently feed the models so as to keep them busy at all times. In the below setup, this is done by using a producer-consumer model. We maintain a common python queue shared across all the models. The common queue enables feeding data continuously to the models.

```

[ ]: from parallel import NeuronSimpleDataParallel
from bert_benchmark_utils import BertTestDataset, BertResults
import time
import functools

max_length = 128
num_cores = 16
batch_size = 1

tsv_file="glue_mrpc_dev.tsv"

data_set = BertTestDataset( tsv_file=tsv_file, tokenizer=tokenizer, max_length=max_
↪length )
data_loader = torch.utils.data.DataLoader(data_set, batch_size=batch_size, shuffle=True)

#Result aggregation class (code in bert_benchmark_utils.py)
results = BertResults(batch_size, num_cores)
def result_handler(output, result_id, start, end, input_dict):
    correct_count, inference_count = count(output[0], input_dict.pop(result_id))
    elapsed = end - start
    results.add_result(correct_count, inference_count, [elapsed], [end], [start])

parallel_neuron_model = NeuronSimpleDataParallel('bert_neuron.pt', num_cores)

#Starting the inference threads
parallel_neuron_model.start_continuous_inference()

# Warm up the cores
z = torch.zeros( [batch_size, max_length], dtype=torch.long )
batch = (z, z, z)

```

(continues on next page)



(continued from previous page)

```

for _ in range(num_cores*4):
    parallel_neuron_model.infer(batch, -1, None)

input_dict = {}
input_id = 0
for _ in range(30):
    for batch in data_loader:
        batch, quality = get_input_with_padding(batch, batch_size, max_length)
        input_dict[input_id] = quality
        callback_fn = functools.partial(result_handler, input_dict=input_dict)
        parallel_neuron_model.infer(batch, input_id, callback_fn)
        input_id+=1

# Stop inference
parallel_neuron_model.stop()

with open("benchmark.txt", "w") as f:
    results.report(f, window_size=1)

with open("benchmark.txt", "r") as f:
    for line in f:
        print(line)

```

Now recompile with a larger batch size of six sentence pairs

```

[ ]: batch_size = 6

example_inputs_paraphrase = (
    torch.cat([paraphrase['input_ids']] * batch_size,0),
    torch.cat([paraphrase['attention_mask']] * batch_size,0),
    torch.cat([paraphrase['token_type_ids']] * batch_size,0)
)

# Run torch.neuron.trace to generate a TorchScript that is optimized by AWS Neuron
model_neuron_batch = torch.neuron.trace(model, example_inputs_paraphrase)

## Save the batched model
model_neuron_batch.save('bert_neuron_b{}.pt'.format(batch_size))

```

Rerun inference with batch 6

```

[ ]: from parallel import NeuronSimpleDataParallel
from bert_benchmark_utils import BertTestDataset, BertResults
import time
import functools

max_length = 128
num_cores = 16
batch_size = 6

data_set = BertTestDataset( tsv_file=tsv_file, tokenizer=tokenizer, max_length=max_

```

(continues on next page)

(continued from previous page)

```

↪ length )
data_loader = torch.utils.data.DataLoader(data_set, batch_size=batch_size, shuffle=True)

#Result aggregation class (code in bert_benchmark_utils.py)
results = BertResults(batch_size, num_cores)
def result_handler(output, result_id, start, end, input_dict):
    correct_count, inference_count = count(output[0], input_dict.pop(result_id))
    elapsed = end - start
    results.add_result(correct_count, inference_count, [elapsed], [end], [start])

parallel_neuron_model = NeuronSimpleDataParallel('bert_neuron_b{}.pt'.format(batch_size),
↪ num_cores)

#Starting the inference threads
parallel_neuron_model.start_continuous_inference()

# Adding to the input queue to warm all cores
z = torch.zeros( [batch_size, max_length], dtype=torch.long )
batch = (z, z, z)
for _ in range(num_cores*4):
    parallel_neuron_model.infer(batch, -1, None)

input_dict = {}
input_id = 0
for _ in range(30):
    for batch in data_loader:
        batch, quality = get_input_with_padding(batch, batch_size, max_length)
        input_dict[input_id] = quality
        callback_fn = functools.partial(result_handler, input_dict=input_dict)
        parallel_neuron_model.infer(batch, input_id, callback_fn)
        input_id+=1

# Stop inference
parallel_neuron_model.stop()

with open("benchmark_b{}.txt".format(batch_size), "w") as f:
    results.report(f, window_size=1)

with open("benchmark_b{}.txt".format(batch_size), "r") as f:
    for line in f:
        print(line)

```

## Data Parallel HuggingFace Pretrained BERT with Weight Sharing (Deduplication)

### Introduction

In this tutorial we will compile and deploy BERT-base version of HuggingFace Transformers BERT for Inferentia, with additional demonstration of using Weight Sharing (Deduplication) feature.

To use the [Weight Sharing \(Deduplication\)](#) feature, you must set the Neuron Runtime environmental variable `NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS` to “TRUE” together with the [core placement API](#) (`torch_neuron.experimental.neuron_cores_context()`).

This Jupyter notebook should be run on an instance which is `inf1.6xlarge` or larger. The compile part of this tutorial requires `inf1.6xlarge` and not the inference itself. For simplicity we will run this tutorial on `inf1.6xlarge` but in real life scenario the compilation should be done on a compute instance and the deployment on `inf1` instance to save costs.

Verify that this Jupyter notebook is running the Python kernel environment that was set up according to the [PyTorch Installation Guide](#). You can select the kernel from the “Kernel -> Change Kernel” option on the top of this Jupyter notebook page.

### Install Dependencies:

This tutorial requires the following pip packages:

- `torch-neuron`
- `neuron-cc[tensorflow]`
- `transformers`

Most of these packages will be installed when configuring your environment using the Neuron PyTorch setup guide. The additional dependencies must be installed here.

```
[1]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
    ↪detect
    !pip install --upgrade "transformers==4.6.0"
```

### Compile the model into an AWS Neuron optimized TorchScript

This step compiles the model into an AWS Neuron optimized TorchScript, and saves it in the file `bert_neuron.pt`. This step is the same as the pretrained BERT tutorial without Shared Weights feature. We use batch 1 for simplicity.

```
[1]: import tensorflow # to workaround a protobuf version conflict issue
import torch
import torch.neuron
from transformers import AutoTokenizer, AutoModelForSequenceClassification, AutoConfig
import transformers
import os
import warnings

# Build tokenizer and model
tokenizer = AutoTokenizer.from_pretrained("bert-base-cased-finetuned-mrpc")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-cased-finetuned-
    ↪mrpc", return_dict=False)
```

(continues on next page)

(continued from previous page)

```

# Setup some example inputs
sequence_0 = "The company HuggingFace is based in New York City"
sequence_1 = "Apples are especially bad for your health"
sequence_2 = "HuggingFace's headquarters are situated in Manhattan"

max_length=128
paraphrase = tokenizer.encode_plus(sequence_0, sequence_2, max_length=max_length,
    ↪padding='max_length', truncation=True, return_tensors="pt")
not_paraphrase = tokenizer.encode_plus(sequence_0, sequence_1, max_length=max_length,
    ↪padding='max_length', truncation=True, return_tensors="pt")

# Run the original PyTorch model on compilation exaple
paraphrase_classification_logits = model(**paraphrase)[0]

# Convert example inputs to a format that is compatible with TorchScript tracing
example_inputs_paraphrase = paraphrase['input_ids'], paraphrase['attention_mask'],
    ↪paraphrase['token_type_ids']
example_inputs_not_paraphrase = not_paraphrase['input_ids'], not_paraphrase['attention_
    ↪mask'], not_paraphrase['token_type_ids']

# Run torch.neuron.trace to generate a TorchScript that is optimized by AWS Neuron
model_neuron = torch.neuron.trace(model, example_inputs_paraphrase)

# Verify the TorchScript works on both example inputs
paraphrase_classification_logits_neuron = model_neuron(*example_inputs_paraphrase)
not_paraphrase_classification_logits_neuron = model_neuron(*example_inputs_not_
    ↪paraphrase)

# Save the TorchScript for later use
model_neuron.save('bert_neuron.pt')

```

## Deploy the AWS Neuron optimized TorchScript

To deploy the AWS Neuron optimized TorchScript, you may choose to load the saved TorchScript from disk and skip the slow compilation. This step is the same as the pretrained BERT tutorial without Shared Weights feature

```

[2]: # Load TorchScript back
model_neuron = torch.jit.load('bert_neuron.pt')
# Verify the TorchScript works on both example inputs
paraphrase_classification_logits_neuron = model_neuron(*example_inputs_paraphrase)
not_paraphrase_classification_logits_neuron = model_neuron(*example_inputs_not_
    ↪paraphrase)
classes = ['not paraphrase', 'paraphrase']
paraphrase_prediction = paraphrase_classification_logits_neuron[0][0].argmax().item()
not_paraphrase_prediction = not_paraphrase_classification_logits_neuron[0][0].argmax().
    ↪item()
print('BERT says that "{}" and "{}" are {}'.format(sequence_0, sequence_2,
    ↪classes[paraphrase_prediction]))
print('BERT says that "{}" and "{}" are {}'.format(sequence_0, sequence_1, classes[not_
    ↪paraphrase_prediction]))

```

We define two helper functions to pad input and to count correct results.

```
[3]: def get_input_with_padding(batch, batch_size, max_length):
    ## Reformulate the batch into three batch tensors - default batch size batches the
    ↪ outer dimension
    encoded = batch['encoded']
    inputs = torch.squeeze(encoded['input_ids'], 1)
    attention = torch.squeeze(encoded['attention_mask'], 1)
    token_type = torch.squeeze(encoded['token_type_ids'], 1)
    quality = list(map(int, batch['quality']))

    if inputs.size()[0] != batch_size:
        print("Input size = {} - padding".format(inputs.size()))
        remainder = batch_size - inputs.size()[0]
        zeros = torch.zeros( [remainder, max_length], dtype=torch.long )
        inputs = torch.cat( [inputs, zeros] )
        attention = torch.cat( [attention, zeros] )
        token_type = torch.cat( [token_type, zeros] )

    assert(inputs.size()[0] == batch_size and inputs.size()[1] == max_length)
    assert(attention.size()[0] == batch_size and attention.size()[1] == max_length)
    assert(token_type.size()[0] == batch_size and token_type.size()[1] == max_length)

    return (inputs, attention, token_type), quality

def count(output, quality):
    assert output.size(0) >= len(quality)
    correct_count = 0
    count = len(quality)

    batch_predictions = [ row.argmax().item() for row in output ]

    for a, b in zip(batch_predictions, quality):
        if int(a)==int(b):
            correct_count += 1

    return correct_count, count
```

## Data parallel inference

In the below cell, we use the data parallel approach for inference. In this approach, we load multiple models, all of them running in parallel. Each model is loaded onto a single NeuronCore via the core placement API (`torch_neuron.experimental.neuron_cores_context()`). We also set Neuron Runtime environment variable `NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS` to “TRUE” as required to use the Weight Sharing feature.

In the below implementation, we launch 16 models, thereby utilizing all the 16 cores on an inf1.6xlarge.

Note: Now if you try to decrease the `num_cores` in the below cells, please restart the notebook and run `!sudo rmmmod neuron; sudo modprobe neuron` step in cell 2 to clear the Neuron cores.

Since, we can run more than 1 model concurrently, the throughput for the system goes up. To achieve maximum gain in throughput, we need to efficiently feed the models so as to keep them busy at all times. In the below setup, we use parallel threads to feed data continuously to the models.

When running the cell below, you can monitor the Inferentia device activities by running `neuron-top` in another

terminal. You will see that “Device Used Memory” is 1.6GB total, and the model instance loaded onto NeuronDevice 0 NeuronCore 0 uses the most device memory (272MB) while the other model instances loaded onto other NeuronCores use less device memory (92MB). This shows the effect of using Shared Weights as the device memory usage is lower. If you change `NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS` to “FALSE” you will see that “Device Used Memory” is 3.2GB, and the model instances loaded onto NeuronDevice 0 NeuronCore 0 and 1 use the most device memory (360MB) while the other model instances now use 180MB each.

```
[5]: from bert_benchmark_utils import BertTestDataset, BertResults
import time
import functools
import os
import torch.neuron as torch_neuron
from concurrent import futures

# Setting up NeuronCore groups for inf1.6xlarge with 16 cores
num_cores = 16 # This value should be 4 on inf1.xlarge and inf1.2xlarge
os.environ['NEURON_RT_NUM_CORES'] = str(num_cores)
os.environ['NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS'] = 'TRUE'
#os.environ['NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS'] = 'FALSE'

max_length = 128
num_cores = 16
batch_size = 1

tsv_file="glue_mrpc_dev.tsv"

data_set = BertTestDataset( tsv_file=tsv_file, tokenizer=tokenizer, max_length=max_
    ↪length )
data_loader = torch.utils.data.DataLoader(data_set, batch_size=batch_size, shuffle=True)

#Result aggregation class (code in bert_benchmark_utils.py)
results = BertResults(batch_size, num_cores)
def result_handler(output, result_id, start, end, input_dict):
    correct_count, inference_count = count(output[0], input_dict.pop(result_id))
    elapsed = end - start
    results.add_result(correct_count, inference_count, [elapsed], [end], [start])

with torch_neuron.experimental.neuron_cores_context(start_nc=0, nc_count=num_cores):
    model = torch.jit.load('bert_neuron.pt')

# Warm up the cores
z = torch.zeros( [batch_size, max_length], dtype=torch.long )
batch = (z, z, z)
for _ in range(num_cores*4):
    model(*batch)

# Prepare the input data
batch_list = []
for batch in data_loader:
    batch, quality = get_input_with_padding(batch, batch_size, max_length)
    batch_list.append((batch, quality))

# One thread running a model on one core
```

(continues on next page)

(continued from previous page)

```

def one_thread(feed_data, quality):
    start = time.time()
    result = model(*feed_data)
    end = time.time()
    return result[0], quality, start, end

# Launch more threads than models/cores to keep them busy
processes = []
with futures.ThreadPoolExecutor(max_workers=num_cores*2) as executor:
    # extra loops to help you see activities in neuron-top
    for _ in range(10):
        for input_id, (batch, quality) in enumerate(batch_list):
            processes.append(executor.submit(one_thread, batch, quality))

results = BertResults(batch_size, num_cores)
for _ in futures.as_completed(processes):
    (output, quality, start, end) = _.result()
    correct_count, inference_count = count(output, quality)
    results.add_result(correct_count, inference_count, [start - end], [start], [end])

with open("benchmark.txt", "w") as f:
    results.report(f, window_size=1)

with open("benchmark.txt", "r") as f:
    for line in f:
        print(line)

```

[ ]:

## Deploy a pretrained PyTorch BERT model from HuggingFace on Amazon SageMaker with Neuron container

### Overview

In this tutorial we will deploy on SageMaker a pretrained BERT Base model from HuggingFace Transformers, using the [AWS Deep Learning Containers](#). We will use the same model as shown in the [Neuron Tutorial “PyTorch - HuggingFace Pretrained BERT Tutorial”](#). We will compile the model and build a custom AWS Deep Learning Container, to include the HuggingFace Transformers Library.

This Jupyter Notebook should run on a ml.c5.4xlarge SageMaker Notebook instance. You can set up your SageMaker Notebook instance by following the [Get Started with Amazon SageMaker Notebook Instances](#) documentation.

We recommend increasing the size of the base root volume of your SM notebook instance, to accommodate the models and containers built locally. A root volume of 10Gb should suffice.

## Install Dependencies:

This tutorial requires the following pip packages:

- torch-neuron
- neuron-cc[tensorflow]
- transformers

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
    ↪ detect
    !pip install --upgrade --no-cache-dir torch-neuron neuron-cc[tensorflow] torchvision
    ↪ torch --extra-index-url=https://pip.repos.neuron.amazonaws.com
    !pip install --upgrade --no-cache-dir 'transformers==4.6.0'
```

## Compile the model into an AWS Neuron optimized TorchScript

```
[ ]: import torch
    import torch_neuron

    from transformers import AutoTokenizer, AutoModelForSequenceClassification, AutoConfig
```

```
[ ]: # Build tokenizer and model
    tokenizer = AutoTokenizer.from_pretrained("bert-base-cased-finetuned-mrpc")
    model = AutoModelForSequenceClassification.from_pretrained("bert-base-cased-finetuned-
    ↪ mrpc", return_dict=False)

    # Setup some example inputs
    sequence_0 = "The company HuggingFace is based in New York City"
    sequence_1 = "Apples are especially bad for your health"
    sequence_2 = "HuggingFace's headquarters are situated in Manhattan"

    max_length=128
    paraphrase = tokenizer.encode_plus(sequence_0, sequence_2, max_length=max_length,
    ↪ padding='max_length', truncation=True, return_tensors="pt")
    not_paraphrase = tokenizer.encode_plus(sequence_0, sequence_1, max_length=max_length,
    ↪ padding='max_length', truncation=True, return_tensors="pt")

    # Run the original PyTorch model on compilation exaple
    paraphrase_classification_logits = model(**paraphrase)[0]

    # Convert example inputs to a format that is compatible with TorchScript tracing
    example_inputs_paraphrase = paraphrase['input_ids'], paraphrase['attention_mask'],
    ↪ paraphrase['token_type_ids']
    example_inputs_not_paraphrase = not_paraphrase['input_ids'], not_paraphrase['attention_
    ↪ mask'], not_paraphrase['token_type_ids']
```

```
[ ]: %%time
    # Run torch.neuron.trace to generate a TorchScript that is optimized by AWS Neuron
    # This step may need 3-5 min
    model_neuron = torch.neuron.trace(model, example_inputs_paraphrase, verbose=1, compiler_
    ↪ workdir='./compilation_artifacts')
```



You may inspect `model_neuron.graph` to see which part is running on CPU versus running on the accelerator. All native **aten** operators in the graph will be running on CPU.

```
[ ]: # See which part is running on CPU versus running on the accelerator.
print(model_neuron.graph)
```

Save the compiled model, so it can be packaged and sent to S3.

```
[ ]: # Save the TorchScript for later use
model_neuron.save('neuron_compiled_model.pt')
```

### Package the pre-trained model and upload it to S3

To make the model available for the SageMaker deployment, you will TAR the serialized graph and upload it to the default Amazon S3 bucket for your SageMaker session.

```
[ ]: # Now you'll create a model.tar.gz file to be used by SageMaker endpoint
!tar -czvf model.tar.gz neuron_compiled_model.pt
```

```
[ ]: import boto3
import time
from sagemaker.utils import name_from_base
import sagemaker
```

```
[ ]: # upload model to S3
role = sagemaker.get_execution_role()
sess=sagemaker.Session()
region=sess.boto_region_name
bucket=sess.default_bucket()
sm_client=boto3.client('sagemaker')
```

```
[ ]: model_key = '{}'/model/model.tar.gz'.format('inf1_compiled_model')
model_path = 's3://{}/{}'.format(bucket, model_key)
boto3.resource('s3').Bucket(bucket).upload_file('model.tar.gz', model_key)
print("Uploaded model to S3:")
print(model_path)
```

### Build and Push the container

The following shell code shows how to build the container image using docker build and push the container image to ECR using docker push. The Dockerfile in this example is available in the **container** folder. Here's an example of the Dockerfile:

```
FROM 763104351884.dkr.ecr.us-east-1.amazonaws.com/pytorch-inference-neuron:1.7.1-neuron-
    py36-ubuntu18.04
# Install packages
RUN pip install "transformers==4.7.0"
```

```
[ ]: !cat container/Dockerfile
```

Before running the next cell, make sure your SageMaker IAM role has access to ECR. If not, you can attach the role AmazonEC2ContainerRegistryPowerUser to your IAM role ARN, which allows you to upload image layers to ECR.

It takes 5 minutes to build docker images and upload image to ECR

```
[ ]: %%sh

# The name of our algorithm
algorithm_name=neuron-py36-inference

cd container

account=$(aws sts get-caller-identity --query Account --output text)

# Get the region defined in the current configuration (default to us-west-2 if none
↳defined)
region=$(aws configure get region)
region=${region:-us-west-2}

fullname="${account}.dkr.ecr.${region}.amazonaws.com/${algorithm_name}:latest"

# If the repository doesn't exist in ECR, create it.

aws ecr describe-repositories --repository-names "${algorithm_name}" > /dev/null 2>&1

if [ $? -ne 0 ]
then
    aws ecr create-repository --repository-name "${algorithm_name}" > /dev/null
fi

# Get the login command from ECR in order to pull down the SageMaker PyTorch image
aws ecr get-login-password --region us-east-1 | docker login --username AWS --password-
↳stdin 763104351884.dkr.ecr.us-east-1.amazonaws.com
# Build the docker image locally with the image name and then push it to ECR
# with the full name.
docker build -t ${algorithm_name} . --build-arg REGION=${region}
docker tag ${algorithm_name} ${fullname}

# Get the login command from ECR and execute it directly
aws ecr get-login-password --region ${region} | docker login --username AWS --password-
↳stdin ${account}.dkr.ecr.${region}.amazonaws.com
docker push ${fullname}
```

## Deploy Container and run inference based on the pretrained model

To deploy a pretrained PyTorch model, you'll need to use the PyTorch estimator object to create a PyTorchModel object and set a different entry\_point.

You'll use the PyTorchModel object to deploy a PyTorchPredictor. This creates a SageMaker Endpoint – a hosted prediction service that we can use to perform inference.

```
[ ]: import sys

!{sys.executable} -m pip install Transformers

[ ]: import os
import boto3
import sagemaker

role = sagemaker.get_execution_role()
sess = sagemaker.Session()

bucket = sess.default_bucket()
prefix = "inf1_compiled_model/model"

# Get container name in ECR
client=boto3.client('sts')
account=client.get_caller_identity()['Account']

my_session=boto3.session.Session()
region=my_session.region_name

algorithm_name="neuron-py36-inference"
ecr_image='{ }.dkr.ecr.{ }.amazonaws.com/{ }:latest'.format(account, region, algorithm_name)
print(ecr_image)
```

An implementation of *model\_fn* is required for inference script. We are going to implement our own **model\_fn** and **predict\_fn** for Hugging Face Bert, and use default implementations of **input\_fn** and **output\_fn** defined in sagemaker-pytorch-containers.

In this example, the inference script is put in **code** folder. Run the next cell to see it:

```
[ ]: !pygmentize code/inference.py
```

Path of compiled pretrained model in S3:

```
[ ]: key = os.path.join(prefix, "model.tar.gz")
pretrained_model_data = "s3://{}/{ }".format(bucket, key)
print(pretrained_model_data)
```

The model object is defined by using the SageMaker Python SDK's PyTorchModel and pass in the model from the estimator and the entry\_point. The endpoint's entry point for inference is defined by model\_fn as seen in the previous code block that prints out **inference.py**. The model\_fn function will load the model and required tokenizer.

Note, **image\_uri** must be user's own ECR images.

```
[ ]: from sagemaker.pytorch.model import PyTorchModel
```

(continues on next page)

(continued from previous page)

```
pytorch_model = PyTorchModel(
    model_data=pretrained_model_data,
    role=role,
    source_dir="code",
    framework_version="1.7.1",
    entry_point="inference.py",
    image_uri=ecr_image
)

# Let SageMaker know that we've already compiled the model via neuron-cc
pytorch_model._is_compiled_model = True
```

The arguments to the deploy function allow us to set the number and type of instances that will be used for the Endpoint.

Here you will deploy the model to a single **ml.inf1.2xlarge** instance. It may take 6-10 min to deploy.

```
[ ]: %%time

predictor = pytorch_model.deploy(initial_instance_count=1, instance_type="ml.inf1.2xlarge
→")
```

```
[ ]: print(predictor.endpoint_name)
```

Since in the `input_fn` we declared that the incoming requests are json-encoded, we need to use a json serializer, to encode the incoming data into a json string. Also, we declared the return content type to be json string, we need to use a json deserializer to parse the response.

```
[ ]: predictor.serializer = sagemaker.serializers.JSONSerializer()
predictor.deserializer = sagemaker.deserializers.JSONDeserializer()
```

Using a list of sentences, now SageMaker endpoint is invoked to get predictions.

```
[ ]: %%time
result = predictor.predict(
    [
        "Never allow the same bug to bite you twice.",
        "The best part of Amazon SageMaker is that it makes machine learning easy.",
    ]
)
print(result)
```

```
[ ]: %%time
result = predictor.predict(
    [
        "The company HuggingFace is based in New York City",
        "HuggingFace's headquarters are situated in Manhattan",
    ]
)
print(result)
```

## Benchmarking your endpoint

The following cells create a load test for your endpoint. You first define some helper functions: `inference_latency` runs the endpoint request, collects client side latency and any errors, `random_sentence` builds random to be sent to the endpoint.

```
[ ]: import numpy as np
import datetime
import math
import time
import boto3
import matplotlib.pyplot as plt
from joblib import Parallel, delayed
import numpy as np
from tqdm import tqdm
import random
```

```
[ ]: def inference_latency(model,*inputs):
    """
    inference_time is a simple method to return the latency of a model inference.

    Parameters:
        model: torch model object loaded using torch.jit.load
        inputs: model() args

    Returns:
        latency in seconds
    """
    error = False
    start = time.time()
    try:
        results = model(*inputs)
    except:
        error = True
        results = []
    return {'latency':time.time() - start, 'error': error, 'result': results}
```

```
[ ]: def random_sentence():

    s_nouns = ["A dude", "My mom", "The king", "Some guy", "A cat with rabies", "A sloth",
    ↪, "Your homie", "This cool guy my gardener met yesterday", "Superman"]
    p_nouns = ["These dudes", "Both of my moms", "All the kings of the world", "Some guys",
    ↪, "All of a cattery's cats", "The multitude of sloths living under your bed", "Your",
    ↪homies", "Like, these, like, all these people", "Supermen"]
    s_verbs = ["eats", "kicks", "gives", "treats", "meets with", "creates", "hacks",
    ↪, "configures", "spies on", "retards", "meows on", "flees from", "tries to automate",
    ↪, "explodes"]
    p_verbs = ["eat", "kick", "give", "treat", "meet with", "create", "hack", "configure",
    ↪, "spy on", "retard", "meow on", "flee from", "try to automate", "explode"]
    infinitives = ["to make a pie.", "for no apparent reason.", "because the sky is",
    ↪green.", "for a disease.", "to be able to make toast explode.", "to know more about",
    ↪archeology."]
```

(continues on next page)

(continued from previous page)

```

    return (random.choice(s_nouns) + ' ' + random.choice(s_verbs) + ' ' + random.
    ↳choice(s_nouns).lower() or random.choice(p_nouns).lower() + ' ' + random.
    ↳choice(infinitives))

print([random_sentence(), random_sentence()])

```

The following cell creates `number_of_clients` concurrent threads to run `number_of_runs` requests. Once completed, a `boto3` CloudWatch client will query for the server side latency metrics for comparison.

```

[ ]: # Defining Auxiliary variables
number_of_clients = 2
number_of_runs = 1000
t = tqdm(range(number_of_runs), position=0, leave=True)

# Starting parallel clients
cw_start = datetime.datetime.utcnow()

results = Parallel(n_jobs=number_of_clients, prefer="threads")(delayed(inference_
    ↳latency)(predictor.predict, [random_sentence(), random_sentence()]) for mod in t)
avg_throughput = t.total/t.format_dict['elapsed']

cw_end = datetime.datetime.utcnow()

# Computing metrics and print
latencies = [res['latency'] for res in results]
errors = [res['error'] for res in results]
error_p = sum(errors)/len(errors) *100
p50 = np.quantile(latencies[-1000:],0.50) * 1000
p90 = np.quantile(latencies[-1000:],0.95) * 1000
p95 = np.quantile(latencies[-1000:],0.99) * 1000

print(f'Avg Throughput: :{avg_throughput:.1f}\n')
print(f'50th Percentile Latency:{p50:.1f} ms')
print(f'90th Percentile Latency:{p90:.1f} ms')
print(f'95th Percentile Latency:{p95:.1f} ms\n')
print(f'Errors percentage: {error_p:.1f} %\n')

# Querying CloudWatch
print('Getting Cloudwatch:')
cloudwatch = boto3.client('cloudwatch')
statistics=['SampleCount', 'Average', 'Minimum', 'Maximum']
extended=['p50', 'p90', 'p95', 'p100']

# Give 5 minute buffer to end
cw_end += datetime.timedelta(minutes=5)

# Period must be 1, 5, 10, 30, or multiple of 60
# Calculate closest multiple of 60 to the total elapsed time
factor = math.ceil((cw_end - cw_start).total_seconds() / 60)
period = factor * 60
print('Time elapsed: {} seconds'.format((cw_end - cw_start).total_seconds()))

```

(continues on next page)

(continued from previous page)

```

print('Using period of {} seconds\n'.format(period))

cloudwatch_ready = False
# Keep polling CloudWatch metrics until datapoints are available
while not cloudwatch_ready:
    time.sleep(30)
    print('Waiting 30 seconds ...')
    # Must use default units of microseconds
    model_latency_metrics = cloudwatch.get_metric_statistics(MetricName='ModelLatency',
                                                              Dimensions=[{'Name': 'EndpointName',
                                                                      'Value': predictor.endpoint_
↪ name},
                                                              {'Name': 'VariantName',
                                                                      'Value': "AllTraffic"}],
                                                              Namespace="AWS/SageMaker",
                                                              StartTime=cw_start,
                                                              EndTime=cw_end,
                                                              Period=period,
                                                              Statistics=statistics,
                                                              ExtendedStatistics=extended
                                                              )

    # Should be 1000
    if len(model_latency_metrics['Datapoints']) > 0:
        print('{} latency datapoints ready'.format(model_latency_metrics['Datapoints'][0][
↪ 'SampleCount']))
        side_avg = model_latency_metrics['Datapoints'][0]['Average'] / number_of_runs
        side_p50 = model_latency_metrics['Datapoints'][0]['ExtendedStatistics']['p50'] / ↪
↪ number_of_runs
        side_p90 = model_latency_metrics['Datapoints'][0]['ExtendedStatistics']['p90'] / ↪
↪ number_of_runs
        side_p95 = model_latency_metrics['Datapoints'][0]['ExtendedStatistics']['p95'] / ↪
↪ number_of_runs
        side_p100 = model_latency_metrics['Datapoints'][0]['ExtendedStatistics']['p100'] / ↪
↪ number_of_runs

        print(f'50th Percentile Latency:{side_p50:.1f} ms')
        print(f'90th Percentile Latency:{side_p90:.1f} ms')
        print(f'95th Percentile Latency:{side_p95:.1f} ms\n')

    cloudwatch_ready = True

```

## Cleanup

Endpoints should be deleted when no longer in use, to avoid costs.

```
[ ]: predictor.delete_endpoint(predictor.endpoint)
```

```
[ ]:
```

*This document is relevant for: Inf1*

## BERT TorchServe Tutorial

### Table of Contents

- [Overview](#)
- [Run the tutorial](#)
- [Setup TorchServe](#)
- [Run TorchServe](#)
- [Benchmark TorchServe](#)

## Overview

This tutorial demonstrates the use of [TorchServe](#) with Neuron, the SDK for Amazon Inf1 instances. By the end of this tutorial, you will understand how TorchServe can be used to serve a model backed by EC2 Inf1 instances. We will use a pretrained BERT-Base model to determine if one sentence is a paraphrase of another.

Verify that this tutorial is running in a virtual environment that was set up according to the *Torch-Neuronx Installation Guide* <<https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/setup/torch-neuronx.html#setup-torch-neuronx>> or *Torch-Neuron Installation Guide* <<https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/setup/torch-neuron.html#setup-torch-neuron>>

## Run the tutorial

Open a terminal, log into your remote instance, and activate a Pytorch virtual environment setup (see the Pytorch Installation Guide). To complete this tutorial, you will need a compiled BERT model. If you have already completed the HuggingFace Pretrained BERT tutorial [\[html\]](#) [\[notebook\]](#) then you already have the necessary file. Otherwise, you can setup your environment as shown below and then run `trace_bert_neuron.py` to obtain a traced BERT model.

You should now have a compiled `bert_neuron_b6.pt` file, which is required going forward.

Open a shell on the instance you prepared earlier, create a new directory named `torchserve`. Copy your compiled model from the previous tutorial into this new directory.

```
cd torchserve
python trace_bert_neuronx.py
ls
```



```
bert_neuron_b6.pt
```

Prepare a new Python virtual environment with the necessary Neuron and TorchServe components. Use a virtual environment to keep (most of) the various tutorial components isolated from the rest of the system in a controlled way.

```
pip install transformers==4.20.1 torchserve==0.7.0 torch-model-archiver==0.7.0 captum==0.6.0
```

Install the system requirements for TorchServe.

### Amazon Linux 2 DLAMI Base

```
sudo yum install jq java-11-amazon-corretto-headless
sudo alternatives --config java
sudo alternatives --config javac
```

### Ubuntu 20 DLAMI Base

```
sudo apt install openjdk-11-jdk -y
```

```
java -version
```

```
openjdk version "11.0.17" 2022-10-18
OpenJDK Runtime Environment (build 11.0.17+8-post-Ubuntu-1ubuntu218.04)
OpenJDK 64-Bit Server VM (build 11.0.17+8-post-Ubuntu-1ubuntu218.04, mixed mode, sharing)
```

```
javac -version
```

```
javac 11.0.17
```

Verify that TorchServe is now available.

```
torchserve --version
```

```
TorchServe Version is 0.7.0
```

### Setup TorchServe

During this tutorial you will need to download a few files onto your instance. The simplest way to accomplish this is to paste the download links provided above each file into a `wget` command. (We don't provide the links directly because they are subject to change.) For example, right-click and copy the download link for `config.json` shown below.

Listing 2.6: `config.json`

```
{
  "model_name": "bert-base-cased-finetuned-mrpc",
  "max_length": 128,
  "batch_size": 6
}
```

Now execute the following in your shell:

```
wget <paste link here>
ls
```

```
bert_neuron_b6.pt  config.json
```

Download the [custom handler script](#) that will eventually respond to inference requests.

Listing 2.7: handler\_bert.py

```

1  import os
2  import json
3  import sys
4  import logging
5  from abc import ABC
6
7  import torch
8  import torch_neuron
9
10 from transformers import AutoTokenizer
11 from ts.torch_handler.base_handler import BaseHandler
12
13
14 # one core per worker
15 os.environ['NEURON_RT_NUM_CORES'] = '1'
16
17 logger = logging.getLogger(__name__)
18
19 class BertEmbeddingHandler(BaseHandler, ABC):
20     """
21     Handler class for Bert Embedding computations.
22     """
23     def __init__(self):
24         super(BertEmbeddingHandler, self).__init__()
25         self.initialized = False
26
27     def initialize(self, ctx):
28         self.manifest = ctx.manifest
29         properties = ctx.system_properties
30         self.device = 'cpu'
31         model_dir = properties.get('model_dir')
32         serialized_file = self.manifest['model']['serializedFile']
33         model_pt_path = os.path.join(model_dir, serialized_file)
34
35         # point sys.path to our config file
36         with open('config.json') as fp:
37             config = json.load(fp)
38             self.max_length = config['max_length']
39             self.batch_size = config['batch_size']
40             self.classes = ['not paraphrase', 'paraphrase']
41
42         self.model = torch.jit.load(model_pt_path)

```

(continues on next page)

(continued from previous page)

```

43     logger.debug(f'Model loaded from {model_dir}')
44     self.model.to(self.device)
45     self.model.eval()
46
47     self.tokenizer = AutoTokenizer.from_pretrained(config['model_name'])
48     self.initialized = True
49
50     def preprocess(self, input_data):
51         """
52         Tokenization pre-processing
53         """
54
55         input_ids = []
56         attention_masks = []
57         token_type_ids = []
58         for row in input_data:
59             seq_0 = row['seq_0'].decode('utf-8')
60             seq_1 = row['seq_1'].decode('utf-8')
61             logger.debug(f'Received text: "{seq_0}", "{seq_1}"')
62
63             inputs = self.tokenizer.encode_plus(
64                 seq_0,
65                 seq_1,
66                 max_length=self.max_length,
67                 padding='max_length',
68                 truncation=True,
69                 return_tensors='pt'
70             )
71
72             input_ids.append(inputs['input_ids'])
73             attention_masks.append(inputs['attention_mask'])
74             token_type_ids.append(inputs['token_type_ids'])
75
76         batch = (torch.cat(input_ids, 0),
77                 torch.cat(attention_masks, 0),
78                 torch.cat(token_type_ids, 0))
79
80         return batch
81
82     def inference(self, inputs):
83         """
84         Predict the class of a text using a trained transformer model.
85         """
86
87         # sanity check dimensions
88         assert(len(inputs) == 3)
89         num_inferences = len(inputs[0])
90         assert(num_inferences <= self.batch_size)
91
92         # insert padding if we received a partial batch
93         padding = self.batch_size - num_inferences
94         if padding > 0:

```

(continues on next page)

(continued from previous page)

```

95         pad = torch.nn.ConstantPad1d((0, 0, 0, padding), value=0)
96         inputs = [pad(x) for x in inputs]
97
98         outputs = self.model(*inputs)[0]
99         predictions = []
100         for i in range(num_inferences):
101             prediction = self.classes[outputs[i].argmax().item()]
102             predictions.append([prediction])
103             logger.debug("Model predicted: '%s'", prediction)
104         return predictions
105
106     def postprocess(self, inference_output):
107         return inference_output

```

Next, we need to associate the handler script with the compiled model using `torch-model-archiver`. Run the following commands in your terminal:

```

mkdir model_store
MAX_LENGTH=$(jq '.max_length' config.json)
BATCH_SIZE=$(jq '.batch_size' config.json)
MODEL_NAME=bert-max_length$MAX_LENGTH-batch_size$BATCH_SIZE
torch-model-archiver --model-name "$MODEL_NAME" --version 1.0 --serialized-file ./bert_
neuron_b6.pt --handler "./handler_bert_neuronx.py" --extra-files "./config.json" --
export-path model_store

```

**Note:** If you modify your model or a dependency, you will need to rerun the archiver command with the `-f` flag appended to update the archive.

The result of the above will be a `mar` file inside the `model_store` directory.

```
ls model_store
```

```
bert-max_length128-batch_size6.mar
```

This file is essentially an archive associated with a fixed version of your model along with its dependencies (e.g. the handler code).

**Note:** The version specified in the `torch-model-archiver` command can be appended to REST API requests to access a specific version of your model. For example, if your model was hosted locally on port 8080 and named “bert”, the latest version of your model would be available at `http://localhost:8080/predictions/bert`, while version 1.0 would be accessible at `http://localhost:8080/predictions/bert/1.0`. We will see how to perform inference using this API in Step 6.

Create a `custom config` file to set some parameters. This file will be used to configure the server at launch when we run `torchserve --start`.

Listing 2.8: `torchserve.config`

```

# bind inference API to all network interfaces with SSL enabled
inference_address=http://0.0.0.0:8080

```

(continues on next page)

(continued from previous page)

```
default_workers_per_model=1
```

**Note:** This will cause TorchServe to bind on all interfaces. For security in real-world applications, you'll probably want to use port 8443 and [enable SSL](#).

## Run TorchServe

It's time to start the server. Typically we'd want to launch this in a separate console, but for this demo we'll just redirect output to a file.

```
torchserve --start --ncs --model-store model_store --ts-config torchserve.config 2>&1 >
↪ torchserve.log
```

Verify that the server seems to have started okay.

```
curl http://127.0.0.1:8080/ping
```

```
{
  "status": "Healthy"
}
```

**Note:** If you get an error when trying to ping the server, you may have tried before the server was fully launched. Check `torchserve.log` for details.

Use the Management API to instruct TorchServe to load our model.

```
MAX_BATCH_DELAY=5000 # ms timeout before a partial batch is processed
INITIAL_WORKERS=2 # Number from table above
curl -X POST "http://localhost:8081/models?url=$MODEL_NAME.mar&batch_size=$BATCH_SIZE&
↪ initial_workers=$INITIAL_WORKERS&max_batch_delay=$MAX_BATCH_DELAY"
```

```
{
  "status": "Model \"bert-max_length128-batch_size6\" Version: 1.0 registered with 4
↪ initial workers"
}
```

**Note:** Any additional attempts to configure the model after the initial curl request will cause the server to return a 409 error. You'll need to stop/start/configure the server to realize any changes.

The `MAX_BATCH_DELAY` is a timeout value that determines how long to wait before processing a partial batch. This is why the handler code needs to check the batch dimension and potentially add padding. TorchServe will instantiate the number of model handlers indicated by `INITIAL_WORKERS`, so this value controls how many models we will load onto Inferentia in parallel. This tutorial was performed on an `inf1.xlarge` instance (one Inferentia chip), so there are four NeuronCores available. If you want to control worker scaling more dynamically, [see the docs](#).

**Warning:** If you attempt to load more models than NeuronCores available, one of two things will occur. Either the extra models will fit in device memory but performance will suffer, or you will encounter an error on your initial inference. You shouldn't set INITIAL\_WORKERS above the number of NeuronCores. However, you may want to use fewer cores if you are using the *NeuronCore Pipeline* feature.

It looks like everything is running successfully at this point, so it's time for an inference.

Create the `infer_bert.py` file below on your instance.

Listing 2.9: `infer_bert.py`

```

1 import json
2 import concurrent.futures
3 import requests
4
5 with open('config.json') as fp:
6     config = json.load(fp)
7     max_length = config['max_length']
8     batch_size = config['batch_size']
9     name = f'bert-max_length{max_length}-batch_size{batch_size}'
10
11 # dispatch requests in parallel
12 url = f'http://localhost:8080/predictions/{name}'
13 paraphrase = {'seq_0': "HuggingFace's headquarters are situated in Manhattan",
14               'seq_1': "The company HuggingFace is based in New York City"}
15 not_paraphrase = {'seq_0': paraphrase['seq_0'], 'seq_1': 'This is total nonsense.'}
16
17 with concurrent.futures.ThreadPoolExecutor(max_workers=batch_size) as executor:
18     def worker_thread(worker_index):
19         # we'll send half the requests as not_paraphrase examples for sanity
20         data = paraphrase if worker_index < batch_size//2 else not_paraphrase
21         try:
22             response = requests.post(url, data=data)
23
24             # Check if the response status code indicates success
25             if response.status_code == 200:
26                 print(worker_index, response.json())
27             else:
28                 # If the response is not successful, raise an exception with the status_
29                 ↪ code and error message
29                 error_message = response.json().get('message', 'Unknown Error')
30                 raise Exception(f"Failed request with status code {response.status_code}:
31                 ↪ {error_message}")
31             except Exception as e:
32                 # Catch all other exceptions that may be raised
33                 print(f"An unexpected error occurred: {e}")
34                 raise
35
36     for worker_index in range(batch_size):
37         executor.submit(worker_thread, worker_index)

```

This script will send a `batch_size` number of requests to our model. In this example, we are using a model that estimates the probability that one sentence is a paraphrase of another. The script sends positive examples in the first

half of the batch and negative examples in the second half.

Execute the script in your terminal.

```
python infer_bert.py
```

```
1 ['paraphrase']
3 ['not paraphrase']
4 ['not paraphrase']
0 ['paraphrase']
5 ['not paraphrase']
2 ['paraphrase']
```

We can see that the first three threads (0, 1, 2) all report `paraphrase`, as expected. If we instead modify the script to send an incomplete batch and then wait for the timeout to expire, the excess padding results will be discarded.

## Benchmark TorchServe

We've seen how to perform a single batched inference, but how many inferences can we process per second? A separate upcoming tutorial will document performance tuning to maximize throughput. In the meantime, we can still perform a simple naïve stress test. The code below will spawn 64 worker threads, with each thread repeatedly sending a full batch of data to process. A separate thread will periodically print throughput and latency measurements.

Listing 2.10: `benchmark_bert.py`

```
1 import os
2 import argparse
3 import time
4 import numpy as np
5 import requests
6 import sys
7 from concurrent import futures
8
9 import torch
10
11
12 parser = argparse.ArgumentParser()
13 parser.add_argument('--url', help='Torchserve model URL', type=str, default=f'http://127.
    ↳ 0.0.1:8080/predictions/bert-max_length128-batch_size6')
14 parser.add_argument('--num_thread', type=int, default=64, help='Number of threads_
    ↳ invoking the model URL')
15 parser.add_argument('--batch_size', type=int, default=6)
16 parser.add_argument('--sequence_length', type=int, default=128)
17 parser.add_argument('--latency_window_size', type=int, default=1000)
18 parser.add_argument('--throughput_time', type=int, default=300)
19 parser.add_argument('--throughput_interval', type=int, default=10)
20 args = parser.parse_args()
21
22 data = { 'seq_0': 'A completely made up sentence.',
23         'seq_1': 'Well, I suppose they are all made up.' }
24 live = True
25 num_infer = 0
26 latency_list = []
```

(continues on next page)

(continued from previous page)

```

27
28
29 def one_thread(pred, feed_data):
30     global latency_list
31     global num_infer
32     global live
33     session = requests.Session()
34     while True:
35         start = time.time()
36         result = session.post(pred, data=feed_data)
37         latency = time.time() - start
38         latency_list.append(latency)
39         num_infer += 1
40         if not live:
41             break
42
43
44 def current_performance():
45     last_num_infer = num_infer
46     for _ in range(args.throughput_time // args.throughput_interval):
47         current_num_infer = num_infer
48         throughput = (current_num_infer - last_num_infer) / args.throughput_interval
49         p50 = 0.0
50         p90 = 0.0
51         if latency_list:
52             p50 = np.percentile(latency_list[-args.latency_window_size:], 50)
53             p90 = np.percentile(latency_list[-args.latency_window_size:], 90)
54         print('pid {}: current throughput {}, latency p50={:.3f} p90={:.3f}'.format(os.
↳ getpid(), throughput, p50, p90))
55         sys.stdout.flush()
56         last_num_infer = current_num_infer
57         time.sleep(args.throughput_interval)
58     global live
59     live = False
60
61
62 with futures.ThreadPoolExecutor(max_workers=args.num_thread+1) as executor:
63     executor.submit(current_performance)
64     for _ in range(args.num_thread):
65         executor.submit(one_thread, args.url, data)

```

Run the benchmarking script.

```
python benchmark_bert.py
```

```

pid 28523: current throughput 0.0, latency p50=0.000 p90=0.000
pid 28523: current throughput 617.7, latency p50=0.092 p90=0.156
pid 28523: current throughput 697.3, latency p50=0.082 p90=0.154
pid 28523: current throughput 702.8, latency p50=0.081 p90=0.149
pid 28523: current throughput 699.1, latency p50=0.085 p90=0.147
pid 28523: current throughput 703.8, latency p50=0.083 p90=0.148
pid 28523: current throughput 699.3, latency p50=0.083 p90=0.148

```

(continues on next page)



(continued from previous page)

...

**Congratulations!** By now you should have successfully served a batched model over TorchServe.

You can now shutdown torchserve.

```
torchserve --stop
```

*This document is relevant for: Inf1*

## Transformers MarianMT Tutorial

In this tutorial, you will deploy the [HuggingFace MarianMT](#) model for text translation.

This Jupyter notebook should be run on an inf1.6xlarge instance since you will be loading and compiling several large models.

Verify that this Jupyter notebook is running the Python kernel environment that was set up according to the [PyTorch Installation Guide](#). You can select the kernel from the “Kernel -> Change Kernel” option on the top of this Jupyter notebook page.

To generate text, you will be using the beam search algorithm to incrementally generate token candidates until the full output text has been created. Unlike simple single-pass models, this algorithm divides the work into two distinct phases:

- **Encoder:** Convert the input text into an encoded representation. (Executed once)
- **Decoder:** Use the encoded representation of the input text and the current output tokens to incrementally generate the set of next best candidate tokens. (Executed many times)

In this tutorial you will perform the following steps:

- **Compile:** Compile both the Encoder and Decoder for Neuron using simplified interfaces for inference.
- **Infer:** Run on CPU and Neuron and compare results.

Finally, a completely unrolled decoder will be built which simplifies the implementation at the cost of performing fixed-length inferences.

## Install Dependencies:

This tutorial has the following dependencies:

- `transformers==4.25.1`
- `torch-neuron`
- `sentencepiece`
- `neuron-cc[tensorflow]`

The following will install the required `transformers` version. Note that encoder/decoder API changes across different minor versions requires that you are specific about the version used. Also note that the `torch-neuron` version is pinned due to transformer compatibility issues.

```
[ ]: !pip install sentencepiece transformers==4.26.1
```

## Parameters

The parameters of a generative model can be tuned for different use-cases. In this example, you'll tailor the parameters to a single inference beam search for an on-demand inference use-case. See the [MarianConfig](#) for parameter details.

Rather than varying the encoder/decoder token sizes at runtime, you must define these parameters prior to compilation. The encoder/decoder token sizes are important tunable parameters as a large token sequence will offer greater sentence length flexibility but perform worse than a small token sequence.

To maximize performance on Neuron, the `num_beams`, `max_encode_length` and `max_decoder_length` should be made as small as possible for the use-case.

For this tutorial you will use a model that translates sentences of up to 32 token from English to German.

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
    ↪ detect
model_name = "Helsinki-NLP/opus-mt-en-de" # English -> German model
num_texts = 1                               # Number of input texts to decode
num_beams = 4                               # Number of beams per input text
max_encoder_length = 32                     # Maximum input token length
max_decoder_length = 32                     # Maximum output token length
```

## CPU Model Inference

Start by executing the model on CPU to test its execution.

The following defines the inference function which will be used to compare the Neuron and CPU output. In this example you will display all beam search sequences that were generated. For a real on-demand use case, set the `num_beams` to 1 to return only the top result.

```
[ ]: def infer(model, tokenizer, text):

    # Truncate and pad the max length to ensure that the token size is compatible with
    ↪ fixed-sized encoder (Not necessary for pure CPU execution)
    batch = tokenizer(text, max_length=max_decoder_length, truncation=True, padding='max_
    ↪ length', return_tensors="pt")
    output = model.generate(**batch, max_length=max_decoder_length, num_beams=num_beams,
    ↪ num_return_sequences=num_beams)
    results = [tokenizer.decode(t, skip_special_tokens=True) for t in output]

    print('Texts:')
    for i, summary in enumerate(results):
        print(i + 1, summary)
```

Note that after loading the model, we also set the maximum length. This will later be used to limit the size of the compiled model.

```
[ ]: from transformers import MarianMTModel, MarianTokenizer

model_cpu = MarianMTModel.from_pretrained(model_name)
model_cpu.config.max_length = max_decoder_length
model_cpu.eval()

tokenizer = MarianTokenizer.from_pretrained(model_name)
```

(continues on next page)

(continued from previous page)

```
sample_text = "I am a small frog."
```

```
[ ]: infer(model_cpu, tokenizer, sample_text)
```

## Padded Model

In order to perform inference on Neuron, the model must be changed in a way that it supports tracing and fixed-sized inputs. One way in which this is possible is to use a pad the model inputs to the maximum possible tensor sizes. The benefit of using a padded model is that it supports variable length text generation up to a specified length `max_decoder_length`. A consequence of padding is that it can negatively impact performance due to large data transfers.

## PaddedEncoder & PaddedDecoder Modules

Here you will define wrappers around the encoder and decoder portions of the generation model that are compatible with `torch.jit.trace` as well as fixed-sized inputs.

The following are important features which are distinct from the default configuration:

1. Disabled `return_dict`. When this is enabled, the network uses `dataclass` type outputs which are not compatible with `torch.jit.trace`.
2. Disabled `use_cache`. When this option is enabled, the network expects a collection of cache tensors which grow upon each iteration. Since Neuron requires fixed sized inputs, this must be disabled.
3. The `GenerationMixin:beam_search` implementation uses only the logits for the current iteration index from the original decoder layer output. Since inputs must be padded, performance can be improved by selecting only a subset of the hidden state prior to the final linear layer. For efficiency on Neuron, this reduction uses an elementwise-multiply to mask out the unused hidden values and then sums along an axis.
4. Since a reduction step is inserted between the decoder output and the final logit calculation, the original `model` attribute is not used. Instead the `PaddedDecoder` class combines the decoder, reducer, and linear layers into a combined forward pass. In the original model there is a clear distinction between the decoder layer and the final linear layer. These layers are fused together to get one large fully optimized graph.

```
[ ]: import torch
from torch.nn import functional as F

class PaddedEncoder(torch.nn.Module):

    def __init__(self, model):
        super().__init__()
        self.encoder = model.model.encoder
        self.main_input_name = 'input_ids'

    def forward(self, input_ids, attention_mask):
        return self.encoder(input_ids, attention_mask=attention_mask, return_dict=False)

class PaddedDecoder(torch.nn.Module):
```

(continues on next page)

(continued from previous page)

```

def __init__(self, model):
    super().__init__()
    self.weight = model.model.shared.weight.clone().detach()
    self.bias = model.final_logits_bias.clone().detach()
    self.decoder = model.model.decoder

def forward(self, input_ids, attention_mask, encoder_outputs, index):

    # Invoke the decoder
    hidden, = self.decoder(
        input_ids=input_ids,
        encoder_hidden_states=encoder_outputs,
        encoder_attention_mask=attention_mask,
        return_dict=False,
        use_cache=False,
    )

    _, n_length, _ = hidden.shape

    # Create selection mask
    mask = torch.arange(n_length, dtype=torch.float32) == index
    mask = mask.view(1, -1, 1)

    # Broadcast mask
    masked = torch.multiply(hidden, mask)

    # Reduce along 1st dimension
    hidden = torch.sum(masked, 1, keepdims=True)

    # Compute final linear layer for token probabilities
    logits = F.linear(
        hidden,
        self.weight,
        bias=self.bias
    )
    return logits

```

### PaddedGenerator - GenerationMixin Class

On text generation tasks, HuggingFace Transformers defines a [GenerationMixin](#) base class which provides standard methods and algorithms to generate text. For this tutorial, you will be using the beam search algorithm on encoder/decoder architectures.

To be able to use these methods, you will be defining your own class derived from the [GenerationMixin](#) class to run a beam search. This will invoke the encoder and decoder layers in a way that is compatible with fixed sized inputs and traced modules. This means you must import the base class and the output objects ([Seq2SeqLMOutput](#), [BaseModelOutput](#)) used by the [beam\\_search](#) algorithm.

The [GenerationMixin:generate](#) method will use [GenerationMixin:beam\\_search](#) which requires that you to define your own class implementation that invokes the [PaddedEncoder](#) and [PaddedDecoder](#) modules using padded

inputs. The standard generator model implementation will not work by default because it is intended to infer with variable-sized (growing) input tensors.

The `from_model` method is defined to create the `PaddedGenerator` from an existing pretrained generator class.

To invoke the Encoder and Decoder traced modules in a way that is compatible with the `GenerationMixin: beam_search` implementation, the `get_encoder`, `__call__`, and `prepare_inputs_for_generation` methods are overridden.

Lastly, the class defines methods for serialization so that the model can be easily saved and loaded.

```
[ ]: import os

from transformers import GenerationMixin, AutoConfig
from transformers.modeling_outputs import Seq2SeqLMOutput, BaseModelOutput
from transformers.modeling_utils import PreTrainedModel

class PaddedGenerator(PreTrainedModel, GenerationMixin):

    @classmethod
    def from_model(cls, model):
        generator = cls(model.config)
        generator.encoder = PaddedEncoder(model)
        generator.decoder = PaddedDecoder(model)
        return generator

    def prepare_inputs_for_generation(
        self,
        input_ids,
        encoder_outputs=None,
        attention_mask=None,
        **kwargs,
    ):
        # Pad the inputs for Neuron
        current_length = input_ids.shape[1]
        pad_size = self.config.max_length - current_length
        return dict(
            input_ids=F.pad(input_ids, (0, pad_size)),
            attention_mask=attention_mask,
            encoder_outputs=encoder_outputs.last_hidden_state,
            current_length=torch.tensor(current_length - 1),
        )

    def get_encoder(self):
        def encode(input_ids, attention_mask, **kwargs):
            output, = self.encoder(input_ids, attention_mask)
            return BaseModelOutput(
                last_hidden_state=output,
            )
        return encode

    def forward(self, input_ids, attention_mask, encoder_outputs, current_length,
        ↪ **kwargs):
        logits = self.decoder(input_ids, attention_mask, encoder_outputs, current_length)
```

(continues on next page)

(continued from previous page)

```

        return Seq2SeqLMOutput(logits=logits)

    @property
    def device(self): # Attribute required by beam search
        return torch.device('cpu')

    def save_pretrained(self, directory):
        if os.path.isfile(directory):
            print(f"Provided path ({directory}) should be a directory, not a file")
            return
        os.makedirs(directory, exist_ok=True)
        torch.jit.save(self.encoder, os.path.join(directory, 'encoder.pt'))
        torch.jit.save(self.decoder, os.path.join(directory, 'decoder.pt'))
        self.config.save_pretrained(directory)

    @classmethod
    def from_pretrained(cls, directory):
        config = AutoConfig.from_pretrained(directory)
        obj = cls(config)
        obj.encoder = torch.jit.load(os.path.join(directory, 'encoder.pt'))
        obj.decoder = torch.jit.load(os.path.join(directory, 'decoder.pt'))
        setattr(obj.encoder, 'main_input_name', 'input_ids') # Attribute required by ↵
↵ beam search
        return obj

```

## Padded CPU Inference

To start, it is important to ensure that the transformations we have made to the model were successful. Using the classes defined above we can test that the padded model execution on CPU is identical to the original output also running on CPU.

```
[ ]: padded_model_cpu = PaddedGenerator.from_model(model_cpu)
infer(padded_model_cpu, tokenizer, sample_text)
```

## Padded Neuron Tracing & Inference

Now that the padded version of model is confirmed to produce the same outputs as the non-padded version, the model can be compiled for Neuron.

```
[ ]: import torch
import torch_neuron

def trace(model, num_texts, num_beams, max_decoder_length, max_encoder_length):
    """
    Traces the encoder and decoder modules for use on Neuron.

    This function fixes the network to the given sizes. Once the model has been

```

(continues on next page)

(continued from previous page)

compiled to a given size, the inputs to these networks must always be of fixed size.

Args:

```
model (PaddedGenerator): The padded generator to compile for Neuron
num_texts (int): The number of input texts to translate at once
num_beams (int): The number of beams to compute per text
max_decoder_length (int): The maximum number of tokens to be generated
max_encoder_length (int): The maximum number of input tokens that will be encoded
"""
```

```
# Trace the encoder
```

```
inputs = (
    torch.ones((num_texts, max_encoder_length), dtype=torch.long),
    torch.ones((num_texts, max_encoder_length), dtype=torch.long),
)
encoder = torch_neuron.trace(model.encoder, inputs)
```

```
# Trace the decoder (with expanded inputs)
```

```
batch_size = num_texts * num_beams
```

```
inputs = (
    torch.ones((batch_size, max_decoder_length), dtype=torch.long),
    torch.ones((batch_size, max_encoder_length), dtype=torch.long),
    torch.ones((batch_size, max_encoder_length, model.config.d_model), dtype=torch.
↳ float),
    torch.tensor(0),
)
decoder = torch_neuron.trace(model.decoder, inputs)
```

```
traced = PaddedGenerator(model.config)
```

```
traced.encoder = encoder
```

```
traced.decoder = decoder
```

```
setattr(encoder, 'main_input_name', 'input_ids') # Attribute required by beam search
```

```
return traced
```

```
[ ]: padded_model_neuron = trace(padded_model_cpu, num_texts, num_beams, max_decoder_length,
↳ max_encoder_length)
```

Comparing the Neuron execution to the original CPU implementation, you will see the exact same generated text.

```
[ ]: # CPU execution for comparison
infer(padded_model_neuron, tokenizer, sample_text)
```

## Padded Neuron Serialization

Finally, we can test that we can serialize and reload the model so that it can be used later in its precompiled format.

```
[ ]: padded_model_neuron.save_pretrained('NeuronPaddedMarianMT')
padded_model_loaded = PaddedGenerator.from_pretrained('NeuronPaddedMarianMT')
infer(padded_model_loaded, tokenizer, sample_text)
```

## Greedy Unrolled Model

An unrolled version of the model can achieve better performance in some cases since all operations will be executed on the Neuron hardware without returning to CPU. The consequence of this type of model is that since the generation loop execution never returns to CPU, the entire sequence up to `max_decoder_length` is performed in a single forward pass.

The following module performs greedy text generation. Unlike the original beam search text generation, this implementation always selects the most probable token and does not generate multiple result texts.

## GreedyUnrolledGenerator Module

```
[ ]: class GreedyUnrolledGenerator(torch.nn.Module):

    def __init__(self, model):
        super().__init__()
        self.config = model.config
        self.model = model

    def forward(self, input_ids, attention_mask):

        # Generate the encoder state for the input tokens. This is only done once and
        ↳ the state is reused.
        encoder_outputs, = self.model.model.encoder(input_ids, attention_mask=attention_
        ↳ mask, return_dict=False)

        # Set the initial state for the decode loop. This will grow per decoder iteration
        tokens = torch.full((input_ids.size(0), 2), self.config.decoder_start_token_id)

        # Iteratively invoke the decoder on incrementally generated `tokens` to generate
        ↳ a `next_token`.
        # Note that unlike the GeneratorMixin.generate function, there is no early-exit
        ↳ if the stop token
        # has been reached. This will always run a fixed number of iterations.
        for i in range(self.config.max_length):

            hidden, = self.model.model.decoder(
                input_ids=tokens,
                encoder_hidden_states=encoder_outputs,
                encoder_attention_mask=attention_mask,
                return_dict=False,
                use_cache=False,
            ) # size: [batch, current_length, vocab_size]
```

(continues on next page)



(continued from previous page)

```

        logits = F.linear(
            hidden[:, -1, :],
            self.model.model.shared.weight,
            bias=self.model.final_logits_bias
        )
        next_tokens = torch.argmax(logits, dim=1, keepdims=True)
        tokens = torch.cat([tokens, next_tokens], dim=1)

    return tokens

```

## Greedy CPU Inference

The inference code must be updated since the `generate` method is no longer used. This is because the entire generative inference loop occurs within the `GreedyUnrolledGenerator.forward` method.

```

[ ]: def infer_greedy(model, tokenizer, text):
    batch = tokenizer(text, max_length=max_decoder_length, truncation=True, padding='max_
    ↪length', return_tensors="pt")
    inputs = batch['input_ids'], batch['attention_mask']
    tokens = greedy_cpu(*inputs)
    print('Texts:')
    for i, t in enumerate(tokens):
        result = tokenizer.decode(t, skip_special_tokens=True)
        print(i + 1, result)

```

Like in previous section of this tutorial, first the greedy model is executed on CPU to validate that the correct results were produced. In this example, the generated text matches the first result of the original beam search.

```

[ ]: model_cpu.config.max_length = 8 # This controls the number of decoder loops. Reduced to 8
    ↪improve compilation speed.
    greedy_cpu = GreedyUnrolledGenerator(model_cpu)
    infer_greedy(greedy_cpu, tokenizer, sample_text)

```

## Greedy Neuron Tracing & Inference

Similarly the tracing is simplified since now the `GreedyUnrolledGenerator.forward` can be compiled as a single unit.

For compilation efficiency, two changes will be made compared to normal compilation:

- `torch.jit.freeze` is used because it can *sometimes* speed up compilation by in the case where a module is re-used multiple times. In this case, it is more efficient because the `self.model.model.decoder` is used in a loop.
- The `torch_neuron.trace` option `fallback` is set to `False`. This forces all operations to execute on Neuron. Most of the time this is not recommended or efficient. In this case, it is more efficient because it means a single subgraph is produced rather than many. Usually one subgraph would be produced per decoder iteration since `aten::embedding` is executed in a loop. The `aten::embedding` operation is otherwise executed on CPU by default since this is usually more efficient than executing on Neuron.

You may notice that compilation will take significantly longer with the unrolled model since the model inserts new operations into the compute graph for every single decoder iteration. This creates a much larger model graph even though the weights are re-used.

```
[ ]: example = (
    torch.ones((num_texts, max_encoder_length), dtype=torch.long),
    torch.ones((num_texts, max_encoder_length), dtype=torch.long),
)
greedy_cpu.eval()
greedy_trace = torch.jit.trace(greedy_cpu, example)
greedy_frozen = torch.jit.freeze(greedy_trace)
greedy_neuron = torch_neuron.trace(greedy_frozen, example, fallback=False)
```

```
[ ]: infer_greedy(greedy_neuron, tokenizer, sample_text)
```

### Greedy Neuron Serialization

Unlike the previous version of the model that used the `GenerationMixin` base class. This greedy version of the model can be serialized using the regular `torch.jit.save` and `torch.jit.load` utilities since it is a pure torchscript module.

```
[ ]: torch.jit.save(greedy_neuron, 'greedy_neuron.pt')
loaded_greedy_neuron = torch.jit.load('greedy_neuron.pt')
infer_greedy(loaded_greedy_neuron, tokenizer, sample_text)
```

## Appendix

### BART (Mask Filling Task)

These `PaddedGenerator` class can be applied to the BART model for the task of filling in mask tokens.

```
[ ]: from transformers import BartForConditionalGeneration, BartTokenizer
bart_name = "facebook/bart-large"
bart_model = BartForConditionalGeneration.from_pretrained(bart_name)
bart_model.config.max_length = max_decoder_length
bart_tokenizer = BartTokenizer.from_pretrained(bart_name)
bart_text = "UN Chief Says There Is No <mask> in Syria"
```

```
[ ]: # CPU Execution
infer(bart_model, bart_tokenizer, bart_text)
```

```
[ ]: # Neuron Execution
paddded_bart = PaddedGenerator.from_model(bart_model)
bart_neuron = trace(paddded_bart, num_texts, num_beams, max_decoder_length, max_encoder_
    ↪length)
infer(bart_neuron, bart_tokenizer, bart_text)
```

## Pegasus (Summarization Task)

These PaddedGenerator class can be applied to the Pegasus model for summarization.

```
[ ]: from transformers import PegasusForConditionalGeneration, PegasusTokenizer
    pegasus_name = 'google/pegasus-xsum'
    pegasus_model = PegasusForConditionalGeneration.from_pretrained(pegasus_name)
    pegasus_model.config.max_length = max_decoder_length
    pegasus_tokenizer = PegasusTokenizer.from_pretrained(pegasus_name)
    pegasus_text = "PG&E stated it scheduled the blackouts in response to forecasts for high_
    ↪ winds amid dry conditions. The aim is to reduce the risk of wildfires."

[ ]: # CPU Execution
    infer(pegasus_model, pegasus_tokenizer, pegasus_text)

[ ]: # Neuron Execution
    padded_pegasus = PaddedGenerator.from_model(pegasus_model)
    pegasus_neuron = trace(padded_pegasus, num_texts, num_beams, max_decoder_length, max_
    ↪ encoder_length)
    infer(pegasus_neuron, pegasus_tokenizer, pegasus_text)
```

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Utilizing Neuron Capabilities Tutorials

- BERT TorchServe tutorial [\[html\]](#)
- NeuronCore Pipeline tutorial [\[html\]](#) [\[notebook\]](#)

## Using NeuronCore Pipeline with PyTorch

In this tutorial you compile a pretrained BERT base model from HuggingFace Transformers, using the NeuronCore Pipeline feature of the AWS Neuron SDK. You benchmark model latency of the pipeline parallel mode and compare with the usual data parallel (multi-worker) deployment.

This tutorial is intended to run in an inf1.6xlarge, running the latest AWS Deep Learning AMI (DLAMI). The inf1.6xlarge instance size has AWS Inferentia chips for a total of 16 NeuronCores.

Verify that this Jupyter notebook is running the Python or Conda kernel environment that was set up according to the [PyTorch Installation Guide](#). You can select the kernel from the “Kernel -> Change Kernel” option on the top of this Jupyter notebook page.

**Note:** Do not execute this tutorial using “Run -> Run all cells” option.

## Install Dependencies:

This tutorial requires the following pip packages:

- torch-neuron
- neuron-cc[tensorflow]
- transformers

Most of these packages will be installed when configuring your environment using the Neuron PyTorch setup guide. The additional HuggingFace Transformers dependency must be installed here.

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
↪ detect
!pip install --upgrade "transformers==4.6.0"
```

## Compiling a BERT base model for a single NeuronCore

To run a HuggingFace `BERTModel` on Inferentia, you only need to add a single extra line of code to the usual Transformers PyTorch implementation, after importing the `torch_neuron` framework.

Add the argument `return_dict=False` to the BERT transformers model so it can be traced with `TorchScript`. TorchScript is a way to create serializable and optimizable models from PyTorch code.

Enable padding to a maximum sequence length of 128, to test the model's performance with a realistic payload size. You can adapt this sequence length to your application's requirement.

You can adapt the original example on the [BertModel forward pass docstring](#) according to the following cell

```
[ ]: import torch
import torch_neuron
from transformers import BertTokenizer, BertModel

from joblib import Parallel, delayed
import numpy as np
from tqdm import tqdm

import os
import time

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained('bert-base-uncased',return_dict=False)

inputs = tokenizer("Hello, my dog is cute",return_tensors="pt",max_length=128,padding=
↪ 'max_length',truncation=True)
```

The one extra line required is the call to `torch.neuron.trace()` method. This call compiles the model and returns the forward method of the `torch.nn.Model` method, which you can use to run inference.

The compiled graph can be saved using the `torch.jit.save` function and restored using `torch.jit.load` function for inference on Inf1 instances. During inference, the previously compiled artifacts will be loaded into the Neuron Runtime for inference execution.

```
[ ]: neuron_model = torch.neuron.trace(model,
                                     example_inputs = (inputs['input_ids'],inputs[
↳ 'attention_mask']),
                                     verbose=1)
```

### Running the BERT base model on a single NeuronCore

With the model already available in memory, you can time one execution and check for the latency on the single inference call. You will load the model into Inferentia with a single inference call. A large “wall time” is expected when you first run the next cell, running the cell twice will show the actual inference latency:

```
[ ]: %%time
# The following line tests inference and should be executed on Inf1 instance family.
outputs = neuron_model(*(inputs['input_ids'],inputs['attention_mask']))
```

You can also check for the throughput of the single model running on a single NeuronCore.

The sequential inference test (for loop) does not measure all the performance one can achieve in an instance with multiple NeuronCores. To improve hardware utilization you can run parallel inference requests over multiple model workers, which you’ll test in the Data Parallel Bonus Section below.

```
[ ]: %%time
for _ in tqdm(range(100)):
    outputs = neuron_model(*(inputs['input_ids'],inputs['attention_mask']))
```

Save the compiled model for later use:

```
[ ]: neuron_model.save('bert-base-uncased-neuron.pt')
```

### Compiling a BERT base model for 16 NeuronCores

Our next step is to compile the same model for all 16 NeuronCores available in the inf1.6xlarge and check the performance difference when running pipeline parallel inferences..

```
[ ]: import torch
import torch_neuron
from transformers import BertTokenizer, BertModel

from joblib import Parallel, delayed
import numpy as np
from tqdm import tqdm

import os
import time

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained('bert-base-uncased',return_dict=False)

inputs = tokenizer("Hello, my dog is cute",return_tensors="pt",max_length=128,padding=
```

(continues on next page)

(continued from previous page)

```
↪ 'max_length', truncation=True)
```

To enable pipeline mode during compilation, you need only to add the compiler flag `--neuroncore-pipeline-cores` and set the number of desired cores. The cell below sets up a `neuroncore_pipeline_cores` string, which you can set for the available number of NeuronCores on the instance: *inf1.6xlarge* has 16 NeuronCores in 4 Inferentia chips.

```
[ ]: # Number of Cores in the Pipeline Mode
neuroncore_pipeline_cores = 16 # This string should be '4' on an inf1.xlarge

# Compiling for neuroncore-pipeline-cores='16'
neuron_pipeline_model = torch.neuron.trace(model,
                                           example_inputs = (inputs['input_ids'], inputs[
↪ 'attention_mask']),
                                           verbose=1,
                                           compiler_args = ['--neuroncore-pipeline-cores
↪ ', str(neuroncore_pipeline_cores)]
                                           )
```

## Running the BERT base model on 16 NeuronCores

Next, time one execution and check for the latency on the single inference call over 16 cores. You will load the model into Inferentia with a single inference call. A large “wall time” is expected when you first run the next cell, running the cell twice will show the actual inference latency:

```
[ ]: %%time
# The following line tests inference and should be executed on Inf1 instance family.
outputs = neuron_pipeline_model(*(inputs['input_ids'], inputs['attention_mask']))
```

Check also for the throughput of the single model running over a 16 NeuronCores.

The sequential inference test (for loop) does not measure all the performance one can achieve with Pipeline mode. As the inference runs in streaming fashion, at least 15 cores are waiting for a new call until the last one processes the first call. This results in low NeuronCore utilization. To improve hardware utilization you will require parallel inference requests, which you’ll test in the next section.

```
[ ]: for _ in tqdm(range(100)):
      outputs = neuron_pipeline_model(*(inputs['input_ids'], inputs['attention_mask']))
```

## Load Testing the Pipeline Parallel Mode

To put the 16 NeuronCores group to test, a client has to run concurrent requests to the model. In this Notebook setup you achieve it by creating a thread pool with `Joblib.Parallel`, with all workers on the pool running one inference call.

You can define a new method called `inference_latency()` so that you measure the amount of time each inference calls take.

```
[ ]: def inference_latency(model, *inputs):
      """
      infetence_time is a simple method to return the latency of a model inference.
```

(continues on next page)

(continued from previous page)

```

Parameters:
    model: torch model onbject loaded using torch.jit.load
    inputs: model() args

Returns:
    latency in seconds
"""
start = time.time()
_ = model(*inputs)
return time.time() - start

```

Use `tqdm` to measure total throughput of your experiment, with a nice side-effect of “cool progress bar!”. The total throughput is expected to be high, so set your experiment range to a large number, here 30k inferences.

To calculate the latency statistics over the returned 30k list of latencies use `numpy.quantile()` method.

```

[ ]: t = tqdm(range(30000), position=0, leave=True)
latency = Parallel(n_jobs=12,prefer="threads")(delayed(inference_latency)(neuron_
→pipeline_model,*(inputs['input_ids'],inputs['attention_mask']))) for i in t

p50 = np.quantile(latency[-10000:],0.50) * 1000
p95 = np.quantile(latency[-10000:],0.95) * 1000
p99 = np.quantile(latency[-10000:],0.99) * 1000
avg_throughput = t.total/t.format_dict['elapsed']
print(f'Avg Throughput: :{avg_throughput:.1f}')
print(f'50th Percentile Latency:{p50:.1f} ms')
print(f'95th Percentile Latency:{p95:.1f} ms')
print(f'99th Percentile Latency:{p99:.1f} ms')

```

Save compile model for later use:

```

[ ]: # Save the TorchScript graph
neuron_pipeline_model.save('bert-base-uncased-neuron-pipeline.pt')

```

## Bonus Section - Load Testing Data Parallel Mode

```

[ ]: import torch
import torch_neuron
from transformers import BertTokenizer

from joblib import Parallel, delayed
import numpy as np
from tqdm import tqdm

import os
import time

def inference_latency(model,*inputs):
    """
    infetence_time is a simple method to return the latency of a model inference.

```

(continues on next page)

(continued from previous page)

```

    Parameters:
        model: torch model object loaded using torch.jit.load
        inputs: model() args

    Returns:
        latency in seconds
    """
    start = time.time()
    _ = model(*inputs)
    return time.time() - start

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

inputs = tokenizer("Hello, my dog is cute",return_tensors="pt",max_length=128,padding=
↳ 'max_length',truncation=True)

```

You use the 'NEURON\_RT\_NUM\_CORES' environment variable to define how many Neuron cores to be used. Set the environment variable to the number of individual workers you want to test in parallel.

`torch_neuron` will load one model per NeuronCore group until it runs out of cores. At that point, if the Python process continues to spawn more model object using `torch.jit.load`, `torch_neuron` will start stacking more than one model per core, until the Inferentia chip memory is full.

Inferentia is able to run inference over all the loaded models, but only one at a time. The Neuron Runtime takes care of dynamically switching the model context as requests come in, no extra worker process management required. Use 1 model per NeuronCore to achieve maximum performance.

The following cell creates a list with as many models as NeuronCore Groups and execute one single dummy inference to load the models into Inferentia.

```

[ ]: import warnings
    # Number of data parallel workers
    number_of_workers=16 # This number should be 4 on an inf1.xlarge

    # Setting up a data parallel group
    os.environ['NEURON_RT_NUM_CORES'] = str(number_of_workers)

    # Loading 'number_of_workers' amount of models in Python memory
    model_list = [torch.jit.load('bert-base-uncased-neuron.pt') for _ in range(number_of_
↳ workers)]

    # Dummy inference to load models to Inferentia
    _ = [mod(*(inputs['input_ids'],inputs['attention_mask'])) for mod in model_list]

```

Adapt the call to `joblib.Parallel()` iterating over a concatenated version of the `model_list`, to run 'round-robin' calls to each of the model workers.

```

[ ]: t = tqdm(model_list*1500,position=0, leave=True)
    latency = Parallel(n_jobs=number_of_workers,prefer="threads")(delayed(inference_
↳ latency)(mod,*(inputs['input_ids'],inputs['attention_mask'])) for mod in t)

```

(continues on next page)



(continued from previous page)

```

p50 = np.quantile(latency[-10000:],0.50) * 1000
p95 = np.quantile(latency[-10000:],0.95) * 1000
p99 = np.quantile(latency[-10000:],0.99) * 1000
avg_throughput = t.total/t.format_dict['elapsed']
print(f'Avg Throughput: :{avg_throughput:.1f}')
print(f'50th Percentile Latency:{p50:.1f} ms')
print(f'95th Percentile Latency:{p95:.1f} ms')
print(f'99th Percentile Latency:{p99:.1f} ms')

```

For this model, despite the larger number of workers, the per-worker latency increases when running a single model per core, which in turn reduces the total throughput.

This behavior may not repeat if the model memory footprint or the input payload size changes, i.e batch size > 1. We encourage you to experiment with the data parallel and pipeline parallel modes to optimize your application performance.

*This document is relevant for:* Inf1

## Computer Vision Tutorials

- ResNet-50 tutorial [\[html\]](#) [\[notebook\]](#)
- PyTorch YOLOv4 tutorial [\[html\]](#) [\[notebook\]](#)

## Natural Language Processing (NLP) Tutorials

- HuggingFace pretrained BERT tutorial [\[html\]](#) [\[notebook\]](#)
- HuggingFace pretrained BERT tutorial with shared weights [\[html\]](#) [\[notebook\]](#)
- Bring your own HuggingFace pretrained BERT container to Sagemaker Tutorial [\[html\]](#) [\[notebook\]](#)
- LibTorch C++ tutorial [\[html\]](#)
- TorchServe tutorial [\[html\]](#)
- HuggingFace MarianMT tutorial [\[html\]](#) [\[notebook\]](#)

## Utilizing Neuron Capabilities Tutorials

- BERT TorchServe tutorial [\[html\]](#)
- NeuronCore Pipeline tutorial [\[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Additional Examples (torch-neuron)

- [AWS Neuron Samples GitHub Repository](#)

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## API Reference Guide (torch-neuron)

*This document is relevant for: Inf1*

### PyTorch-Neuron trace python API

The PyTorch-Neuron trace Python API provides a method to generate PyTorch models for execution on Inferentia, which can be serialized as TorchScript. It is analogous to `torch.jit.trace()` function in PyTorch.

`torch_neuron.trace(model, example_inputs, **kwargs)`

The `torch_neuron.trace()` method sends operations to the Neuron-Compiler (`neuron-cc`) for compilation and embeds compiled artifacts in a TorchScript graph.

Compilation can be done on any EC2 machine with sufficient memory and compute resources. c5.4xlarge or larger is recommended.

Options can be passed to Neuron compiler via the compile function. See [Neuron compiler CLI Reference Guide \(neuron-cc\)](#) for more information about compiler options.

This function partitions nodes into operations that are supported by Neuron and operations which are not. Operations which are not supported by Neuron are run on CPU. Graph partitioning can be controlled by the `subgraph_builder_function`, `minimum_segment_size`, and `fallback` parameters (See below). By default all supported operations are compiled and run on Neuron.

The compiled graph can be saved using the `torch.jit.save()` function and restored using `torch.jit.load()` function for inference on Inf1 instances. During inference, the previously compiled artifacts will be loaded into the Neuron Runtime for inference execution.

#### *Required Arguments*

##### **Parameters**

- **model** (*Module, callable*) – The functions that that will be run with `example_inputs` arguments. The arguments and return types must compatible with `torch.jit.trace()`. When a Module is passed to `torch_neuron.trace()`, only the `forward()` method is run and traced.
- **example\_inputs** (*tuple*) – A tuple of example inputs that will be passed to the `model` while tracing. The resulting trace can be run with inputs of different types and shapes assuming the traced operations support those types and shapes. This parameter may also be a single `torch.Tensor` in which case it is automatically wrapped in a tuple.

#### *Optional Keyword Arguments*

##### **Keyword Arguments**

- **compiler\_args** (*list[str]*) – List of strings representing `neuron-cc` compiler arguments. Note that these arguments apply to all subgraphs generated by allowlist partitioning. For example, use `compiler_args=['--neuroncore-pipeline-cores', '4']` to set number of NeuronCores per subgraph to 4. See [Neuron compiler CLI Reference Guide \(neuron-cc\)](#) for more information about compiler options.

- **compiler\_timeout** (*int*) – Timeout in seconds for waiting neuron-cc to complete. Exceeding this timeout will cause a `subprocess.TimeoutExpired` exception.
- **compiler\_workdir** (*str*) – Work directory used by neuron-cc. Useful for debugging and/or inspecting neuron-cc logs/IRs.
- **subgraph\_builder\_function** (*callable*) – A function which is evaluated on each node during graph partitioning. This takes in a torch graph operator node and returns a `bool` value of whether it should be included in the fused Neuron graph or not. By default the partitioner selects all operators which are supported by Neuron.
- **minimum\_segment\_size** (*int*) – A parameter used during partitioning. This specifies the minimum number of graph nodes which should be compiled into a Neuron graph (default=2). If the number of nodes is smaller than this size, the operations will run on CPU.
- **single\_fusion\_ratio\_threshold** (*float*) – A parameter used during partitioning. During partitioning, if a single partition contains a fraction of operations greater than this threshold, only one graph partition will be compiled (default=0.6). This is used to avoid compiling many small Neuron graphs. To force compilation of all graphs to Neuron (even when they are very small), a value of 1.0 can be used.
- **fallback** (*bool*) – A function parameter to turn off graph partitioning. Indicates whether to attempt to fall back to CPU operations if an operation is not supported by Neuron. By default this is True. If this is set to False and an operation is not supported by Neuron, this will fail compilation and raise an `AttributeError`.
- **dynamic\_batch\_size** (*bool*) – A flag to allow Neuron graphs to consume variable sized batches of data. Dynamic sizing is restricted to the 0th dimension of a tensor.
- **optimizations** (*list*) – A list of `Optimization` passes to apply to the model.
- **separate\_weights** (*bool*) – A flag to enable compilation of models with over 1.9GB of constant parameters. By default this flag is False. If this is set to True and the compiler version is not new enough to support the flag, this will raise an `NotImplementedError`.
- **\*\*kwargs** – All other keyword arguments will be forwarded directly to `torch.jit.trace()`. This supports flags like `strict=False` in order to allow dictionary outputs.

### Returns

The traced `ScriptModule` with embedded compiled neuron sub-graphs. Operations in this module will run on Neuron unless they are not supported by Neuron or manually partitioned to run on CPU.

Note that in `torch<1.8` This would return a `ScriptFunction` if the input was function type.

### Return type

*ScriptModule, ScriptFunction*

### `class torch_neuron.Optimization`

A set of optimization passes that can be applied to the model.

#### `FLOAT32_TO_FLOAT16`

A post-processing pass that converts all `torch.float32` tensors to `torch.float16` tensors. The advantage to this optimization pass is that input/output tensors will be type cast. This reduces the amount of data that will be copied to and from Inferentia hardware. The resulting traced model will accept both `torch.float32` and `torch.float16` inputs where the model used `torch.float32` inputs during tracing. It is only beneficial to enable this optimization if the throughput of a model is highly dependent upon data transfer speed. This optimization is not recommended if the final application will use `torch.float32` inputs since the `torch.float16` type cast will occur on CPU during inference.

## Example Usage

### Function Compilation

```
import torch
import torch_neuron

def foo(x, y):
    return 2 * x + y

# Run `foo` with the provided inputs and record the tensor operations
traced_foo = torch.neuron.trace(foo, (torch.rand(3), torch.rand(3)))

# `traced_foo` can now be run with the TorchScript interpreter or saved
# and loaded in a Python-free environment
torch.jit.save(traced_foo, 'foo.pt')
traced_foo = torch.jit.load('foo.pt')
```

### Module Compilation

```
import torch
import torch_neuron
import torch.nn as nn

class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        return self.conv(x) + 1

n = Net()
n.eval()

inputs = torch.rand(1, 1, 3, 3)

# Trace a specific method and construct `ScriptModule` with
# a single `forward` method
neuron_forward = torch.neuron.trace(n.forward, inputs)

# Trace a module (implicitly traces `forward`) and constructs a
# `ScriptModule` with a single `forward` method
neuron_net = torch.neuron.trace(n, inputs)
```

## Pre-Trained Model Compilation

The following is an example usage of the compilation Python API, with default compilation arguments, using a pre-trained `torch.nn.Module`:

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)
```

## Compiling models with `torch.jit.trace` kwargs

This example uses the `strict=False` flag to compile a model with dictionary outputs. Similarly, any other keyword argument of `torch.jit.trace()` can be passed directly to `torch_neuron.trace()` so that it is passed to the underlying trace call.

```
import torch
import torch_neuron
import torch.nn as nn

class Model(nn.Module):
    def __init__(self):
        super(Model, self).__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        return {'conv': self.conv(x) + 1}

model = Model()
model.eval()

inputs = torch.rand(1, 1, 3, 3)

# use the strict=False kwarg to compile a model with dictionary outputs
# the model output format does not change
model_neuron = torch.neuron.trace(model, inputs, strict=False)
```

## Dynamic Batching

This example uses the optional `dynamic_batch_size` option in order to support variable sized batches at inference time.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input of batch size 1
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image, dynamic_batch_size=True)

# Execute with a batch of 7 images
batch = torch.rand([7, 3, 224, 224])
results = model_neuron(batch)
```

## Manual Partitioning

The following example uses the optional `subgraph_builder_function` parameter to ensure that only a specific convolution layer is compiled to Neuron. The remaining operations are executed on CPU.

```
import torch
import torch_neuron
import torch.nn as nn

class ExampleConvolutionLayer(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        return self.conv(x) + 1

class Model(nn.Module):
    def __init__(self):
        super().__init__()
        self.layer = ExampleConvolutionLayer()

    def forward(self, x):
        return self.layer(x) * 100

def subgraph_builder_function(node) -> bool:
    """Select if the node will be included in the Neuron graph"""

    # Node names are tuples of Module names.
    if 'ExampleConvolutionLayer' in node.name:
        return True
```

(continues on next page)

(continued from previous page)

```

    # Ignore all operations not in the example convolution layer
    return False

model = Model()
model.eval()

inputs = torch.rand(1, 1, 3, 3)

# Log output shows that `aten::_convolution` and `aten::add` are compiled
# but `aten::mul` is not. This will seamlessly switch between Neuron/CPU
# execution in a single graph.
neuron_model = torch_neuron.trace(
    model,
    inputs,
    subgraph_builder_function=subgraph_builder_function
)

```

## Separate Weights

This example uses the optional `separate_weights` option in order to support compilation of models greater than 1.9GB.

```

import torch
import torch_neuron
from torchvision import models

# Load the model
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
#the models' output format does not change
model_neuron = torch.neuron.trace(model, image, separate_weights=True)

```

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## torch.neuron.DataParallel API

The `torch.neuron.DataParallel()` Python API implements data parallelism on `ScriptModule` models created by the *PyTorch-Neuron trace python API*. This function is analogous to `DataParallel` in PyTorch. The *Data Parallel Inference on Torch Neuron* application note provides an overview of how `torch.neuron.DataParallel()` can be used to improve the performance of inference workloads on Inferentia.

`torch.neuron.DataParallel(model, device_ids=None, dim=0)`

Applies data parallelism by replicating the model on available NeuronCores and distributing data across the different NeuronCores for parallelized inference.

By default, DataParallel will use all available NeuronCores allocated for the current process for parallelism. DataParallel will apply parallelism on `dim=0` if `dim` is not specified.

DataParallel automatically enables *dynamic batching* on eligible models if `dim=0`. Dynamic batching can be disabled using `torch.neuron.DataParallel.disable_dynamic_batching()`. If dynamic batching is not enabled, the batch size at compilation-time must be equal to the batch size at inference-time divided by the number of NeuronCores being used. Specifically, the following must be true when dynamic batching is disabled: `input.shape[dim] / len(device_ids) == compilation_input.shape[dim]`. DataParallel will throw a warning if dynamic batching cannot be enabled.

DataParallel will try load all of a model's NEFFs onto a single NeuronCore, only if all of the NEFFs can fit on a single NeuronCore. DataParallel does not currently support models that have been compiled with *NeuronCore Pipeline*.

`torch.neuron.DataParallel()` requires PyTorch  $\geq 1.8$ .

#### Required Arguments

##### Parameters

**model** (*ScriptModule*) – Model created by the *PyTorch-Neuron trace python API* to be parallelized.

#### Optional Arguments

##### Parameters

- **device\_ids** (*list*) – List of `int` or `'nc:#'` that specify the NeuronCores to use for parallelization (default: all NeuronCores). Refer to the *device\_ids note* for a description of how `device_ids` indexing works.
- **dim** (*int*) – Dimension along which the input tensor is scattered across NeuronCores (default `dim=0`).

#### Attributes

##### Parameters

- **num\_workers** (*int*) – Number of worker threads used for multithreaded inference (default: `2 * number of NeuronCores`).
- **split\_size** (*int*) – Size of the input chunks (default: `max(1, input.shape[dim] // number of NeuronCores)`).

`torch.neuron.DataParallel.disable_dynamic_batching()`

Disables automatic dynamic batching on the DataParallel module. See *Dynamic batching disabled* for example of how DataParallel can be used with dynamic batching disabled. Use as follows:

```
>>> model_parallel = torch.neuron.DataParallel(model_neuron)
>>> model_parallel.disable_dynamic_batching()
```

**Note:** `device_ids` uses per-process NeuronCore granularity and zero-based indexing. Per-process granularity means that each Python process “sees” its own view of the world. Specifically, this means that `device_ids` only “sees” the NeuronCores that are allocated for the current process. Zero-based indexing means that each Python process will index its allocated NeuronCores starting at 0, regardless of the “global” index of the NeuronCores. Zero-based indexing makes it possible to redeploy the exact same code unchanged in different process. This behavior is analogous to the `device_ids` argument in the PyTorch `DataParallel` function.

As an example, assume DataParallel is run on an `inf1.6xlarge`, which contains four Inferentia chips each of which contains four NeuronCores:



- If `NEURON_RT_VISIBLE_CORES` is not set, a single process can access all 16 NeuronCores. Thus specifying `device_ids=["nc:0"]` will correspond to `chip0:core0` and `device_ids=["nc:14"]` will correspond to `chip3:core2`.
- However, if two processes are launched where: process 1 has `NEURON_RT_VISIBLE_CORES=0-6` and process 2 has `NEURON_RT_VISIBLE_CORES=7-15`, `device_ids=["nc:14"]` cannot be specified in either process. Instead, `chip3:core2` can only be accessed in process 2. Additionally, `chip3:core2` is specified in process 2 with `device_ids=["nc:7"]`. Furthermore, in process 1, `device_ids=["nc:0"]` would correspond to `chip0:core0`; in process 2 `device_ids=["nc:0"]` would correspond to `chip1:core3`.

## Examples

The following sections provide example usages of the `torch.neuron.DataParallel()` module.

### Default usage

The default `DataParallel` use mode will replicate the model on all available NeuronCores in the current process. The inputs will be split on `dim=0`.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module
model_parallel = torch.neuron.DataParallel(model_neuron)

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

## Specifying NeuronCores

The following example uses the `device_ids` argument to use the first three NeuronCores for DataParallel inference.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module, run on the first three NeuronCores
# Equivalent to model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2])
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=['nc:0', 'nc:1', 'nc:2'])

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

## DataParallel with dim != 0

In this example we run DataParallel inference using four NeuronCores and `dim = 2`. Because `dim != 0`, dynamic batching is not enabled. Consequently, the DataParallel inference-time batch size must be four times the compile-time batch size. DataParallel will generate a warning that dynamic batching is disabled because `dim != 0`.

```
import torch
import torch_neuron

# Create an example model
class Model(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = torch.nn.Conv2d(3, 3, 3)

    def forward(self, x):
        return self.conv(x) + 1

model = Model()
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 8, 8])
```

(continues on next page)

(continued from previous page)

```

model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module using 4 NeuronCores and dim = 2
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2, 3], dim=2)

# Create a batched input
# Note that image_batched.shape[dim] / len(device_ids) == image.shape[dim]
batch_size = 4 * 8
image_batched = torch.rand([1, 3, batch_size, 8])

# Run inference with a batched input
output = model_parallel(image_batched)

```

## Dynamic batching

In the following example, we use the `torch.neuron.DataParallel()` module to run inference using several different batch sizes without recompiling the Neuron model.

```

import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module
model_parallel = torch.neuron.DataParallel(model_neuron)

# Create batched inputs and run inference on the same model
batch_sizes = [2, 3, 4, 5, 6]
for batch_size in batch_sizes:
    image_batched = torch.rand([batch_size, 3, 224, 224])

    # Run inference with a batched input
    output = model_parallel(image_batched)

```

## Dynamic batching disabled

In the following example, we use `torch.neuron.DataParallel.disable_dynamic_batching()` to disable dynamic batching. We provide an example of a batch size that will not work when dynamic batching is disabled as well as an example of a batch size that does work when dynamic batching is disabled.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module and use 4 NeuronCores
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2, 3], dim=0)

# Disable dynamic batching
model_parallel.disable_dynamic_batching()

# Create a batched input (this won't work)
batch_size = 8
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will fail because dynamic batching is disabled and
# image_batched.shape[dim] / len(device_ids) != image.shape[dim]
# output = model_parallel(image_batched)

# Create a batched input (this will work)
batch_size = 4
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will work because
# image_batched.shape[dim] / len(device_ids) == image.shape[dim]
output = model_parallel(image_batched)
```

## Full tutorial with torch.neuron.DataParallel

For an end-to-end tutorial that uses DataParallel, see the PyTorch Resnet Tutorial.

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## PyTorch Neuron (torch-neuron) Core Placement API [Beta]

**Warning:** The following functionality is beta and **will not be supported** in future releases of the Neuron SDK. This module serves only as a preview for future functionality. In future releases, equivalent functionality may be moved directly to the `torch_neuron` module and will no longer be available in the `torch_neuron.experimental` module.

Functions which enable placement of `torch.jit.ScriptModule` to specific NeuronCores. Two sets of functions are provided which can be used interchangeably but have different performance characteristics and advantages:

- The `multicore_context()` & `neuron_cores_context()` functions are context managers that allow a model to be placed on a given NeuronCore at `torch.jit.load()` time. These functions are the most efficient way of loading a model since the model is loaded directly to a NeuronCore. The alternative functions described below require that a model is unloaded from one core and then reloaded to another.
- The `set_multicore()` & `set_neuron_cores()` functions allow a model that has already been loaded to a NeuronCore to be moved to a different NeuronCore. This functionality is less efficient than directly loading a model to a NeuronCore within a context manager but allows device placement to be fully dynamic at runtime. This is analogous to the `torch.nn.Module.to()` function for device placement.

**Important:** A prerequisite to enable placement functionality is that the loaded `torch.jit.ScriptModule` has already been compiled with the `torch_neuron.trace()` API. Attempting to place a regular `torch.nn.Module` onto a NeuronCore prior to compilation will do nothing.

### `torch_neuron.experimental.multicore_context()`

A context which loads all Neuron subgraphs to all visible NeuronCores.

This loads each Neuron subgraph within a `torch.jit.ScriptModule` to multiple NeuronCores without requiring multiple calls to `torch.jit.load()`. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores for concurrent threadsafe inferences. Executions use a round-robin strategy to distribute across NeuronCores.

Any calls to `torch.jit.load()` will cause any underlying Neuron subgraphs to load to the specified NeuronCores within this context. This context manager only needs to be used during the model load. After loading, inferences do not need to occur in this context in order to use the correct NeuronCores.

Note that this context is *not* threadsafe. Using multiple core placement contexts from multiple threads may not correctly place models.

#### Raises

`RuntimeError` – If the Neuron runtime cannot be initialized.

### Examples

*Multiple Core Replication:* Directly load a model to all visible NeuronCores. This allows a single `torch.jit.ScriptModule` to use all NeuronCores by running round-robin executions.

```
>>> with torch_neuron.experimental.multicore_context():
>>>     model = torch.jit.load('example_neuron_model.pt')
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 1
>>> model(example) # Executes on NeuronCore 2
```

`torch_neuron.experimental.neuron_cores_context(start_nc: int = -1, nc_count: int = -1)`

A context which sets the NeuronCore start/count for all Neuron subgraphs.

Any calls to `torch.jit.load()` will cause any underlying Neuron subgraphs to load to the specified NeuronCores within this context. This context manager only needs to be used during the model load. After loading, inferences do not need to occur in this context in order to use the correct NeuronCores.

Note that this context is *not* threadsafe. Using multiple core placement contexts from multiple threads may not correctly place models.

#### Parameters

- **start\_nc** – The starting NeuronCore index where the Module is placed. The value -1 automatically loads to the optimal NeuronCore (least used). Note that this index is always relative to NeuronCores visible to this process.
- **nc\_count** – The number of NeuronCores to use. The value -1 will load a model to exactly the number of cores required by that model (1 for most models, >1 when using NeuronCore Pipeline). If `nc_count` is greater than the number of NeuronCores required by the model, the model will be replicated across multiple NeuronCores. (`replications = floor(nc_count / cores_per_model)`)

#### Raises

- **RuntimeError** – If the Neuron runtime cannot be initialized.
- **ValueError** – If the `nc_count` is an invalid number of NeuronCores.

#### Examples

*Single Load:* Directly load a model from disk to the first visible NeuronCore.

```
>>> with torch_neuron.experimental.neuron_cores_context(start_nc=0, nc_count=1):
>>>     model = torch.jit.load('example_neuron_model.pt')
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 0
```

*Multiple Core Replication:* Directly load a model from disk to 2 NeuronCores. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores by running round-robin executions.

```
>>> with torch_neuron.experimental.neuron_cores_context(start_nc=2, nc_count=2):
>>>     model = torch.jit.load('example_neuron_model.pt')
>>> model(example) # Executes on NeuronCore 2
>>> model(example) # Executes on NeuronCore 3
>>> model(example) # Executes on NeuronCore 2
```

*Multiple Model Load:* Directly load 2 models from disk and pin them to separate NeuronCores. This causes each `torch.jit.ScriptModule` to always execute on a specific NeuronCore.

```
>>> with torch_neuron.experimental.neuron_cores_context(start_nc=2):
>>>     model1 = torch.jit.load('example_neuron_model.pt')
>>> with torch_neuron.experimental.neuron_cores_context(start_nc=0):
>>>     model2 = torch.jit.load('example_neuron_model.pt')
>>> model1(example) # Executes on NeuronCore 2
>>> model1(example) # Executes on NeuronCore 2
```

(continues on next page)

(continued from previous page)

```
>>> model2(example) # Executes on NeuronCore 0
>>> model2(example) # Executes on NeuronCore 0
```

`torch_neuron.experimental.set_multicore(trace: torch.jit.ScriptModule)`

Loads all Neuron subgraphs in a torch Module to all visible NeuronCores.

This loads each Neuron subgraph within a `torch.jit.ScriptModule` to multiple NeuronCores without requiring multiple calls to `torch.jit.load()`. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores for concurrent threadsafe inferences. Executions use a round-robin strategy to distribute across NeuronCores.

This will unload the model from an existing NeuronCore if it is already loaded.

Requires Torch 1.8+

#### Parameters

**trace** – A torch module which contains one or more Neuron subgraphs.

#### Raises

**RuntimeError** – If the Neuron runtime cannot be initialized.

### Examples

*Multiple Core Replication:* Move a model across all visible NeuronCores after loading. This allows a single `torch.jit.ScriptModule` to use all NeuronCores by running round-robin executions.

```
>>> model = torch.jit.load('example_neuron_model.pt')
>>> torch_neuron.experimental.set_multicore(model)
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 1
>>> model(example) # Executes on NeuronCore 2
```

`torch_neuron.experimental.set_neuron_cores(trace: torch.jit.ScriptModule, start_nc: int = -1, nc_count: int = -1)`

Set the NeuronCore start/count for all Neuron subgraphs in a torch Module.

This will unload the model from an existing NeuronCore if it is already loaded.

Requires Torch 1.8+

#### Parameters

- **trace** – A torch module which contains one or more Neuron subgraphs.
- **start\_nc** – The starting NeuronCore index where the Module is placed. The value -1 automatically loads to the optimal NeuronCore (least used). Note that this index is always relative to NeuronCores visible to this process.
- **nc\_count** – The number of NeuronCores to use. The value -1 will load a model to exactly the number of cores required by that model (1 for most models, >1 when using NeuronCore Pipeline). If `nc_count` is greater than the number of NeuronCores required by the model, the model will be replicated across multiple NeuronCores. (`replications = floor(nc_count / cores_per_model)`)

#### Raises

- **RuntimeError** – If the Neuron runtime cannot be initialized.
- **ValueError** – If the `nc_count` is an invalid number of NeuronCores.

## Examples

*Single Load:* Move a model to the first visible NeuronCore after loading.

```
>>> model = torch.jit.load('example_neuron_model.pt')
>>> torch_neuron.experimental.set_neuron_cores(model, start_nc=0, nc_count=1)
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 0
>>> model(example) # Executes on NeuronCore 0
```

*Multiple Core Replication:* Replicate a model to 2 NeuronCores after loading. This allows a single `torch.jit.ScriptModule` to use multiple NeuronCores by running round-robin executions.

```
>>> model = torch.jit.load('example_neuron_model.pt')
>>> torch_neuron.experimental.set_neuron_cores(model, start_nc=2, nc_count=2)
>>> model(example) # Executes on NeuronCore 2
>>> model(example) # Executes on NeuronCore 3
>>> model(example) # Executes on NeuronCore 2
```

*Multiple Model Load:* Move and pin 2 models to separate NeuronCores. This causes each `torch.jit.ScriptModule` to always execute on a specific NeuronCore.

```
>>> model1 = torch.jit.load('example_neuron_model.pt')
>>> torch_neuron.experimental.set_neuron_cores(model1, start_nc=2)
>>> model2 = torch.jit.load('example_neuron_model.pt')
>>> torch_neuron.experimental.set_neuron_cores(model2, start_nc=0)
>>> model1(example) # Executes on NeuronCore 2
>>> model1(example) # Executes on NeuronCore 2
>>> model2(example) # Executes on NeuronCore 0
>>> model2(example) # Executes on NeuronCore 0
```

*This document is relevant for:* Inf1

- *PyTorch Neuron trace Python API*
- *torch.neuron.DataParallel API*
- *PyTorch Neuron (torch-neuron) Core Placement API [Beta]*

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Developer Guide (torch-neuron)

*This document is relevant for:* Inf1



## Running inference on variable input shapes with bucketing

### Table of contents

- *Introduction*
- *Applications that benefit from bucketing*
- *Implementing bucketing*
  - *Creating bucketed models*
  - *Running inference with bucketing*
- *Examples*
  - *Computer vision bucketing*
  - *End-to-end computer vision bucketing example*
  - *Natural language processing bucketing*
  - *End-to-end natural language processing bucketing example*

### Introduction

With Inferentia, the shape of every input must be fixed at compile time. For applications that require multiple input sizes, we recommend using padding or bucketing techniques. Padding requires you to compile your model with the largest expected input size and pad every input to this maximum size. If the performance of your model using padding is not within your targets, you can consider implementing bucketing.

This guide introduces bucketing, a technique to run inference on inputs with variable shapes on Inferentia. The following sections explain how bucketing can improve the performance of inference workloads on Inferentia. It covers an overview of how bucketing works and provides examples of using bucketing in *computer vision* and *natural language processing* applications.

### Applications that benefit from bucketing

Bucketing refers to compiling your model multiple times with different target input shapes to create “bucketed models.” *Creating bucketed models* provides an overview on selecting the input shapes that you use to create bucketed models. At inference time, each input is padded until its shape matches the next largest bucket shape. The padded input is then passed into the corresponding bucketed model for inference. By compiling the same model with multiple different input shapes, the amount of input padding is reduced compared to padding every input to the maximum size in your dataset. This technique minimizes the compute overhead and improves inference performance compared to padding every image to the maximum shape in your dataset.

Bucketing works best when multiple different bucketed models are created to efficiently cover the full range of input shapes. You can fine-tune the model performance by experimenting with different bucket sizes that correspond to the distribution of input shapes in your dataset.

Bucketing can only be used if there is an upper bound on the shape of the inputs. If necessary, an upper bound on the input shape can be enforced using resizing and other forms of preprocessing.

The upper bound on the number of bucketed models that you use is dictated by the total size of the compiled bucketed models. Each Inferentia chip has 8GB of DRAM, or 2GB of DRAM per NeuronCore. An inf1.xlarge and inf1.2xlarge have 1 Inferentia chip, an inf1.6xlarge has 4 Inferentia chips, and an inf1.24xlarge has 16 Inferentia chips. Thus,

you should limit the total size of all bucketed models to around 8GB per Inferentia chip or 2GB per NeuronCore. The following formula provides an approximation for the number of compiled bucketed models you can fit on each NeuronCore:

```
number-of-buckets = round(10^9 / number-of-weights-in-model)
```

We recommend using *neuron-top* to monitor the memory usage on your inf1 instance as you load multiple bucketed models.

## Implementing bucketing

Implementing bucketing consists of two main parts: creating multiple bucketed models at compile-time and running inference using the bucketed models on (padded) inputs. The following sections describe how to implement bucketing to run inference in applications that have variable input shapes.

### Creating bucketed models

Before running inference, models should be compiled for different input shapes that are representative of the input dataset. The input shapes that are used to compile the models determine the bucket shapes that are used during inference. The bucket shapes should be chosen to minimize the amount of padding on each new input. Additionally, there should always be a bucket that's large enough to handle the maximum input shape in the dataset. The limit on the number of compiled bucketed models that can be used is described in this [section](#).

### Running inference with bucketing

At inference time, each input should be padded to match the size of the next largest bucket, such that the height and width (or sequence length) of the padded input equals the size of the bucket. Then, the padded input should be passed into the corresponding bucket for inference. If necessary, it's important to remove and/or crop any aberrant predictions that occur in the padded region. For example, in object detection applications, bounding box predictions that occur in the padded regions should be removed to avoid erroneous predictions.

## Examples

The following sections provide examples of applying the bucketing technique to run inference in applications that have variable input shapes.

### Computer vision bucketing

As an example of implementing bucketing for computer vision models, consider an application where the height and width of images in dataset are uniformly distributed between  $[400, 400]$  and  $[800, 800]$ . Given that every input shape between  $[400, 400]$  and  $[800, 800]$  is equally likely, it could make sense to create bucketed models that divide up the range of input shapes into equally sized chunks. For example, we could create bucketed models for the input shapes  $[500, 500]$ ,  $[600, 600]$ ,  $[700, 700]$ , and  $[800, 800]$ .

As an example of running inference with bucketing, let's assume that we created bucketed models for the input shapes  $[500, 500]$ ,  $[600, 600]$ ,  $[700, 700]$ , and  $[800, 800]$ . If we receive an input with shape  $[640, 640]$ , we would pad the input to the next largest bucket,  $[700, 700]$ , and use this bucket for inference. If we receive an input with shape  $[440, 540]$ , we would need to pad the input to the bucket size,  $[600, 600]$ , and use this bucket for inference.

As another example of creating bucketed models, consider a computer vision application where the dataset is not uniformly distributed. As before, let's assume the input shapes range between  $[400, 400]$  to  $[800, 800]$ . Now, let's assume the data shape distribution is bimodal, such that  $[540, 540]$  and  $[720, 720]$  are the two most common input shapes. In this example, it might make sense to create bucketed models for input shapes  $[540, 540]$ ,  $[720, 720]$ , and  $[800, 800]$  to target the most common shapes while still including the entire range of input shapes.

### End-to-end computer vision bucketing example

In this example, we run inference in a computer vision application that has variable shaped images that range in shape from  $[400, 400]$  to  $[800, 800]$ . We create bucketed models for the input shapes  $[500, 500]$ ,  $[600, 600]$ ,  $[700, 700]$ , and  $[800, 800]$  to handle the variable input shapes.

```
import numpy as np
import torch
from torchvision import models
import torch_neuron

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Define the bucket sizes that will be used for compilation and inference
bucket_sizes = [(500, 500), (600, 600), (700, 700), (800, 800)]

# Create the bucketed models by compiling a model for each bucket size
buckets = {}
for bucket_size in bucket_sizes:
    # Create an example input that is the desired bucket size
    h, w = bucket_size
    image = torch.rand([1, 3, h, w])

    # Compile with the example input to create the bucketed model
    model_neuron = torch.neuron.trace(model, image)

    # Run a warm up inference to load the model into Inferentia memory
    model_neuron(image)

    # Add the bucketed model based on its bucket size
    buckets[bucket_size] = model_neuron

def get_bucket_and_pad_image(image):
    # Determine which bucket size to use
    oh, ow = image.shape[-2:]
    target_bucket = None
    for bucket_size in bucket_sizes:
        # Choose a bucket that's larger in both the height and width dimensions
        if oh <= bucket_size[0] and ow <= bucket_size[1]:
            target_bucket = bucket_size
            break

    # Pad the image to match the size of the bucket
```

(continues on next page)

(continued from previous page)

```

h_delta = target_bucket[0] - oh
w_delta = target_bucket[1] - ow

b_pad = h_delta # Bottom padding
l_pad = 0 # Left padding
t_pad = 0 # Top padding
r_pad = w_delta # Right padding

# Pad the height and width of the image
padding_amounts = (l_pad, r_pad, t_pad, b_pad)
image_padded = torch.nn.functional.pad(image, padding_amounts, value=0)

return image_padded, target_bucket

# Run inference on inputs with different shapes
for _ in range(10):
    # Create an image with a random height and width in range [400, 400] to [800, 800]
    h = int(np.random.uniform(low=400, high=800))
    w = int(np.random.uniform(low=400, high=800))
    image = torch.rand(1, 3, h, w)

    # Determine bucket and pad the image
    image_padded, target_bucket = get_bucket_and_pad_image(image)

    # Use the corresponding bucket to run inference
    output = buckets[target_bucket](image_padded)

```

## Natural language processing bucketing

As an example of implementing bucketing for natural language processing models, consider an application where the lengths of tokenized sequences in a dataset are uniformly distributed between 0 and 128 tokens. Given that every tokenized sequence length between 0 and 128 is equally likely, it might make sense to create bucketed models that divide up the range of tokenized sequence lengths into equally sized chunks. For example, we could create bucketed models for tokenized sequence lengths 64 and 128.

As an example of running inference with bucketing, let's assume that we created bucketed models for the input tokenized sequence lengths 64 and 128. If we receive a tokenized sequence with length 55, we would need to pad it to the bucket size 64 and use this bucket for inference. If we receive a tokenized sequence with length 112, we would need to pad it to the bucket size 128 and use this bucket for inference.

## End-to-end natural language processing bucketing example

In this example, we run inference in a natural language processing application that has variable length tokenized sequences that range from 0 to 128. We create bucketed models for lengths 64 and 128 to handle the variable input lengths.

```

import numpy as np
import torch
from transformers import AutoTokenizer, AutoModelForSequenceClassification

```

(continues on next page)

(continued from previous page)

```

import torch_neuron

# Build tokenizer and model
tokenizer = AutoTokenizer.from_pretrained("bert-base-cased-finetuned-mrpc")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-cased-finetuned-
↳mrpc", return_dict=False)
model.eval()

# Define the bucket sizes that will be used for compilation and inference
bucket_sizes = [64, 128]

# Create the bucketed models by compiling a model for each bucket size
buckets = {}
for bucket_size in bucket_sizes:
    # Setup some example inputs
    sequence_0 = "The company HuggingFace is based in New York City"
    sequence_1 = "HuggingFace's headquarters are situated in Manhattan"

    # Create an example input that is the desired bucket size
    paraphrase = tokenizer.encode_plus(sequence_0,
                                       sequence_1,
                                       max_length=bucket_size,
                                       padding='max_length',
                                       truncation=True,
                                       return_tensors="pt")

    # Convert example inputs to a format that is compatible with TorchScript tracing
    example_inputs_paraphrase = paraphrase['input_ids'], paraphrase['attention_mask'],
↳paraphrase['token_type_ids']

    # Compile with the example input to create the bucketed model
    model_neuron = torch.neuron.trace(model, example_inputs_paraphrase)

    # Run a warm up inference to load the model into Inferentia memory
    model_neuron(*example_inputs_paraphrase)

    # Add the bucketed model based on its bucket size
    buckets[bucket_size] = model_neuron

def get_bucket_and_pad_paraphrase(paraphrase):
    # Determine which bucket size to use
    inputs = paraphrase['input_ids']
    attention = paraphrase['attention_mask']
    token_type = paraphrase['token_type_ids']
    paraphrase_len = inputs.shape[1]
    target_bucket = None
    for bucket_size in bucket_sizes:
        if paraphrase_len <= bucket_size:
            target_bucket = bucket_size
            break

```

(continues on next page)

(continued from previous page)

```

# Pad the paraphrase to match the size of the bucket
delta = target_bucket - paraphrase_len
zeros = torch.zeros([1, delta], dtype=torch.long)
inputs = torch.cat([inputs, zeros], dim=1)
attention = torch.cat([attention, zeros], dim=1)
token_type = torch.cat([token_type, zeros], dim=1)

paraphrase_padded = inputs, attention, token_type
return paraphrase_padded, target_bucket

# Create two sample sequences
sequence_0 = ("The only other bear similar in size to the polar bear is the "
              "Kodiak bear, which is a subspecies of the brown bear. Adult male "
              "polar bears weigh 350-700 kg and measure 2.4-3 meters in total "
              "length. All bears are short-tailed, the polar bear's tail is "
              "relatively the shortest amongst living bears.")
sequence_1 = ("Around the Beaufort Sea, however, mature males reportedly "
              "average 450 kg. Adult females are roughly half the size of males "
              "and normally weigh 150-250 kg, measuring 1.8-2.4 meters in length. "
              "The legs are stocky and the ears and tail are small.")

# Run inference on inputs with different shapes
# We create the variable shapes by randomly cropping the sequences
for _ in range(10):
    # Get random sequence lengths between 0 and 128
    paraphrase_len = int(np.random.uniform(128))

    # Crop the paraphrase
    paraphrase_cropped = tokenizer.encode_plus(sequence_0,
                                                sequence_1,
                                                max_length=paraphrase_len,
                                                padding='max_length',
                                                truncation=True,
                                                return_tensors="pt")

    # Determine bucket and pad the paraphrase
    paraphrase_padded, target_bucket = get_bucket_and_pad_paraphrase(paraphrase_cropped)

    # Use the corresponding bucket to run inference
    output = buckets[target_bucket](*paraphrase_padded)

```

This document is relevant for: Inf1

This document is relevant for: Inf1

## Data Parallel Inference on Torch Neuron

### Table of Contents

- *Introduction*
- *Data parallel inference*
- *torch.neuron.DataParallel*
  - *NeuronCore selection*
  - *Batch dim*
  - *Dynamic batching*
  - *Performance optimizations*
- *Examples*
  - *Default usage*
  - *Specifying NeuronCores*
  - *DataParallel with dim != 0*
  - *Dynamic batching*
  - *Dynamic batching disabled*
  - *Full tutorial with torch.neuron.DataParallel*

### Introduction

This guide introduces `torch.neuron.DataParallel()`, a Python API that implements data parallelism on `ScriptModule` models created by the `/neuron-guide/neuron-frameworks/pytorch-neuron/api-compilation-python-api.rst`. The following sections explain how data parallelism can improve the performance of inference workloads on Inferentia, including how `torch.neuron.DataParallel()` uses dynamic batching to run inference on variable input sizes. It covers an overview of the `torch.neuron.DataParallel()` module and provides a few *example data parallel applications*.

### Data parallel inference

Data Parallelism is a form of parallelization across multiple devices or cores, referred to as nodes. Each node contains the same model and parameters, but data is distributed across the different nodes. By distributing the data across multiple nodes, data parallelism reduces the total execution time of large batch size inputs compared to sequential execution. Data parallelism works best for smaller models in latency sensitive applications that have large batch size requirements.

## `torch.neuron.DataParallel`

To fully leverage the Inferentia hardware, we want to use all available NeuronCores. An `inf1.xlarge` and `inf1.2xlarge` have four NeuronCores, an `inf1.6xlarge` has 16 NeuronCores, and an `inf1.24xlarge` has 64 NeuronCores. For maximum performance on Inferentia hardware, we can use `torch.neuron.DataParallel()` to utilize all available NeuronCores.

`torch.neuron.DataParallel()` implements data parallelism at the module level by replicating the Neuron model on all available NeuronCores and distributing data across the different cores for parallelized inference. This function is analogous to `DataParallel` in PyTorch. `torch.neuron.DataParallel()` requires PyTorch  $\geq 1.8$ .

The following sections provide an overview of some of the features of `torch.neuron.DataParallel()` that enable maximum performance on Inferentia.

### NeuronCore selection

By default, `DataParallel` will try to use all NeuronCores allocated to the current process to fully saturate the Inferentia hardware for maximum performance. It is more efficient to make the batch dimension divisible by the number of NeuronCores. This will ensure that NeuronCores are not left idle during parallel inference and the Inferentia hardware is fully utilized.

In some applications, it is advantageous to use a subset of the available NeuronCores for `DataParallel` inference. `DataParallel` has a `device_ids` argument that accepts a list of `int` or `'nc: #'` that specify the NeuronCores to use for parallelization. See *Specifying NeuronCores* for an example of how to use `device_ids` argument.

### Batch dim

`DataParallel` accepts a `dim` argument that denotes the batch dimension used to split the input data for distributed inference. By default, `DataParallel` splits the inputs on `dim = 0` if the `dim` argument is not specified. For applications with a non-zero batch dim, the `dim` argument can be used to specify the inference-time input batch dimension. *DataParallel with `dim != 0`* provides an example of data parallel inference on inputs with batch `dim = 2`.

### Dynamic batching

Batch size has a direct impact on model performance. The Inferentia chip is optimized to run with small batch sizes. This means that a Neuron compiled model can outperform a GPU model, even if running single digit batch sizes.

As a general best practice, we recommend optimizing your model's throughput by compiling the model with a small batch size and gradually increasing it to find the peak throughput on Inferentia.

Dynamic batching is a feature that allows you to use tensor batch sizes that the Neuron model was not originally compiled against. This is necessary because the underlying Inferentia hardware will always execute inferences with the batch size used during compilation. Fixed batch size execution allows tuning the input batch size for optimal performance. For example, batch size 1 may be best suited for an ultra-low latency on-demand inference application, while batch size  $> 1$  can be used to maximize throughput for offline inferencing. Dynamic batching is implemented by slicing large input tensors into chunks that match the batch size used during the `torch_neuron.trace()` compilation call.

The `torch.neuron.DataParallel()` class automatically enables dynamic batching on eligible models. This allows us to run inference in applications that have inputs with a variable batch size without needing to recompile the model. See *Dynamic batching* for an example of how `DataParallel` can be used to run inference on inputs with a dynamic batch size without needing to recompile the model.



Dynamic batching using small batch sizes can result in sub-optimal throughput because it involves slicing tensors into chunks and iteratively sending data to the hardware. Using a larger batch size at compilation time can use the Inferentia hardware more efficiently in order to maximize throughput. You can test the tradeoff between individual request latency and total throughput by fine-tuning the input batch size.

Automatic batching in the DataParallel module can be disabled using the `disable_dynamic_batching()` function as follows:

```
>>> model_parallel = torch.neuron.DataParallel(model_neuron)
>>> model_parallel.disable_dynamic_batching()
```

If dynamic batching is disabled, the compile-time batch size must be equal to the inference-time batch size divided by the number of NeuronCores. *DataParallel with dim != 0* and *Dynamic batching disabled* provide examples of running DataParallel inference with dynamic batching disabled.

## Performance optimizations

The DataParallel module has a `num_workers` attribute that can be used to specify the number of worker threads used for multithreaded inference. By default, `num_workers = 2 * number of NeuronCores`. This value can be fine tuned to optimize DataParallel performance.

DataParallel has a `split_size` attribute that dictates the size of the input chunks that are distributed to each Neuron-Core. By default, `split_size = max(1, input.shape[dim] // number of NeuronCores)`. This value can be modified to optimally match the inference input chunk size with the compile-time batch size.

## Examples

The following sections provide example usages of the `torch.neuron.DataParallel()` module.

### Default usage

The default DataParallel use mode will replicate the model on all available NeuronCores in the current process. The inputs will be split on `dim=0`.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module
model_parallel = torch.neuron.DataParallel(model_neuron)

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])
```

(continues on next page)

(continued from previous page)

```
# Run inference with a batched input
output = model_parallel(image_batched)
```

## Specifying NeuronCores

The following example uses the `device_ids` argument to use the first three NeuronCores for DataParallel inference.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module, run on the first three NeuronCores
# Equivalent to model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2])
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=['nc:0', 'nc:1', 'nc:2'])

# Create a batched input
batch_size = 5
image_batched = torch.rand([batch_size, 3, 224, 224])

# Run inference with a batched input
output = model_parallel(image_batched)
```

## DataParallel with dim != 0

In this example we run DataParallel inference using four NeuronCores and `dim = 2`. Because `dim != 0`, dynamic batching is not enabled. Consequently, the DataParallel inference-time batch size must be four times the compile-time batch size. DataParallel will generate a warning that dynamic batching is disabled because `dim != 0`.

```
import torch
import torch_neuron

# Create an example model
class Model(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.conv = torch.nn.Conv2d(3, 3, 3)

    def forward(self, x):
```

(continues on next page)

(continued from previous page)

```

        return self.conv(x) + 1

model = Model()
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 8, 8])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module using 4 NeuronCores and dim = 2
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2, 3], dim=2)

# Create a batched input
# Note that image_batched.shape[dim] / len(device_ids) == image.shape[dim]
batch_size = 4 * 8
image_batched = torch.rand([1, 3, batch_size, 8])

# Run inference with a batched input
output = model_parallel(image_batched)

```

## Dynamic batching

In the following example, we use the `torch.neuron.DataParallel()` module to run inference using several different batch sizes without recompiling the Neuron model.

```

import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module
model_parallel = torch.neuron.DataParallel(model_neuron)

# Create batched inputs and run inference on the same model
batch_sizes = [2, 3, 4, 5, 6]
for batch_size in batch_sizes:
    image_batched = torch.rand([batch_size, 3, 224, 224])

    # Run inference with a batched input
    output = model_parallel(image_batched)

```

## Dynamic batching disabled

In the following example, we use `torch.neuron.DataParallel.disable_dynamic_batching()` to disable dynamic batching. We provide an example of a batch size that will not work when dynamic batching is disabled as well as an example of a batch size that does work when dynamic batching is disabled.

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input
image = torch.rand([1, 3, 224, 224])
model_neuron = torch.neuron.trace(model, image)

# Create the DataParallel module and use 4 NeuronCores
model_parallel = torch.neuron.DataParallel(model_neuron, device_ids=[0, 1, 2, 3], dim=0)

# Disable dynamic batching
model_parallel.disable_dynamic_batching()

# Create a batched input (this won't work)
batch_size = 8
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will fail because dynamic batching is disabled and
# image_batched.shape[dim] / len(device_ids) != image.shape[dim]
# output = model_parallel(image_batched)

# Create a batched input (this will work)
batch_size = 4
image_batched = torch.rand([batch_size, 3, 224, 224])

# This will work because
# image_batched.shape[dim] / len(device_ids) == image.shape[dim]
output = model_parallel(image_batched)
```

## Full tutorial with torch.neuron.DataParallel

For an end-to-end tutorial that uses DataParallel, see the PyTorch Resnet Tutorial.

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Developer Guide - PyTorch Neuron (torch-neuron) LSTM Support

The *torch-neuron* package can support LSTM operations and yield high performance on both fixed-length and variable-length sequences. Most network configurations can be supported, with the exception of those that require PackedSequence usage outside of LSTM or `pad_packed_sequence()` operations. Neuron must guarantee that the shapes can remain fixed throughout the network.

The following sections describe which scenarios can and cannot be supported.

### Supported Usage

#### Fixed-Length Sequences

In normal usage of an LSTM, the inputs and outputs are expected to be a fixed size sequence length. This is the most basic usage of an LSTM but may not be applicable to applications where the input sequence length may vary.

```
import torch
import torch_neuron

class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs):
        output, (ht, ct) = self.lstm(inputs)
        return output, (ht, ct)

# Example Inputs
seq_len, batch_size, input_size = 5, 2, 3
inputs = torch.rand(seq_len, batch_size, input_size)

# Trace
torch_neuron.trace(Network(), (inputs,))
```

#### Packed Input, Padded Output, *Pre-Sorted* Inputs

A common usage of an LSTM is when the input sequence sizes vary according to an input sequence lengths (such as tokens).

For example, the following sentences could result in two different sequence lengths after tokenization:

```
# Input
text = [
    'Hello, sailor',
    'Example',
]

# ... Tokenization ...
```

(continues on next page)

(continued from previous page)

```
# Result
tokens = [
    [101, 7592, 1010, 11803, 102],
    [101, 2742, 102, 0, 0],
]
lengths = [5, 3]
```

Because the lengths are different, the final LSTM state will be dependent upon the lengths of each sequence in the batch. Torch provides a way to deal with these types of sequences by densely packing batches into a `PackedSequence`. The most common way this is constructed is by using the `pack_padded_sequence()` utility function prior to feeding inputs into the LSTM.

Packing the above sequences would result in the following data and batch size tensors.

```
data = [101, 101, 7592, 2742, 1010, 102, 11803, 102]
batch_sizes = [2, 2, 2, 1, 1]
```

In addition to correctly computing final LSTM state, using a packed sequence instead of a padded sequence also improves model performance on CPU. On Neuron, where computation is fixed to the maximum length ahead of time, **this does not improve performance**.

When an LSTM is processing a `PackedSequence`, it must do so in a descending sorted length order. To ensure that sequences are sorted, `pack_padded_sequence()` provides an `enforce_sorted` flag. When `enforce_sorted` is `True`, the input is *already expected* to contain sequences sorted by length in a decreasing order along the batch dimension. Note that this must be enforced in the application-level code but is only relevant when batch size > 1.

The following network can compile successfully because the input and output to the network are guaranteed to be a fixed shape. The input shape is expected to be a padded tensor and the output tensor is expected to be padded to the maximum sequence length using the `pad_packed_sequence()` function call:

```
import torch
import torch_neuron

class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs, lengths):
        packed_input = torch.nn.utils.rnn.pack_padded_sequence(
            inputs,
            lengths=lengths,
            enforce_sorted=True,
        )
        packed_result, (ht, ct) = self.lstm(packed_input)
        padded_result, _ = torch.nn.utils.rnn.pad_packed_sequence(packed_result)
        return padded_result, ht, ct

# Example Inputs
seq_len, batch_size, input_size = 5, 2, 3
inputs = torch.rand(seq_len, batch_size, input_size)
lengths = torch.tensor([seq_len] * batch_size)
```

(continues on next page)

(continued from previous page)

```
# Trace
torch_neuron.trace(Network(), (inputs, lengths))
```

### Packed Input, Padded Output, *Unsorted* Inputs

When `enforce_sorted` is `False`, the input will be sorted unconditionally. This causes some CPU overhead on Neuron because unsupported operators will be inserted into the graph such as `aten::sort` and `aten::scatter_`. The `aten::lstm` operation can still be supported, but it will be less efficient than when `enforce_sorted` is `True`.

The following code is able to be traced, but results in the sorting operations running on CPU. This is not problematic in this case because the `aten::sort` and `aten::scatter_` are executed on CPU at the very beginning of the graph just prior to Neuron execution.

Like the previous example, the call to `pad_packed_sequence()` ensures that the output is a fixed-shape based on the maximum sequence length.

```
import torch
import torch_neuron

class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs, lengths):
        packed_input = torch.nn.utils.rnn.pack_padded_sequence(
            inputs,
            lengths=lengths,
            enforce_sorted=False,
        )
        packed_result, (ht, ct) = self.lstm(packed_input)
        padded_result, _ = torch.nn.utils.rnn.pad_packed_sequence(packed_result)
        return padded_result, ht, ct

# Example Inputs
seq_len, batch_size, input_size = 5, 2, 3
inputs = torch.rand(seq_len, batch_size, input_size)
lengths = torch.tensor([seq_len] * batch_size)

# Trace
trace = torch_neuron.trace(Network(), (inputs, lengths))
```

## Packed Inputs, Final Hidden & Cell State Only

When **only** the final LSTM hidden & cell state is used, it does not matter if the inputs are packed or unpacked since these state tensors will not vary in size.

```
import torch
import torch_neuron

class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs, lengths):
        packed_input = torch.nn.utils.rnn.pack_padded_sequence(
            inputs,
            lengths=lengths,
            enforce_sorted=True,
        )
        packed_output, (ht, ct) = self.lstm(packed_input)
        return ht, ct

# Example Inputs
seq_len, batch_size, input_size = 5, 2, 3
inputs = torch.rand(seq_len, batch_size, input_size)
lengths = torch.tensor([seq_len] * batch_size)

# Trace
trace = torch_neuron.trace(Network(), (inputs, lengths))
```

Note that when the `packed_output` is unused, it does not need to be passed to the `pack_padded_sequence()` to enable the LSTM to be compiled.

## Unsupported Usage

Neuron does not support the use of a `PackedSequence` outside of the LSTM operation and the `pack_padded_sequence()` operation. This is because the shape of a `PackedSequence` can vary depending on the input data. This is incompatible with the Neuron restriction that all tensor sizes must be known at compilation time. When a `PackedSequence` is used only by an LSTM or `pack_padded_sequence()` operation, Neuron *can guarantee* the size of the intermediary tensors by padding on behalf of the application.

This means that If the `PackedSequence` is either used by a different operation or returned from the network this would result in all of the LSTM operations to be executed on CPU or the network compilation will fail.



## PackedSequence Returned

The following is unsupported because the PackedSequence result of the LSTM is returned by the network:

```
class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs, lengths):
        packed_input = torch.nn.utils.rnn.pack_padded_sequence(
            inputs,
            lengths=lengths,
            enforce_sorted=False,
        )
        packed_result, (ht, ct) = self.lstm(packed_input)
        return packed_result.data, ht, ct
```

**Behavior:** In this case, compilation fails and the following warning is generated:

Operator "aten::lstm" consuming a PackedSequence input can only be supported when its  
 ↳ corresponding PackedSequence output is unused or unpacked using "aten::\_pad\_packed\_  
 ↳ input". Found usage by "prim::Return"

**Resolution:** To avoid this error, the packed\_result should be padded prior to being returned from the network by using pad\_packed\_sequence()

## Invalid PackedSequence Usage

The following is unsupported because the PackedSequence result of the LSTM is used by a non-LSTM operator:

```
class Network(torch.nn.Module):

    def __init__(self):
        super().__init__()
        self.lstm = torch.nn.LSTM(input_size=3, hidden_size=7)

    def forward(self, inputs, lengths):
        packed_input = torch.nn.utils.rnn.pack_padded_sequence(
            inputs,
            lengths=lengths,
            enforce_sorted=False,
        )
        packed_result, (ht, ct) = self.lstm(packed_input)
        return torch.max(packed_result.data)
```

**Behavior:** In this case, compilation fails and the following warning is generated:

Operator "aten::lstm" consuming a PackedSequence input can only be supported when its  
 ↳ corresponding PackedSequence output is unused or unpacked using "aten::\_pad\_packed\_  
 ↳ input". Found usage by "aten::max"

**Resolution:** To avoid this error, the `packed_result` should be padded prior to being used in the `max()` from the network by using `pad_packed_sequence()`.

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## PyTorch Neuron (torch-neuron) Core Placement

This programming guide describes the available techniques and APIs to be able to allocate NeuronCores to a process and place models onto specific NeuronCores. In order of precedence, the current recommendation is to use the following placement techniques:

1. For most regular models, default core placement should be used in conjunction with `NEURON_RT_NUM_CORES` (*Default Core Allocation & Placement*).
2. For more specific core placement for NeuronCore Pipelined models, then `NEURONCORE_GROUP_SIZES` should be used (*NEURONCORE\_GROUP\_SIZES*).
3. Finally, for even more granular control, then the beta explicit placement APIs may be used (*Explicit Core Placement [Beta]*).

### Table of Contents

- *PyTorch Neuron (torch-neuron) Core Placement*
  - *NeuronCore Pipeline*
  - *Default Core Allocation & Placement*
    - \* *Example: Default*
    - \* *Example: NEURON\_RT\_NUM\_CORES*
    - \* *Example: NEURON\_RT\_VISIBLE\_CORES*
    - \* *Example: Overlapping Models*
    - \* *Example: Multiple Processes*
  - *NEURONCORE\_GROUP\_SIZES*
    - \* *Example: Single NeuronCore Group*
    - \* *Example: Multiple NeuronCore Groups*
    - \* *Issue: Overlapping Models with Differing Model Sizes*
    - \* *Issue: Incompatible Model Sizes*
    - \* *Issue: Multiple Model Copies*
    - \* *Issue Summary*
  - *Explicit Core Placement [Beta]*
    - \* *Example: Manual Core Selection*
    - \* *Example: Automatic Multicore*
    - \* *Example: Explicit Replication*

The following guide will assume a machine with 8 NeuronCores:

- NeuronCores will use the notation `nc0`, `nc1`, etc.
- NeuronCore Groups will use the notation `ncg0`, `ncg1` etc.
- Models will use the notation `m0`, `m1` etc.

NeuronCores, NeuronCore Groups, and model allocations will be displayed in the following format:

Note that the actual cores that are visible to the process can be adjusted according to the [NeuronX Runtime Configuration](#).

## NeuronCore Pipeline

A key concept to understand the intent behind certain core placement strategies is NeuronCore Pipelining (See [NeuronCore Pipeline](#)). NeuronCore Pipelining allows a model to be automatically split into pieces and executed on different NeuronCores.

For most models only 1 NeuronCore will be required for execution. A model will **only** require more than one NeuronCore when using NeuronCore Pipeline. When model pipelining is enabled, the model is split between multiple NeuronCores and data is transferred between them. For example, if the compiler flag `--neuroncore-pipeline-cores 4` is used, this splits the model into 4 pieces to be executed on 4 separate NeuronCores.

## Default Core Allocation & Placement

The most basic requirement of an inference application is to be able to place a single model on a single NeuronCore. More complex applications may use multiple NeuronCores or even multiple processes each executing different models. The important thing to note about designing an inference application is that a single NeuronCore will always be allocated to a single process. *Processes do not share NeuronCores*. Different configurations can be used to ensure that an application process has enough NeuronCores allocated to execute its model(s):

- **Default:** A process will attempt to take ownership of **all NeuronCores** visible on the instance. This should be used when an instance is only running a single inference process since no other process will be allowed to take ownership of any NeuronCores.
- **NEURON\_RT\_NUM\_CORES:** Specify the **number of NeuronCores** to allocate to the process. This places no restrictions on which NeuronCores will be used, however, the resulting NeuronCores will always be contiguous. This should be used in multi-process applications where each process should only use a subset of NeuronCores.
- **NEURON\_RT\_VISIBLE\_CORES:** Specifies exactly **which NeuronCores** are allocated to the process by index. Similar to **NEURON\_RT\_NUM\_CORES**, this can be used in multi-process applications where each process should only use a subset of NeuronCores. This provides more fined-grained controls over the exact NeuronCores that are allocated to a given process.
- **NEURONCORE\_GROUP\_SIZES:** Specifies a number of **NeuronCore Groups** which are allocated to the process. This is described in more detail in the [NEURONCORE\\_GROUP\\_SIZES](#) section.

See the [NeuronX Runtime Configuration](#) for more environment variable details.

### Example: Default

#### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
m1 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc1
```

With no environment configuration, the process will take ownership of all NeuronCores. In this example, only two of the NeuronCores are used by the process and the remaining are allocated but left idle.

### Example: NEURON\_RT\_NUM\_CORES

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '2'
```

#### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
m1 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc1
```

Since there is no other process on the instance, only the first 2 NeuronCores will be acquired by the process. Models load in a simple linear order to the least used NeuronCores.

### Example: NEURON\_RT\_VISIBLE\_CORES

#### Environment Setup:

```
export NEURON_RT_VISIBLE_CORES = '4-5'
```

#### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc4
m1 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc5
```

Unlike NEURON\_RT\_NUM\_CORES, setting the visible NeuronCores allows the process to take control of a specific contiguous set. This allows an application to have a more fine-grained control of where models will be placed.

## Example: Overlapping Models

### Environment Setup:

```
export NEURON_RT_VISIBLE_CORES = '0-1'
```

### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
m1 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc0-nc1
m2 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc1
```

This shows how models may share NeuronCores but the default model placement will attempt to evenly distribute NeuronCore usage rather than overlapping all models on a single NeuronCore.

## Example: Multiple Processes

### Environment Setup:

```
export NEURON_RT_NUM_CORES = '2'
```

### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
m1 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc1
```

In this example, if the script is run **twice**, the following allocations will be made:

Note that each process will take ownership of as many NeuronCores as is specified by the NEURON\_RT\_NUM\_CORES configuration.

## NEURONCORE\_GROUP\_SIZES

**Important:** The use of explicit core placement should only be used when a specific performance goal is required. By default torch-neuron places models on the **least used** NeuronCores. This should be optimal for most applications.

Secondly, NEURONCORE\_GROUP\_SIZES is being deprecated in a future release and should be avoided in favor of newer placement methods. Use NEURON\_RT\_NUM\_CORES or NEURON\_RT\_VISIBLE\_CORES with default placement if possible (See [Default Core Allocation & Placement](#))

In the current release of NeuronSDK, the most well-supported method of placing models onto specific NeuronCores is to use the NEURONCORE\_GROUP\_SIZES environment variable. This will define a set of “NeuronCore Groups” for the application process.

NeuronCore Groups are *contiguous sets of NeuronCores* that are allocated to a given process. Creating groups allows an application to ensure that a model has a defined set of NeuronCores that will always be allocated to it.

Note that NeuronCore Groups *can* be used to allocate non-pipelined models (those requiring exactly 1 NeuronCore) to specific NeuronCores but this is not the primary intended use. The intended use of NeuronCore Groups is to ensure pipelined models (those requiring >1 NeuronCore) have exclusive access to a specific set of contiguous NeuronCores.

In the cases where models are being used *without* NeuronCore Pipeline, the general recommendation is to use default placement (See [Default Core Allocation & Placement](#)).

The following section demonstrates how `NEURONCORE_GROUP_SIZES` can be used and the issues that may arise.

### Example: Single NeuronCore Group

In the example where one model requires 4 NeuronCores, the correct environment configuration would be:

#### Environment Setup:

```
export NEURONCORE_GROUP_SIZES = '4'
```

#### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-4-neuron-pipeline-cores.pt') # Loads to nc0-nc3
```

This is the most basic usage of a NeuronCore Group. The environment setup causes the process to take control of 4 NeuronCores and then the script loads a model compiled with a NeuronCore Pipeline size of 4 to the first group.

### Example: Multiple NeuronCore Groups

With more complicated configurations, the intended use of `NEURONCORE_GROUP_SIZES` is to create 1 Group per model with the correct size to ensure that the models are placed on the intended NeuronCores. Similarly, the environment would need to be configured to create a NeuronCore Group for each model:

#### Environment Setup:

```
export NEURONCORE_GROUP_SIZES = '3,4,1'
```

#### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-3-neuron-pipeline-cores.pt') # Loads to nc0-nc2
m1 = torch.jit.load('model-with-4-neuron-pipeline-cores.pt') # Loads to nc3-nc6
m2 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc7
```

## Issue: Overlapping Models with Differing Model Sizes

When multiple models are loaded to a single NeuronCore Group, this can cause unintended inefficiencies. A single model is only intended to span a single NeuronCore Group. Applications with many models of varying sizes can be restricted by NeuronCore Group configurations since the most optimal model layout may require more fine-grained controls.

### Environment Setup:

```
export NEURONCORE_GROUP_SIZES = '2,2'
```

### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc0-nc1
m1 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc2-nc3
m2 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
m3 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc2
m4 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc0
```

Here the NEURONCORE\_GROUP\_SIZES does not generate an optimal layout because placement strictly follows the layout of NeuronCore Groups. A potentially more optimal layout would be to place m4 onto nc1. In this case, since a pipelined model will not be able to have exclusive access to a set of NeuronCores, the default NeuronCore placement (no NeuronCore Groups specified) would more evenly distribute the models.

Also note here that this is an example of where the order of model loads affects which model is assigned to which NeuronCore Group. If the order of the load statements is changed, models may be assigned to different NeuronCore Groups.

## Issue: Incompatible Model Sizes

Another problem occurs when attempting to place a model which does not evenly fit into a single group:

### Environment Setup:

```
export NEURONCORE_GROUP_SIZES = '2,2'
```

### Python Script:

```
import torch
import torch_neuron

m0 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc0-nc1
m1 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc2-nc3
m2 = torch.jit.load('model-with-3-neuron-pipeline-cores.pt') # Loads to nc0-nc2
```

The model will be placed *across* NeuronCore Groups since there is no obvious group to assign the model to according to the environment variable configuration. Depending on the individual model and application requirements, the placement here may not be optimal.

## Issue: Multiple Model Copies

It is common in inference serving applications to use multiple replicas of a single model across different NeuronCores. This allows the hardware to be fully utilized to maximize throughput. In this scenario, when using NeuronCore Groups, the only way to replicate a model on multiple NeuronCores is to create a *new model* object. In the example below, 4 models loads are performed to place a model in each NeuronCore Group.

### Environment Setup:

```
export NEURONCORE_GROUP_SIZES = '2,2,2,2'
```

### Python Script:

```
import torch
import torch_neuron

models = list()
for _ in range(4):
    model = torch.jit.load('model-with-2-neuron-pipeline-cores.pt')
    models.append(model)
```

The largest consequence of this type of model allocation is that the application code is responsible for routing inference requests to models. There are a variety of ways to implement the inference switching but in all cases routing logic needs to be implemented in the application code.

## Issue Summary

The use of NEURONCORE\_GROUP\_SIZES has the following problems:

- **Variable Sized Models:** Models which require crossing NeuronCore Group boundaries may be placed poorly. This means group configuration limits the size of which models can be loaded.
- **Model Load Order:** Models are loaded to NeuronCore Groups greedily. This means that the order of model loads can potentially negatively affect application performance by causing unintentional overlap.
- **Implicit Placement:** NeuronCore Groups cannot be explicitly chosen in the application code.
- **Manual Replication:** When loading multiple copies of a model to different NeuronCore Groups, this requires that multiple model handles are used.

## Explicit Core Placement [Beta]

To address the limitations of NEURONCORE\_GROUP\_SIZES, a new set of APIs has been added which allows specific NeuronCores to be chosen by the application code. These can be found in the `torch_neuron_core_placement_api` documentation.



### Example: Manual Core Selection

The most direct usage of the placement APIs is to manually select the start NeuronCore that each model is loaded to. This will automatically use as many NeuronCores as is necessary for that model (1 for most models, >1 for NeuronCore Pipelines models).

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '4'
```

#### Python Script:

```
import torch
import torch_neuron

# NOTE: Order of loads does NOT matter

with torch_neuron.experimental.neuron_cores_context(2):
    m1 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc2-nc3

with torch_neuron.experimental.neuron_cores_context(0):
    m2 = torch.jit.load('model-with-3-neuron-pipeline-cores.pt') # Loads to nc0-nc2

with torch_neuron.experimental.neuron_cores_context(0):
    m0 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads to nc0-nc1

with torch_neuron.experimental.neuron_cores_context(3):
    m3 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads to nc3
```

Note that this directly solves the NEURONCORE\_GROUP\_SIZES issues of:

- **Variable Sized Models:** Now since models are directly placed on the NeuronCores requested by the application, there is no disconnect between the model sizes and NeuronCore Group sizes.
- **Model Load Order:** Since the NeuronCores are explicitly selected, there is no need to be careful about the order in which models are loaded since they can be placed deterministically regardless of the load order.
- **Implicit Placement:** Similarly, explicit placement means there is no chance that a model will end up being allocated to an incorrect NeuronCore Group.

### Example: Automatic Multicore

Using explicit core placement it is possible to replicate a model to multiple NeuronCores simultaneously. This means that a single model object within python can utilize all available NeuronCores (or NeuronCores allocated to the process).

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '8'
```

#### Python Script:

```
import torch
import torch_neuron

with torch_neuron.experimental.multicore_context():
```

(continues on next page)

(continued from previous page)

```
m0 = torch.jit.load('model-with-1-neuron-pipeline-cores.pt') # Loads replications.
↪to nc0-nc7
```

This addresses the last NEURONCORE\_GROUP\_SIZES issue of:

- **Manual Replication:** Since models can be automatically replicated to multiple NeuronCores, this means that applications no longer need to implement routing logic and perform multiple loads.

This API has a secondary benefit that the exact same loading logic can be used on an `inf1.xlarge` or an `inf1.6xlarge`. In either case, it will use all of the NeuronCores that are visible to the process. This means that no special logic needs to be coded for different instance types.

### Example: Explicit Replication

Replication is also possible with the `neuron_cores_context()` API. The number of replications is chosen by `replications = floor(nc_count / cores_per_model)`.

#### Environment Setup:

```
export NEURON_RT_NUM_CORES = '8'
```

#### Python Script:

```
import torch
import torch_neuron

with torch_neuron.experimental.neuron_cores_context(start_nc=2, nc_count=4):
    m0 = torch.jit.load('model-with-2-neuron-pipeline-cores.pt') # Loads replications.
↪to nc2-nc5
```

This document is relevant for: `Inf1`

- *Running Inference on Variable Input Shapes with Bucketing*
- *Data Parallel Inference on PyTorch Neuron*
- *Developer Guide - PyTorch Neuron (torch-neuron) LSTM Support*
- *PyTorch Neuron (torch-neuron) Core Placement*

This document is relevant for: `Inf1`

This document is relevant for: `Inf1`

### Misc (torch-neuron)

This document is relevant for: `Inf1`

## PyTorch Neuron (torch-neuron) Supported operators

Current operator lists may be generated with these commands inside python:

```
import torch.neuron
print(*torch.neuron.get_supported_operations(), sep='\n')
```

### PyTorch Neuron release [package version 1.\*\*.2.9.1.0, SDK 2.13.0]

Date: 08/28/2023

Added support for new operators:

- aten::clamp\_min
- aten::clamp\_max

### PyTorch Neuron release [2.9.0.0]

Date: 03/28/2023

Added support for new operators:

- aten::tensordot
- aten::adaptive\_avg\_pool1d
- aten::prelu
- aten::reflection\_pad2d
- aten::baddbmm
- aten::repeat

### PyTorch Neuron release [2.5.0.0]

Date: 11/23/2022

Added support for new operators:

- aten::threshold
- aten::roll
- aten::instance\_norm
- aten::amin
- aten::amax
- aten::new\_empty
- aten::new\_ones
- aten::tril
- aten::triu
- aten::zero\_

- `aten::all`
- `aten::broadcast_tensors`
- `aten::broadcast_to`
- `aten::logical_and`
- `aten::logical_not`
- `aten::logical_or`
- `aten::logical_xor`
- `aten::_convolution_mode`

Added **limited** support for new operators:

- LSTM Operations. See: [Developer Guide - PyTorch Neuron \(torch-neuron\) LSTM Support](#)
  - `aten::lstm`
  - `aten::_pack_padded_sequence`
  - `aten::_pad_packed_sequence`
- `aten::norm`: Supported when `p` argument is one of (1, 2, inf, -inf, 'fro')

### PyTorch Neuron release [2.2.0.0]

Date: 03/25/2022

Added support for new operators:

- `aten::max_pool2d_with_indices`: Fully supported (Was previously supported only when indices were unused).

### PyTorch Neuron release [2.1.7.0]

Date: 01/20/2022

Added support for new operators:

- `aten::bucketize`
- `aten::any`
- `aten::remainder`
- `aten::clip`
- `aten::repeat_interleave`
- `aten::tensor_split`
- `aten::split_with_sizes`
- `aten::isnan`
- `aten::embedding_renorm_`
- `aten::dot`
- `aten::mv`
- `aten::hardsigmoid`

- `aten::hardswish`
- `aten::trunc`
- **`aten::one_hot`: Supported when `num_classes` is known at trace time.**  
The dynamic version of this operation when `num_classes = -1` is not supported.
- `aten::adaptive_max_pool1d`
- `aten::adaptive_max_pool2d`

## PyTorch Neuron Release [2.0.536.0]

- The following are operators with limited support on Neuron. Unlike fully supported operators, these operators are not returned when using `torch_neuron.get_supported_operations()`. See each operator description for conditional support:
  - `aten::max_pool2d_with_indices` - Supported when indices outputs are not used by a downstream operation. This allows the operation to be compiled to Neuron when it is equivalent to an `aten::max_pool2d`.
  - `aten::max_pool3d_with_indices` - Supported when indices outputs are not used by a downstream operation. This allows the operation to be compiled to Neuron when it is equivalent to an `aten::max_pool3d`.
  - `aten::where` - Supported when used as a conditional selection (3-argument variant). Unsupported when used to generate a dynamic list of indices (1-argument variant). See `torch.where()`.

## PyTorch Neuron Release [2.0.318.0]

Added support for new operators:

- `aten::empty_like`
- `aten::log`
- `aten::type_as`
- `aten::movedim`
- `aten::einsum`
- `aten::argmax`
- `aten::min`
- `aten::argmin`
- `aten::abs`
- `aten::cos`
- `aten::sin`
- `aten::linear`
- `aten::pixel_shuffle`
- `aten::group_norm`
- `aten::_weight_norm`

### PyTorch Neuron Release [1.5.21.0]

No change

### PyTorch Neuron Release [1.5.7.0]

Added support for new operators:

- `aten::erf`
- `prim::DictConstruct`

### PyTorch Neuron Release [1.4.1.0]

No change

### PyTorch Neuron Release [1.3.5.0]

Added support for new operators:

- `aten::numel`
- `aten::ones_like`
- `aten::reciprocal`
- `aten::topk`

### PyTorch Neuron Release [1.2.16.0]

No change

### PyTorch Neuron Release [1.2.15.0]

No change

### PyTorch Neuron Release [1.2.3.0]

Added support for new operators:

- `aten::silu`
- `aten::zeros_like`

### PyTorch Neuron Release [1.1.7.0]

Added support for new operators:

- `aten::_shape_as_tensor`
- `aten::chunk`
- `aten::empty`
- `aten::masked_fill`

### PyTorch Neuron Release [1.0.24045.0]

Added support for new operators:

- `aten::__and__`
- `aten::bmm`
- `aten::clone`
- `aten::expand_as`
- `aten::fill_`
- `aten::floor_divide`
- `aten::full`
- `aten::hardtanh`
- `aten::hardtanh_`
- `aten::le`
- `aten::leaky_relu`
- `aten::lt`
- `aten::mean`
- `aten::ne`
- `aten::softplus`
- `aten::unbind`
- `aten::upsample_bilinear2d`

### PyTorch Neuron Release [1.0.1720.00]

Added support for new operators:

- `aten::constant_pad_nd`
- `aten::meshgrid`

### PyTorch Neuron Release [1.0.1532.0]

Added support for new operators:

- `aten::ones`

### PyTorch Neuron Release [1.0.1522.0]

No change

### PyTorch Neuron Release [1.0.1386.0]

Added support for new operators:

- `aten::ceil`
- `aten::clamp`
- `aten::eq`
- `aten::exp`
- `aten::expand_as`
- `aten::flip`
- `aten::full_like`
- `aten::ge`
- `aten::gt`
- `aten::log2`
- `aten::log_softmax`
- `aten::max`
- `aten::neg`
- `aten::relu`
- `aten::rsqrt`
- `aten::scalarImplicit`
- `aten::sqrt`
- `aten::squeeze`
- `aten::stack`
- `aten::sub`
- `aten::sum`
- `aten::true_divide`
- `aten::upsample_nearest2d`
- `prim::Constant`
- `prim::GetAttr`
- `prim::ImplicitTensorToNum`



- `prim::ListConstruct`
- `prim::ListUnpack`
- `prim::NumToTensor`
- `prim::TupleConstruct`
- `prim::TupleUnpack`

Please note, primitives are included in this list from this release.

### PyTorch Neuron Release [1.0.1168.0]

Added support for new operators:

- `aten::ScalarImplicit`

### PyTorch Neuron Release [1.0.1001.0]

Added support for new operators:

- `aten::detach`
- `aten::floor`
- `aten::gelu`
- `aten::pow`
- `aten::sigmoid`
- `aten::split`

Remove support for operators:

- `aten::embedding`: Does not meet **performance** criteria
- `aten::erf`: Error function does not meet **accuracy** criteria
- `aten::tf_dtype_from_torch`: Internal support function, not an operator

### PyTorch Neuron Release [1.0.825.0]

No change

### PyTorch Neuron Release [1.0.763.0]

Added support for new operators:

- `aten::Int`
- `aten::arange`
- `aten::contiguous`
- `aten::div`
- `aten::embedding`
- `aten::erf`

- `aten::expand`
- `aten::eye`
- `aten::index_select`
- `aten::layer_norm`
- `aten::matmul`
- `aten::mm`
- `aten::permute`
- `aten::reshape`
- `aten::rsub`
- `aten::select`
- `aten::size`
- `aten::slice`
- `aten::softmax`
- `aten::tf_dtype_from_torch`
- `aten::to`
- `aten::transpose`
- `aten::unsqueeze`
- `aten::view`
- `aten::zeros`

Remove support for operators:

- `aten::tf_broadcastable_slice`: Internal support function, not an operator
- `aten::tf_padding`: Internal support function, not an operator

These operators were already supported previously:

- `aten::_convolution`
- `aten::adaptive_avg_pool2d`
- `aten::add`
- `aten::add_`
- `aten::addmm`
- `aten::avg_pool2d`
- `aten::batch_norm`
- `aten::cat`
- `aten::dimension_value`
- `aten::dropout`
- `aten::flatten`
- `aten::max_pool2d`
- `aten::mul`

- aten::relu\_
- aten::t
- aten::tanh
- aten::values
- prim::Constant
- prim::GetAttr
- prim::ListConstruct
- prim::ListUnpack
- prim::TupleConstruct
- prim::TupleUnpack

### PyTorch Neuron Release [1.0.672.0]

No change

### PyTorch Neuron Release [1.0.552.0]

Added support for new operators:

- aten::\_convolution
- aten::adaptive\_avg\_pool2d
- aten::add
- aten::add\_
- aten::addmm
- aten::avg\_pool2d
- aten::batch\_norm
- aten::cat
- aten::dimension\_value
- aten::dropout
- aten::flatten
- aten::max\_pool2d
- aten::mul
- aten::relu\_
- aten::t
- aten::tanh
- aten::tf\_broadcastable\_slice
- aten::tf\_padding
- aten::values

- `prim::Constant`
- `prim::GetAttr`
- `prim::ListConstruct`
- `prim::ListUnpack`
- `prim::TupleConstruct`
- `prim::TupleUnpack`

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Troubleshooting Guide for PyTorch Neuron (torch-neuron)

### General Torch-Neuron issues

If you see an error about “Unknown builtin op: neuron::forward\_1” like below, please ensure that import line “import torch\_neuron” (to register the Neuron custom operation) is in the inference script before using `torch.jit.load`.

```
Unknown builtin op: neuron::forward_1.  
Could not find any similar ops to neuron::forward_1. This op may not exist or may not be  
currently supported in TorchScript.
```

### TorchVision related issues

If you encounter an error like below, it is because latest torchvision version  $\geq 0.7$  is not compatible with Torch-Neuron 1.5.1. Please downgrade torchvision to version 0.6.1:

```
E AttributeError: module 'torch.jit' has no attribute '_script_if_tracing'
```

### 2GB protobuf limit related issues

If you encounter an error like below, it is because the model size is larger than 2GB. To compile such large models, use the `separate_weights=True` flag. Note, ensure that you have the latest version of compiler installed to support this flag. You can upgrade neuron-cc using `python3 -m pip install neuron-cc[tensorflow] -U --force-extra-index-url=https://pip.repos.neuron.amazonaws.com`

```
E google.protobuf.message.DecodeError: Error parsing message with type 'tensorflow.  
GraphDef'
```

## torch.jit.trace issues

The `/neuron-guide/neuron-frameworks/pytorch-neuron/api-compilation-python-api.rst` uses the PyTorch `torch.jit.trace()` function to generate `ScriptModule` models for execution on Inferentia. Due to that, to execute your PyTorch model on Inferentia it must be torch-jit-traceable, otherwise you need to make sure your model is torch-jit-traceable. You can try modifying your underlying PyTorch model code to make it traceable. If it's not possible to change your model code, you can *write a wrapper around your model* that makes it torch-jit-traceable to compile it for Inferentia.

Please visit `torch.jit.trace()` to review the properties that a model must have to be torch-jit-traceable. The PyTorch-Neuron trace API `torch_neuron.trace()` accepts `**kwargs` for `torch.jit.trace()`. For example, you can use the `strict=False` flag to *compile models with dictionary outputs*.

## Compiling models with outputs that are not torch-jit-traceable

To enable compilation of models with non torch-jit-traceable outputs, you can use a technique that involves writing a wrapper that converts the model's output into a form that is torch-jit-traceable. You can then compile the wrapped model for Inferentia using `torch_neuron.trace()`.

The following example uses a wrapper to compile a model with non torch-jit-traceable outputs. This model cannot be compiled for Inferentia in its current form because it outputs a list of tuples and tensors, which is not torch-jit-traceable.

```
import torch
import torch_neuron
import torch.nn as nn

class Model(nn.Module):
    def __init__(self):
        super(Model, self).__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
        a = self.conv(x) + 1
        b = self.conv(x) + 2
        c = self.conv(x) + 3
        # An output that is a list of tuples and tensors is not torch-traceable
        return [(a, b), c]

model = Model()
model.eval()

inputs = torch.rand(1, 1, 3, 3)

# Try to compile the model
model_neuron = torch_neuron.trace(model, inputs) # ERROR: This cannot be traced, we must
↪ change the output format
```

To compile this model for Inferentia, we can write a wrapper around the model to convert its outputs into a tuple of tensors, which is torch-jit-traceable.

```
class NeuronCompatibilityWrapper(nn.Module):
    def __init__(self):
        super(NeuronCompatibilityWrapper, self).__init__()
        self.model = Model()
```

(continues on next page)

(continued from previous page)

```
def forward(self, x):
    out = self.model(x)
    # An output that is a tuple of tuples and tensors is torch-jit-traceable
    return tuple(out)
```

Now, we can successfully compile the model for Inferentia using the `NeuronCompatibilityWrapper` wrapper as follows:

```
model = NeuronCompatibilityWrapper()
model.eval()

# Compile the traceable wrapped model
model_neuron = torch.neuron.trace(model, inputs)
```

If the model's outputs must be in the original form, a second wrapper can be used to transform the outputs after compilation for Inferentia. The following example uses the `OutputFormatWrapper` wrapper to convert the compiled model's output back into the original form of a list of tuples and tensors.

```
class OutputFormatWrapper(nn.Module):
    def __init__(self):
        super(OutputFormatWrapper, self).__init__()
        self.traceable_model = NeuronCompatibilityWrapper()

    def forward(self, x):
        out = self.traceable_model(x)
        # Return the output in the original format of Model()
        return list(out)

model = OutputFormatWrapper()
model.eval()

# Compile the traceable wrapped model
model.traceable_model = torch.neuron.trace(model.traceable_model, inputs)
```

## Compiling a submodule in a model that is not torch-jit-traceable

The following example shows how to compile a submodule that is part of a non torch-jit-traceable model. In this example, the top-level model `Outer` uses a dynamic flag, which is not torch-jit-traceable. However, the submodule `Inner` is torch-jit-traceable and can be compiled for Inferentia.

```
import torch
import torch_neuron
import torch.nn as nn

class Inner(nn.Module) :
    def __init__(self):
        super().__init__()
        self.conv = nn.Conv2d(1, 1, 3)

    def forward(self, x):
```

(continues on next page)

(continued from previous page)

```

        return self.conv(x) + 1

class Outer(nn.Module):
    def __init__(self):
        super().__init__()
        self.inner = Inner()

    def forward(self, x, add_offset: bool = False):
        base = self.inner(x)
        if add_offset:
            return base + 1
        return base

model = Outer()
inputs = torch.rand(1, 1, 3, 3)

# Compile the traceable wrapped submodule
model.inner = torch.neuron.trace(model.inner, inputs)

# TorchScript the model for serialization
script = torch.jit.script(model)
torch.jit.save(script, 'model.pt')

loaded = torch.jit.load('model.pt')

```

Alternatively, for usage scenarios in which the model configuration is static during inference, the dynamic flags can be hardcoded in a wrapper to make the model torch-jit-traceable and enable compiling the entire model for Inferentia. In this example, we assume the `add_offset` flag is always `True` during inference, so we can hardcode this conditional path in the Static wrapper to remove the dynamic behavior and compile the entire model for Inferentia.

```

class Static(nn.Module):
    def __init__(self):
        super().__init__()
        self.outer = Outer()

    def forward(self, x):
        # hardcoded `add_offset=True`
        output = self.outer(x, add_offset=True)
        return output

model = Static()

# We can now compile the entire model because `add_offset=True` is hardcoded in the
# Static wrapper
model_neuron = torch.neuron.trace(model, inputs)

```

This document is relevant for: Inf1

This document is relevant for: Inf1

## PyTorch Neuron (torch-neuron) release notes

**Table of contents**

- *Known Issues and Limitations - Updated 03/21/2023*
- *PyTorch Neuron release [package ver. 1.\*.\*2.11.6.0, SDK ver. 2.20.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.10.12.0, SDK ver. 2.19.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.9.74.0, SDK ver. 2.18.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.9.17.0, SDK ver. 2.16.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.9.6.0, SDK ver. 2.15.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.9.1.0, SDK ver. 2.13.0]*
- *PyTorch Neuron release [package ver. 1.\*.\*2.8.9.0, SDK ver. 2.12.0]*
- *PyTorch Neuron release [2.7.10.0]*
- *PyTorch Neuron release [2.7.1.0]*
- *PyTorch Neuron release [2.6.5.0]*
- *PyTorch Neuron release [2.5.0.0]*
- *PyTorch Neuron release [2.3.0.0]*
- *PyTorch Neuron release [2.2.0.0]*
- *PyTorch Neuron release [2.1.7.0]*
- *PyTorch Neuron release [2.0.536.0]*
- *PyTorch Neuron release [2.0.468.0]*
- *PyTorch Neuron release [2.0.392.0]*
- *PyTorch Neuron release [2.0.318.0]*
- *[1.8.1.1.5.21.0]*
- *[1.8.1.1.5.7.0]*
- *[1.8.1.1.4.1.0]*
- *[1.7.1.1.3.5.0]*
- *[1.7.1.1.2.16.0]*
- *[1.7.1.1.2.15.0]*
- *[1.7.1.1.2.3.0]*
- *[1.1.7.0]*
- *[1.0.1978.0]*
- *[1.0.1721.0]*
- *[1.0.1532.0]*
- *[1.0.1522.0]*
- *[1.0.1386.0]*



- [1.0.1168.0]
- [1.0.1001.0]
- [1.0.825.0]
- [1.0.763.0]
- [1.0.672.0]
- [1.0.627.0]

This document lists the release notes for the Pytorch-Neuron package.

## Known Issues and Limitations - Updated 03/21/2023

### Min & Max Accuracy

The index outputs of the `aten::argmin`, `aten::argmax`, `aten::min`, and `aten::max` operator implementations are sensitive to precision. For models that contain these operators and have `float32` inputs, we recommend using the `--fp32-cast=matmult --fast-math no-fast-relayout` compiler option to avoid numerical imprecision issues. Additionally, the `aten::min` and `aten::max` operator implementations do not currently support `int64` inputs when `dim=0`. For more information on precision and performance-accuracy tuning, see [Mixed precision and performance-accuracy tuning \(neuron-cc\)](#).

### Python 3.5

If you attempt to import `torch.neuron` from Python 3.5 you will see this error in 1.1.7.0 - please use Python 3.6 or greater:

```
File "/tmp/install_test_env/lib/python3.5/site-packages/torch_neuron/__init__.py", line 29
    f'Invalid dependency version torch=={torch.__version__}. '
    ^
SyntaxError: invalid syntax
```

- Torchvision has dropped support for Python 3.5
- HuggingFace transformers has dropped support for Python 3.5

### Torchvision

When versions of `torchvision` and `torch` are mismatched, this can result in exceptions when compiling `torchvision` based models. Specific versions of `torchvision` are built against each release of `torch`. For example:

- `torch==1.5.1` matches `torchvision==0.6.1`
- `torch==1.7.1` matches `torchvision==0.8.2`
- etc.

Simultaneously installing both `torch-neuron` and `torchvision` is the recommended method of correctly resolving versions.

## Dynamic Batching

Dynamic batching does not work properly for some models that use the `aten::size` operator. When this issue occurs, the input batch sizes are not properly recorded at inference time, resulting in an error such as:

```
RuntimeError: The size of tensor a (X) must match the size of tensor b (Y) at non-  
↪singleton dimension 0.
```

This error typically occurs when `aten::size` operators are partitioned to CPU. We are investigating a fix for this issue.

### PyTorch Neuron release [package ver. 1.\*.2.11.6.0, SDK ver. 2.20.0]

Date: 09/16/2024

- Minor updates.

### PyTorch Neuron release [package ver. 1.\*.2.10.12.0, SDK ver. 2.19.0]

Date: 07/03/2024

- Minor updates.

### PyTorch Neuron release [package ver. 1.\*.2.9.74.0, SDK ver. 2.18.0]

Date: 04/01/2024

- Minor updates.

### PyTorch Neuron release [package ver. 1.\*.2.9.17.0, SDK ver. 2.16.0]

Date: 12/21/2023

- Minor updates.

### PyTorch Neuron release [package ver. 1.\*.2.9.6.0, SDK ver. 2.15.0]

Date: 10/26/2023

- Minor updates.

### PyTorch Neuron release [package ver. 1.\*.2.9.1.0, SDK ver. 2.13.0]

Date: 08/28/2023

- Added support for `clamp_min/clamp_max` ATEN operators.

**PyTorch Neuron release [package ver. 1.\*.\*2.8.9.0, SDK ver. 2.12.0]**

Date: 07/19/2023

- Minor updates.

**PyTorch Neuron release [2.7.10.0]**

Date: 06/14/2023

**New in this release**

- Added support for Python 3.10

**Bug fixes**

- torch.pow Operation now correctly handles mismatch between base and exponent data types

**PyTorch Neuron release [2.7.1.0]**

Date: 05/1/2023

- Minor updates.

**PyTorch Neuron release [2.6.5.0]**

Date: 03/28/2023

**New in this release**

- Added support for torch==1.13.1
- New releases of torch-neuron no longer include versions for torch==1.7 and torch==1.8
- Added support for Neuron runtime 2.12
- Added support for new operators:
  - aten::tensordot
  - aten::adaptive\_avg\_pool1d
  - aten::prelu
  - aten::reflection\_pad2d
  - aten::baddbmm
  - aten::repeat
- Added a separate\_weights flag to `torch_neuron.trace()` to support models that are larger than 2GB

## Bug fixes

- Fixed `aten::_convolution` with grouping for:
  - `torch.nn.Conv1d`
  - `torch.nn.Conv3d`
  - `torch.nn.ConvTranspose2d`
- Fixed `aten::linear` to support 1d input tensors
- Fixed an issue where an input could not be directly returned from the network

## PyTorch Neuron release [2.5.0.0]

Date: 11/23/2022

## New in this release

- Added PyTorch 1.12 support
- Added Python 3.8 support
- Added new operators support. See *PyTorch Neuron (torch-neuron) Supported operators*
- Added support for `aten::lstm`. See: *Developer Guide - PyTorch Neuron (torch-neuron) LSTM Support*
- Improved logging:
  - Improved error messages for specific compilation failure modes, including out-of-memory errors
  - Added a warning to show the code location of `prim::PythonOp` operations
  - Removed overly-verbose tracing messages
  - Added improved error messages for `neuron-cc` and `tensorflow` dependency issues
  - Added more debug information when an invalid dynamic batching configuration is used
- Added new beta explicit NeuronCore placement API. See: `torch_neuron_core_placement_api`
- Added new guide for NeuronCore placement. See: *PyTorch Neuron (torch-neuron) Core Placement*
- Improved `torch_neuron.trace()` performance when using large graphs
- Reduced host memory usage of loaded models in `libtorchneuron.so`
- Added `single_fusion_ratio_threshold` argument to `torch_neuron.trace()` to give more fine-grained control of partitioned graphs

## Bug fixes

- Improved handling of tensor mutations which previously caused accuracy issues on certain models (i.e. yolor, yolov5)
- Fixed an issue where `inf` and `-inf` values would cause unexpected NaN values. This could occur with newer versions of `transformers`
- Fixed an issue where `torch.neuron.DataParallel()` would not fully utilize all NeuronCores for specific batch sizes
- Fixed and improved operators:
  - `aten::upsample_bilinear2d`: Improved error messages in cases where the operation cannot be supported
  - `aten::_convolution`: Added support for `output_padding` argument
  - `aten::div`: Added support for `rounding_mode` argument
  - `aten::sum`: Fixed to handle non-numeric data types
  - `aten::expand`: Fixed to handle scalar tensors
  - `aten::permute`: Fixed to handle negative indices
  - `aten::min`: Fixed to support more input types
  - `aten::max`: Fixed to support more input types
  - `aten::max_pool2d`: Fixed to support both 3-dimensional and 4-dimensional input tensors
  - `aten::Int`: Fixed an issue where long values would incorrectly lose precision
  - `aten::constant_pad_nd`: Fixed to correctly use non-0 padding values
  - `aten::pow`: Fixed to support more input types & values
  - `aten::avg_pool2d`: Added support for `count_include_pad` argument. Added support for `ceil_mode` argument if padding isn't specified
  - `aten::zero`: Fixed to handle scalars correctly
  - `prim::Constant`: Fixed an issue where `-inf` was incorrectly handled
  - Improved handling of scalars in arithmetic operators

## PyTorch Neuron release [2.3.0.0]

Date: 04/29/2022

## New in this release

- Added support PyTorch 1.11.
- Updated PyTorch 1.10 to version 1.10.2.
- End of support for torch-neuron 1.5, see eol-pt-15.
- Added support for new operators:
  - `aten::masked_fill_`
  - `aten::new_zeros`

- `aten::frobenius_norm`

### Bug fixes

- Improved `aten::gelu` accuracy
- Updated `aten::meshgrid` to support optional indexing argument introduced in `torch 1.10`, see [PyTorch issue 50276](#)

### PyTorch Neuron release [2.2.0.0]

Date: 03/25/2022

#### New in this release

- Added full support for `aten::max_pool2d_with_indices` - (Was previously supported only when indices were unused).
- Added new torch-neuron packages compiled with `-D_GLIBCXX_USE_CXX11_ABI=1`, the new packages support PyTorch 1.8, PyTorch 1.9, and PyTorch 1.10. To install the additional packages compiled with `-D_GLIBCXX_USE_CXX11_ABI=1` please change the package repo index to <https://pip.repos.neuron.amazonaws.com> (<https://pip.repos.neuron.amazonaws.com/>)/cxx11/

### PyTorch Neuron release [2.1.7.0]

Date: 01/20/2022

#### New in this release

- Added PyTorch 1.10 support
- Added new operators support, see *[PyTorch Neuron \(torch-neuron\) Supported operators](#)*
- Updated `aten::_convolution` to support 2d group convolution
- Updated `neuron::forward` operators to allocate less dynamic memory. This can increase performance on models with many input & output tensors.
- Updated `neuron::forward` to better handle batch sizes when `dynamic_batch_size=True`. This can increase performance at inference time when the input batch size is exactly equal to the traced model batch size.

### Bug fixes

- Added the ability to `torch.jit.trace` a `torch.nn.Module` where a submodule has already been traced with [torch\\_neuron.trace\(\)](#) on a CPU-type instance. Previously, if this had been executed on a CPU-type instance, an initialization exception would have been thrown.
- Fixed `aten::matmul` behavior on 1-dimensional by n-dimensional multiplies. Previously, this would cause a validation error.
- Fixed binary operator type promotion. Previously, in unusual situations, operators like `aten::mul` could produce incorrect results due to invalid casting.

- Fixed `aten::select` when index was -1. Previously, this would cause a validation error.
- Fixed `aten::adaptive_avg_pool2d` padding and striding behavior. Previously, this could generate incorrect results with specific configurations.
- Fixed an issue where dictionary inputs could be incorrectly traced when the tensor values had gradients.

### PyTorch Neuron release [2.0.536.0]

Date: 01/05/2022

#### New in this release

- Added new operator support for specific variants of operations (See *PyTorch Neuron (torch-neuron) Supported operators*)
- Added optional `optimizations` keyword to `torch_neuron.trace()` which accepts a list of *Optimization* passes.

### PyTorch Neuron release [2.0.468.0]

Date: 12/15/2021

#### New in this release

- Added support for `aten::cumsum` operation.
- Fixed `aten::expand` to correctly handle adding new dimensions.

### PyTorch Neuron release [2.0.392.0]

Date: 11/05/2021

- Updated Neuron Runtime (which is integrated within this package) to `libnrt 2.2.18.0` to fix a container issue that was preventing the use of containers when `/dev/neuron0` was not present. See details here [neuron-runtime-release-notes](#).

### PyTorch Neuron release [2.0.318.0]

Date: 10/27/2021

## New in this release

- PyTorch Neuron 1.x now support Neuron Runtime 2.x (`libnrt.so` shared library) only.

---

**Important:**

- You must update to the latest Neuron Driver (`aws-neuron-dkms` version 2.1 or newer) for proper functionality of the new runtime library.
  - Read *Introducing Neuron Runtime 2.x (`libnrt.so`)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.
  - Read *Migrate your application to Neuron Runtime 2.x (`libnrt.so`)* for detailed information of how to migrate your application.
- 

- Introducing PyTorch 1.9.1 support (support for `torch==1.9.1`)
- Added `torch_neuron.DataParallel`, see ResNet-50 tutorial [html] and *Data Parallel Inference on Torch Neuron* application note.
- Added support for tracing on GPUs
- Added support for `ConvTranspose1d`
- Added support for new operators:
  - `aten::empty_like`
  - `aten::log`
  - `aten::type_as`
  - `aten::movedim`
  - `aten::einsum`
  - `aten::argmax`
  - `aten::min`
  - `aten::argmin`
  - `aten::abs`
  - `aten::cos`
  - `aten::sin`
  - `aten::linear`
  - `aten::pixel_shuffle`
  - `aten::group_norm`
  - `aten::_weight_norm`
- Added `torch_neuron.is_available()`



## Resolved Issues

- Fixed a performance issue when using both the `dynamic_batch_size=True` trace option and `--neuron-core-pipeline` compiler option. Dynamic batching now uses OpenMP to execute pipeline batches concurrently.
- Fixed `torch_neuron.trace` issues:
  - Fixed a failure when the same submodule was traced with multiple inputs
  - Fixed a failure where some operations would fail to be called with the correct arguments
  - Fixed a failure where custom operators (torch plugins) would cause a trace failure
- Fixed variants of `aten::upsample_bilinear2d` when `scale_factor=1`
- Fixed variants of `aten::expand` using `dim=-1`
- Fixed variants of `aten::stack` using multiple different input data types
- Fixed variants of `aten::max` using indices outputs

### [1.8.1.1.5.21.0]

Date: 08/12/2021

## Summary

- Minor updates.

### [1.8.1.1.5.7.0]

Date: 07/02/2021

## Summary

- Added support for dictionary outputs using `strict=False` flag. See [/neuron-guide/neuron-frameworks/pytorch-neuron/troubleshooting-guide.rst](#).
- Updated `aten::batch_norm` to correctly implement the affine flag.
- Added support for `aten::erf` and `prim::DictConstruct`. See *PyTorch Neuron (torch-neuron) Supported operators*.
- Added dynamic batch support. See [/neuron-guide/neuron-frameworks/pytorch-neuron/api-compilation-python-api.rst](#).

**[1.8.1.1.4.1.0]**

Date: 5/28/2021

**Summary**

- Added support for PyTorch 1.8.1
  - Models compatibility
    - \* Models compiled with previous versions of PyTorch Neuron (<1.8.1) are compatible with PyTorch Neuron 1.8.1.
    - \* Models compiled with PyTorch Neuron 1.8.1 are not backward compatible with previous versions of PyTorch Neuron (<1.8.1) .
  - Updated tutorials to use Hugging Face Transformers 4.6.0.
  - Added a new set of forward operators (forward\_v2)
  - Host memory allocation when loading the same model on multiple NeuronCores is significantly reduced
  - Fixed an issue where models would not deallocate all memory within a python session after being garbage collected.
  - Fixed a TorchScript/C++ issue where loading the same model multiple times would not use multiple NeuronCores by default.
- Fixed logging to no longer configure the root logger.
- Removed informative messages that were produced during compilations as warnings. The number of warnings reduced significantly.
- Convolution operator support has been extended to include ConvTranspose2d variants.
- Reduce the amount of host memory usage during inference.

**[1.7.1.1.3.5.0]**

Date: 4/30/2021

**Summary**

- ResNext models now functional with new operator support
- YOLOv5 support refer to <https://github.com/aws/aws-neuron-sdk/issues/253> note <https://github.com/ultralytics/yolov5/pull/2953> which optimized YOLOv5 for AWS Neuron
- Convolution operator support has been extended to include most Conv1d and Conv3d variants
- New operator support. Please see *PyTorch Neuron (torch-neuron) Supported operators* for the complete list of operators.

**[1.7.1.1.2.16.0]**

Date: 3/4/2021

**Summary**

- Minor enhancements.

**[1.7.1.1.2.15.0]**

Date: 2/24/2021

**Summary**

- Fix for CVE-2021-3177.

**[1.7.1.1.2.3.0]**

Date: 1/30/2021

**Summary**

- Made changes to allow models with -inf scalar constants to correctly compile
- Added new operator support. Please see *PyTorch Neuron (torch-neuron) Supported operators* for the complete list of operators.

**[1.1.7.0]**

Date: 12/23/2020

**Summary**

- We are dropping support for Python 3.5 in this release
- `torch.neuron.trace` behavior will now throw a `RuntimeError` in the case that no operators are compiled for neuron hardware
- `torch.neuron.trace` will now display compilation progress indicators (dots) as default behavior (neuron-cc must be updated to the December release to greater to see this feature)
- Added new operator support. Please see *PyTorch Neuron (torch-neuron) Supported operators* for the complete list of operators.
- Extended the BERT pretrained tutorial to demonstrate execution on multiple cores and batch modification, updated the tutorial to accommodate changes in the Hugging Face Transformers code for version 4.0
- Added a tutorial for `torch-serve` which extends the BERT tutorial
- Added support for PyTorch 1.7

### [1.0.1978.0]

Date: 11/17/2020

#### Summary

- Fixed bugs in comparison operators, and added remaining variantes (eq, ne, gt, ge, lt, le)
- Added support for `prim::PythonOp` - note that this must be run on CPU and not Neuron. We recommend you replace this code with PyTorch operators if possible
- Support for a series of new operators. Please see *PyTorch Neuron (torch-neuron) Supported operators* for the complete list of operators.
- Performance improvements to the runtime library
- Correction of a runtime library bug which caused models with large tensors to generate incorrect results in some cases

### [1.0.1721.0]

Date: 09/22/2020

#### Summary

- Various minor improvements to the Pytorch autopartitioner feature
- Support for the operators `aten::constant_pad_nd`, `aten::meshgrid`
- Improved performance on various torchvision models. Of note are resnet50 and vgg16

### [1.0.1532.0]

Date: 08/08/2020

#### Summary

- Various minor improvements to the Pytorch autopartitioner feature
- Support for the `aten:ones` operator

### [1.0.1522.0]

Date: 08/05/2020

## Summary

Various minor improvements.

**[1.0.1386.0]**

Date: 07/16/2020

## Summary

This release adds auto-partitioning, model analysis and PyTorch 1.5.1 support, along with a number of new operators

## Major New Features

- Support for Pytorch 1.5.1
- Introduce an automated operator device placement mechanism in `torch.neuron.trace` to run sub-graphs that contain operators that are not supported by the neuron compiler in native PyTorch. This new mechanism is on by default and can be turned off by adding argument `fallback=False` to the compiler arguments.
- Model analysis to find supported and unsupported operators in a model

## Resolved Issues

**[1.0.1168.0]**

Date 6/11/2020

## Summary

## Major New Features

## Resolved Issues

## Known Issues and Limitations

**[1.0.1001.0]**

Date: 5/11/2020

### Summary

Additional PyTorch operator support and improved support for model saving and reloading.

### Major New Features

- Added Neuron Compiler support for a number of previously unsupported PyTorch operators. Please see [:ref:`neuron-cc-ops-pytorch`](#) for the complete list of operators.
- Add support for `torch.neuron.trace` on models which have previously been saved using `torch.jit.save` and then reloaded.

### Resolved Issues

### Known Issues and Limitations

[1.0.825.0]

Date: 3/26/2020

### Summary

### Major New Features

### Resolved Issues

### Known Issues and limitations

[1.0.763.0]

Date: 2/27/2020

### Summary

Added Neuron Compiler support for a number of previously unsupported PyTorch operators. Please see [PyTorch Neuron \(torch-neuron\) Supported operators](#) for the complete list of operators.

### Major new features

- None

**Resolved issues**

- None

**[1.0.672.0]**

Date: 1/27/2020

**Summary****Major new features****Resolved issues**

- Python 3.5 and Python 3.7 are now supported.

**Known issues and limitations****Other Notes****[1.0.627.0]**

Date: 12/20/2019

**Summary**

This is the initial release of torch-neuron. It is not distributed on the DLAMI yet and needs to be installed from the neuron pip repository.

Note that we are currently using a TensorFlow as an intermediate format to pass to our compiler. This does not affect any runtime execution from PyTorch to Neuron Runtime and Inferentia. This is why the neuron-cc installation must include [tensorflow] for PyTorch.

**Major new features****Resolved issues****Known issues and limitations****Models TESTED**

The following models have successfully run on neuron-inferentia systems

1. SqueezeNet
2. ResNet50
3. Wide ResNet50

### Pytorch Serving

In this initial version there is no specific serving support. Inference works correctly through Python on Inf1 instances using the neuron runtime. Future releases will include support for production deployment and serving of models

### Profiler support

Profiler support is not provided in this initial release and will be available in future releases

### Automated partitioning

Automatic partitioning of graphs into supported and non-supported operations is not currently supported. A tutorial is available to provide guidance on how to manually partition a model graph. Please see [pytorch-manual-partitioning-jn-tutorial](#)

### PyTorch dependency

Currently PyTorch support depends on a Neuron specific version of PyTorch v1.3.1. Future revisions will add support for 1.4 and future releases.

### Trace behavior

In order to trace a model it must be in evaluation mode. For examples please see [/src/examples/pytorch/resnet50.ipynb](#)

### Six pip package is required

The Six package is required for the torch-neuron runtime, but it is not modeled in the package dependencies. This will be fixed in a future release.

### Multiple NeuronCore support

If the `num-neuroncores` options is used the number of cores must be manually set in the calling shell environment variable for compilation and inference.

For example: Using the keyword argument `compiler_args=['--num-neuroncores', '4']` in the trace call, requires `NEURONCORE_GROUP_SIZES=4` to be set in the environment at compile time and runtime

### CPU execution

At compilation time a constant output is generated for the purposes of tracing. Running inference on a non neuron instance will generate incorrect results. This must not be used. The following error message is generated to stderr:

```
Warning: Tensor output are ** NOT CALCULATED ** during CPU execution and only
indicate tensor shape
```



## Other notes

- Python version(s) supported:
  - 3.6
- Linux distribution supported:
  - DLAMI Ubuntu 18 and Amazon Linux 2 (using Python 3.6 Conda environments)
  - Other AMIs based on Ubuntu 18
  - For Amazon Linux 2 please install Conda and use Python 3.6 Conda environment

*This document is relevant for: **Inf1***

- *PyTorch Neuron (torch-neuron) Supported operators*
- *Troubleshooting Guide for PyTorch Neuron (torch-neuron)*
- *PyTorch Neuron (torch-neuron) release notes*

*This document is relevant for: **Inf1***

Setup (torch-neuron)

## Tutorials (torch-neuron)

### Computer Vision Tutorials

- ResNet-50 tutorial [\[html\]](#) [\[notebook\]](#)
- PyTorch YOLOv4 tutorial [\[html\]](#) [\[notebook\]](#)

### Natural Language Processing (NLP) Tutorials

- HuggingFace pretrained BERT tutorial [\[html\]](#) [\[notebook\]](#)
- HuggingFace pretrained BERT tutorial with shared weights [\[html\]](#) [\[notebook\]](#)
- Bring your own HuggingFace pretrained BERT container to Sagemaker Tutorial [\[html\]](#) [\[notebook\]](#)
- LibTorch C++ tutorial [\[html\]](#)
- TorchServe tutorial [\[html\]](#)
- HuggingFace MarianMT tutorial [\[html\]](#) [\[notebook\]](#)

### Utilizing Neuron Capabilities Tutorials

- BERT TorchServe tutorial [\[html\]](#)
- NeuronCore Pipeline tutorial [\[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
-

### Additional Examples (torch-neuron)

- [AWS Neuron Samples GitHub Repository](#)

### API Reference Guide (torch-neuron)

- *PyTorch Neuron trace Python API*
- *torch.neuron.DataParallel API*
- *PyTorch Neuron (torch-neuron) Core Placement API [Beta]*

### Developer Guide (torch-neuron)

- *Running Inference on Variable Input Shapes with Bucketing*
- *Data Parallel Inference on PyTorch Neuron*
- *Developer Guide - PyTorch Neuron (torch-neuron) LSTM Support*
- *PyTorch Neuron (torch-neuron) Core Placement*

### Misc (torch-neuron)

- *PyTorch Neuron (torch-neuron) Supported operators*
- *Troubleshooting Guide for PyTorch Neuron (torch-neuron)*
- *PyTorch Neuron (torch-neuron) release notes*

*This document is relevant for:* Inf1

*This document is relevant for:* Trn1, Trn2

### 2.1.4 Training (torch-neuronx)

*This document is relevant for:* Trn1, Trn2

### Tutorials for Training(torch-neuronx)

*This document is relevant for:* Trn1, Trn2

### Hugging Face BERT Pretraining Tutorial (Data-Parallel)

This tutorial explains how to run Hugging Face BERT-Large model pretraining on Trainium using PyTorch Neuron and data-parallel mode.

The Hugging Face BERT pretraining example demonstrates the steps required to perform single-node, multi-accelerator PyTorch model training using the new AWS EC2 Trn1 (Trainium) instances and the AWS Neuron SDK. This tutorial is an adaptation of an existing [BERT example](#) with the following important characteristics:

- Framework: PyTorch/XLA
- Model: Hugging Face BertForPreTraining
- Optimizer: AdamW, LAMB (Layerwise Adaptive Moments optimizer)

- Scheduler: Hugging Face's `get_linear_schedule_with_warmup`
- Allreduce occurs before optimizer step, after gradient accumulations (following DeepSpeed's Smart Gradient Accumulation)
- Training data types: Float32, full BFloat16 and Stochastic Rounding (SR), full BFloat16 with fp32 copy of weights, PyTorch Autocast (Automatic Mixed Precision or AMP)

As done in the original BERT paper, BERT pretraining happens in two phases. In the first phase (phase 1) BERT maximum sequence length is fixed at 128 tokens, while in phase 2 it is fixed at 512 tokens.

Neuron provides access to Trainium devices through an extension of PyTorch/XLA - a library that includes the familiar PyTorch interface along with XLA-specific additions. For additional details relating to PyTorch/XLA, please refer to the [official PyTorch/XLA documentation](#).

## Table of Contents

- *Phase 1 BFloat16 BERT-Large pretraining with AdamW and stochastic rounding*
  - *Setting up the training environment on trn1.32xlarge*
  - *Downloading tokenized and sharded dataset files*
  - *Number of workers*
  - *BFloat16 and stochastic rounding in phase 1*
  - *Pre-compilation*
  - *Initiating a Training Job*
  - *Monitoring Progress of the Training Job*
  - *Monitoring Training Job Progress using neuron-top*
  - *Monitoring Training Job Progress using TensorBoard*
  - *Finishing the tutorial*
- *Phase 1 BERT-Large pretraining with Layerwise Adaptive Moments based optimizer (LAMB)*
- *Phase 1 BFloat16 BERT-Large pretraining with AdamW and FP32 copy of weights*
- *Phase 1 BERT-Large pretraining with AdamW and PyTorch Autocast (Automatic Mixed Precision or AMP)*
- *Phase 1 BERT-Large pretraining on two instances*
- *Phase 2 BERT-Large pretraining*
  - *Training Environment*
  - *Initiating a Training Job*
- *Tools*
  - *neuron-ls*
  - *neuron-top*
  - *Generating tokenized and sharded dataset files*
- *Known issues and limitations*
  - *BERT-large compilation limitations*
  - *BERT-large pretraining with pretokenized dataset hangs when using xm.save*

- *BERT-large two worker pretraining hangs or run out of host memory during checkpointing on trn1.2xlarge*
- *BERT precompilation using neuron\_parallel\_compile hangs when using torchrun*
- *Troubleshooting*
  - *ModuleNotFoundError: No module named 'torch', 'torch\_xla', 'transformers', etc*

---

**Note:** Logs used in tutorials do not present latest performance numbers

For latest performance numbers visit [Neuron performance](#)

---

## Phase 1 BFloat16 BERT-Large pretraining with AdamW and stochastic rounding

### Setting up the training environment on trn1.32xlarge

The BERT training script `dp_bert_large_hf_pretrain_hdf5.py` ([source](#)) can run on a Trainium instance (trn1.32xlarge) that contains the appropriate Neuron runtime and Python dependencies.

First, on a trn1.32xlarge instance, follow the installation instructions at:

*Install PyTorch Neuron on Trn1*

Please set the storage of instance to *512GB* or more if you intent to run multiple experiments and save many checkpoints.

For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

Next, clone the [AWS Neuron Samples repository](#) and install requirements in the BERT tutorial directory `aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain` ([directory link](#)):

```
cd ~/
git clone https://github.com/aws-neuron/aws-neuron-samples.git
```

```
python3 -m pip install -r ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_
↪pretrain/requirements.txt
```

### Downloading tokenized and sharded dataset files

To download the tokenized and sharded dataset files needed for this tutorial, please run the following commands:

```
mkdir -p ~/examples_datasets/
pushd ~/examples_datasets/
aws s3 cp --no-progress s3://neuron-s3/training_datasets/bert_pretrain_wikicorpus_
↪tokenized_hdf5/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar . --no-sign-
↪request
tar -xf bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar
rm bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar
aws s3 cp --no-progress s3://neuron-s3/training_datasets/bert_pretrain_wikicorpus_
```

(continues on next page)

(continued from previous page)

```

↪tokenized_hdf5/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512.tar . --no-sign-
↪request
tar -xf bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512.tar
rm bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512.tar
popd

```

~/examples\_datasets/bert\_pretrain\_wikicorpus\_tokenized\_hdf5\_seqlen128 will now have the tokenized and sharded dataset files for phase 1 pretraining and ~/examples\_datasets/bert\_pretrain\_wikicorpus\_tokenized\_hdf5\_seqlen512 for phase 2 pretraining.

## Number of workers

You will be using `torchrun` (PyTorch's [Elastic Launch](#)) to run some of the commands in this tutorial. When running the training script, you can configure the number of NeuronCores to use for training by using `torchrun`'s `--nproc_per_node` option. In this tutorial, we use 32 NeuronCores on `trn1.32xlarge`.

---

**Note:** Currently Neuron Runtime only support 1 and 2 worker configurations on `trn1.2xlarge` and 1, 2, 8, and 32-worker configurations on `trn1.32xlarge`.

---

## BFloat16 and stochastic rounding in phase 1

Phase 1 pretraining performance can be increased by using BFloat16 casting and stochastic rounding. BFloat16 casting and stochastic rounding can be enabled by moving the model to BFloat16 using `model.to(torch.bfloat16)` expression in the training code and setting the environment variable `NEURON_RT_STOCHASTIC_ROUNDING_EN=1`, both are done in BERT pretraining example `dp_bert_large_hf_pretrain_hdf5.py` by default. Also in the BERT pretraining example, the loss is kept in FP32 to ensure smooth loss curve when loss averaging is used. We also preserve the optimizer states in FP32 using a modified HuggingFace AdamW implementation in order to match FP32 loss with BFloat16. To achieve maximum performance while maintaining loss convergence characteristics, we are using batch size of 16 and gradient accumulation microsteps of 32 to maintain global batch size of 16384 for phase 1. The batch size and gradient accumulation microstep changes can be set by launching the BERT pretraining script `dp_bert_large_hf_pretrain_hdf5.py` with command-line arguments `--batch_size=16 --grad_accum_usteps=32`, as seen in the following steps.

Another option with BFloat16 using PyTorch AutoCast (Automatic Mixed Precision or AMP) is covered in [Phase 1 BERT-Large pretraining with AdamW and PyTorch Autocast \(Automatic Mixed Precision or AMP\)](#).

---

**Note:** `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated starting in `torch-xla` 2.1, and their usage would result in warnings. They will become no-operations in `torch-xla` 2.6. Please switch to using `model.to(torch.bfloat16())` or AMP.

---

## Pre-compilation

PyTorch Neuron evaluates operations lazily during execution of the training loops, which means it builds a symbolic graph in the background and the graph is executed in hardware only when the tensor is printed, transferred to CPU, or `xm.mark_step()` is encountered (`xm.mark_step()` is implicitly called by `pl.MpDeviceLoader/pl.ParallelLoader`). During execution of the training loops, PyTorch Neuron can build multiple graphs depending on the number of conditional paths taken. For BERT-Large pretraining, PyTorch Neuron builds multiple unique graphs that should be compiled before running on the NeuronCores. PyTorch Neuron will compile those graphs only if they are not in the XLA in-memory cache or the persistent cache. To reduce the compilation time of these graphs, you can pre-compile those graphs using the utility `neuron_parallel_compile` (provided by the `libneuronxla` package, a transitive dependency of `torch-neuronx`) as shown:

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
neuron_parallel_compile torchrun --nproc_per_node=32 \
dp_bert_large_hf_pretrain_hdf5.py \
--steps_this_run 10 \
--batch_size 16 \
--grad_accum_usteps 32 | tee compile_log.txt
```

This command performs a fast trial run of the training script to build graphs and then do parallel compilations on those graphs using multiple processes of Neuron Compiler before populating the on-disk persistent cache with compiled graphs. This helps make the actual training run faster because the compiled graphs will be loaded from the persistent cache. Currently it takes ~13 minutes to compile the BERT-Large model training step using the pre-compilation script (compare to ~40 minutes if not using the pre-compilation script). Note that the command above specifies 32 NeuronCores for `trn1.32xlarge` via `--nproc_per_node` option.

The script `run_dp_bert_large_hf_pretrain_bf16_s128.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script using `neuron_parallel_compile ./run_dp_bert_large_hf_pretrain_bf16_s128.sh` to start the precompilation.

The pretokenized dataset is expected to be at `~/examples_datasets/bert_pretrain_wiki_corpus_tokenized_hdf5_seqlen128/` by default (see above for downloading instructions) and can be changed via the `--data_dir` option.

---

**Note:** The trial run during pre-compilation currently outputs invalid loss numbers. Please disregard them.

---

---

**Note:** The command after `neuron_parallel_compile` should match the actual run command, except for the option `--steps_this_run` which shortens the trial run just enough to allow the tool to build all the graphs needed for the actual run.

---

If you interrupt the run and restart the execution without changing model configurations or training hyperparameters, the new run will detect the cached graphs in the persistent cache (on-disk) and reload the compiled graphs for execution, avoiding any recompilation time.

Changes made to the BERT model configuration (layers, hidden size, attention heads in the `get_model` function), batch size (using `--batch_size` option), optimizer or number of workers may trigger graph recompilation. It is best to rerun the pre-compilation step above if these changes are made.

You can adjust the following hyperparameters without changing the model and causing recompilation:

- Number of global steps to run (`--steps_this_run` option)
- Learning rate (`--lr` option)
- Gradient accumulation steps  $> 1$  (`--grad_accum_usteps` option). If 1 then there's no gradient accumulation and the graphs change causing recompilation.

## Initiating a Training Job

After running the pre-compilation step, continue with the actual phase 1 pretraining by running the following set of commands to launch 32 data parallel distributed training workers on trn1.32xlarge:

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 \
dp_bert_large_hf_pretrain_hdf5.py \
--batch_size 16 \
--grad_accum_usteps 32 | tee run_pretrain_log.txt
```

The script `run_dp_bert_large_hf_pretrain_bf16_s128.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script to start the training.

The following messages indicate that the Neuron Runtime is initializing:

```
Using Neuron Runtime
Using Neuron Runtime
Using Neuron Runtime
Using Neuron Runtime
Using Neuron Runtime
...
```

A few moments later, you will see the Training Configuration and Model Configuration in the output:

```
-----TRAINING CONFIG-----
Namespace(batch_size=16, data_dir='~/examples_datasets/
bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128/', debug=False,
enable_pt_autocast=False, grad_accum_usteps=32, local_rank=0, lr=0.0004,
max_pred_len=20, max_steps=28125, metrics_file='/tmp/test_dict.json',
minimal_ckpt=False, num_ckpts_to_keep=1, output_dir='./output',
phase1_end_step=28125, phase2=False, resume_ckpt=False, resume_step=-1,
seed=12349, seq_len=128, shards_per_ckpt=1, steps_this_run=28125, warmup_steps=2000)
```

```
-----MODEL CONFIG-----
BertConfig {
  "_name_or_path": "bert-large-uncased",
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "classifier_dropout": null,
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 1024,
  "initializer_range": 0.02,
  "intermediate_size": 4096,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 16,
  "num_hidden_layers": 24,
  "pad_token_id": 0,
```

(continues on next page)

(continued from previous page)

```

"position_embedding_type": "absolute",
"transformers_version": "4.15.0",
"type_vocab_size": 2,
"use_cache": true,
"vocab_size": 30522
}

```

As the worker processes begin training on the BERT dataset, you will begin to see training metrics and the learning rate logged to the console approximately every training step. The metrics include average\_loss, step\_loss, and throughput:

```

LOG Thu Sep 29 22:30:10 2022 - (0, 78) step_loss : 9.1875  learning_rate : 1.56e-05  ↵
↪ throughput : 2873.14
LOG Thu Sep 29 22:30:16 2022 - (0, 79) step_loss : 8.9375  learning_rate : 1.58e-05  ↵
↪ throughput : 2878.09
LOG Thu Sep 29 22:30:22 2022 - (0, 80) step_loss : 9.0000  learning_rate : 1.60e-05  ↵
↪ throughput : 2875.31
LOG Thu Sep 29 22:30:27 2022 - (0, 81) step_loss : 9.0000  learning_rate : 1.62e-05  ↵
↪ throughput : 2877.35
LOG Thu Sep 29 22:30:33 2022 - (0, 82) step_loss : 8.8750  learning_rate : 1.64e-05  ↵
↪ throughput : 2872.55
LOG Thu Sep 29 22:30:39 2022 - (0, 83) step_loss : 9.0000  learning_rate : 1.66e-05  ↵
↪ throughput : 2876.17
LOG Thu Sep 29 22:30:44 2022 - (0, 84) step_loss : 9.1250  learning_rate : 1.68e-05  ↵
↪ throughput : 2872.48
LOG Thu Sep 29 22:30:50 2022 - (0, 85) step_loss : 9.0000  learning_rate : 1.70e-05  ↵
↪ throughput : 2873.39

```

By default, the training script will store all output files under `~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain/output`. The output files consist of the following:

- PyTorch model checkpoint files, with names containing the global step of the checkpoint (ckpt\_2000.pt, ckpt\_4000.pt, etc.). Currently, the training script saves a checkpoint after every dataset shard. The frequency of saving checkpoint can be reduced by increasing the number of dataset shards per checkpoint, using option `--shards_per_ckpt`. Furthermore, the number of checkpoints kept at a given time is limited by `--num_ckpts_to_keep` option (currently default to 1).
- TensorBoard log files (each training run will store its logs in a subdirectory with prefix `neuron_tblogs_`).

## Monitoring Progress of the Training Job

Using a single Trn1 instance with 32 NeuronCores, the current BERT phase 1 pretraining will finish in about 45 hours. During this time, you will see the average loss metric begin at about 11.2 and ultimately converge to about 1.4.



## Monitoring Training Job Progress using neuron-top

With the training job still running, launch a second SSH connection into the trn1 instance, and use the `neuron-top` command to examine the aggregate NeuronCore utilization. If you have not modified the `--nproc_per_node` option in the run command, you should observe that all 32 NeuronCores are participating in the training job, with utilization fluctuating around 80%.

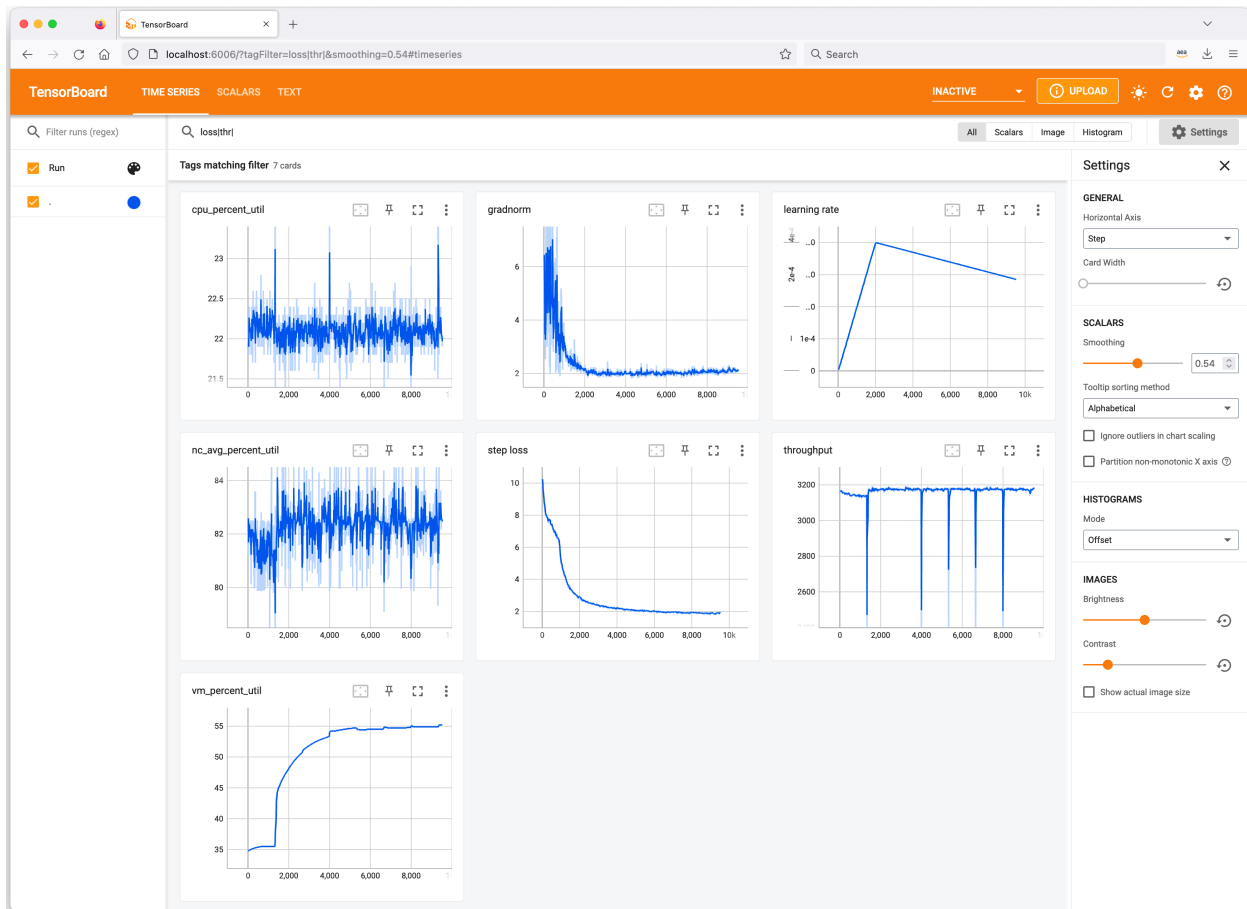
## Monitoring Training Job Progress using TensorBoard

The demo includes TensorBoard-compatible logging, which allows the learning rate and training metrics to be monitored in real-time. By default, the training script logs metrics to the following TensorBoard log directory `~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain/output/neuron_tblogs_<date/time>_<training configs>`.

In order to view your training metrics in TensorBoard, first run the following commands in your SSH session:

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
tensorboard --logdir ./output
```

Once running, open a new SSH connection to the instance and port-forward TCP port 6006 (ex: `ssh -L 6006:127.0.0.1:6006 user_name@remote_ip`). Once the tunnel is established, TensorBoard can then be accessed via web browser at the following URL: <http://localhost:6006>. Please note that you will not be able to access TensorBoard if you disconnect your port-forwarding SSH session to the Trainium instance.



## Finishing the tutorial

Once you are ready, there are a couple of options for finishing the BERT pretraining demo:

1. **Allow the training script to run to completion.** If you would like to observe the training script run to completion, it is recommended to launch the training script from a terminal multiplexer such as `tmux` or `screen`, and then detach the session so that the training script can run in the background. With this approach, you can safely let the training script run unattended, without risk of an SSH disconnection causing the training job to stop running.
2. **Stop the training job early.** To stop the training job early, press CTRL-C in the terminal window in which you launched the training script. In some cases, if you manually cancel a job using CTRL-C and then later want to run the job again, you might first need to execute `sudo rmmmod neuron`; `sudo modprobe neuron` in order to reload/reset the Neuron driver.

## Phase 1 BERT-Large pretraining with Layerwise Adaptive Moments based optimizer (LAMB)

Sometimes, to reduce the training wall time, you can use higher learning rate and larger global batch size. The approach is discussed in [LARGE BATCH OPTIMIZATION FOR DEEP LEARNING: TRAINING BERT IN 76 MINUTES](#). Tranium supports LAMB, and in this tutorial, we use publicly available XLA-friendly LAMB implementation from <https://github.com/rwightman/pytorch-image-models/blob/master/timm/optim/lamb.py>.

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 \
dp_bert_large_hf_pretrain_hdf5.py \
--max_steps 7032 \
--batch_size 8 \
--optimizer LAMB \
--lr 6e-3 \
--grad_accum_usteps 256 | tee run_pretrain_log.txt
```

The command-line argument `--optimizer LAMB` is needed, otherwise, the default optimizer AdamW will be used. Besides, you need to use a set of hyper-parameters supporting the larger global batch size (GBS). In this case, we have 64k as GBS for LAMB and use a set of hyper-params similar to <https://github.com/NVIDIA/DeepLearningExamples/blob/master/PyTorch/LanguageModeling/BERT/README.md>. Given higher GBS from LAMB than AdamW, it takes fewer steps (roughly 7k) to achieve similar level of accuracy as AdamW, which takes more than 28k steps. In addition, you can also use different data types on top of LAMB. Below is an example using the BFloat16 and Stochastic Roundings.

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 \
dp_bert_large_hf_pretrain_hdf5.py \
--max_steps 7032 \
--batch_size 16 \
--optimizer LAMB \
--lr 6e-3 \
--grad_accum_usteps 128 | tee run_pretrain_log.txt
```

The script `run_dp_bert_large_hf_pretrain_bf16_s128_lamb.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script to start the training.

## Phase 1 BFloat16 BERT-Large pretraining with AdamW and FP32 copy of weights

BFloat16 training can be achieved without stochastic rounding when a copy of weights is kept in FP32. To train BERT-Large with AdamW and FP32 copy of weights, specify `--optimizer=AdamW_FP32ParamsCopy` option when calling the BERT pretraining script (stochastic rounding is off):

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 dp_bert_large_hf_pretrain_hdf5.py \
--batch_size 16 \
--optimizer=AdamW_FP32ParamsCopy \
--grad_accum_usteps 32 |& tee run_pretrain_log.txt
```

The script `run_dp_bert_large_hf_pretrain_bf16_s128.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script with `fp32paramscopy` option like `./run_dp_bert_large_hf_pretrain_bf16_s128.sh fp32paramscopy` to start the training with FP32 copy of weights.

## Phase 1 BERT-Large pretraining with AdamW and PyTorch Autocast (Automatic Mixed Precision or AMP)

Besides the *BFloat16 and stochastic rounding in phase 1*, you can also use [PyTorch Autocast for XLA (Automatic Mixed Precision or AMP)](<https://github.com/pytorch/xla/blob/master/docs/source/perf/amp.md>), which automatically converts operations to either a lower precision (like Bfloat16) or Float32. This generally provides better performance over full Float32 due to higher compute density and lower memory footprint (`trn1_training_perf`). With the BERT-Large pretraining scripts you can use AMP by specifying the `--enable_pt_autocast` option without enabling stochastic rounding (`NEURON_RT_STOCHASTIC_ROUNDING_EN` is not set).

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 dp_bert_large_hf_pretrain_hdf5.py \
--batch_size 16 \
--enable_pt_autocast \
--grad_accum_usteps 32 | tee run_pretrain_log.txt
```

The script `run_dp_bert_large_hf_pretrain_bf16_s128.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script with `amp` option like `./run_dp_bert_large_hf_pretrain_bf16_s128.sh amp` to start the training with AMP.

Under the hood, `--enable_pt_autocast` would wrap only the forward pass and loss in the PyTorch autocasting context. The backward pass is NOT in the PyTorch autocasting context. This converts compute operations such as matrix multiply, convolution, activation, and pooling to lower precision such as BFloat16 while keeping numerically sensitive operations such as softmax and cross-entropy in Float32. For information about operations that are autocasted, please see [PyTorch Autocast for XLA AMP guide](<https://github.com/pytorch/xla/blob/master/docs/source/perf/amp.md#supported-operators>).

```
with torch.autocast(enabled=flags.enable_pt_autocast, dtype=torch.bfloat16, device_type=
    ↪ 'xla'):
    outputs = model(input_ids=input_ids,
                    attention_mask=input_mask,
                    token_type_ids=segment_ids,
                    labels=masked_lm_labels,
                    next_sentence_label=next_sentence_labels)
    loss = outputs.loss / flags.grad_accum_usteps
```

(continues on next page)

(continued from previous page)

```
loss.backward()
running_loss += loss.detach()
```

## Phase 1 BERT-Large pretraining on two instances

If you have two trn1.32xlarge instances with EFA-enabled interfaces, using [EFA-enabled security group](#), and setup using Install PyTorch Neuron on Trn1, you can run multi-instance BERT-Large pretraining. The following example demonstrate running BERT phase 1 pretraining on two instances. To ensure that the global batch size remains at 16384 for phase 1, the gradient accumulation microstep count is reduced by half when the number of instances is 2. NOTE: To run on multiple instances, you will need to use trn1.32xlarge instances and using all 32 NeuronCores on each instance.

On the rank-0 Trn1 host (root), run with `--node_rank=0` using `torchrun` utility, and `--master_addr` set to rank-0 host's IP address:

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
export FI_EFA_USE_DEVICE_RDMA=1
export FI_PROVIDER=efa
export BUCKET_CAP_MB=512
export XLA_TRANSFER_SEED_ASYNC=1
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=0 --master_addr=<root IP> --master_
↪port=2020 \
dp_bert_large_hf_pretrain_hdf5.py \
--batch_size 16 \
--grad_accum_usteps 16 |& tee run_pretrain_log.txt
```

On another Trn1 host, run with `--node_rank=1`, and `--master_addr` also set to rank-0 host's IP address:

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
export FI_EFA_USE_DEVICE_RDMA=1
export FI_PROVIDER=efa
export BUCKET_CAP_MB=512
export XLA_TRANSFER_SEED_ASYNC=1
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=1 --master_addr=<root IP> --master_
↪port=2020 \
dp_bert_large_hf_pretrain_hdf5.py \
--batch_size 16 \
--grad_accum_usteps 16 |& tee run_pretrain_log.txt
```

It is important to launch rank-0 worker with `--node_rank=0` to avoid hang.

To train on multiple instances, it is recommended to use a `ParallelCluster`. For a `ParallelCluster` example, please see [Train a model on AWS Trn1 ParallelCluster](#).

## Phase 2 BERT-Large pretraining

As mentioned above, BERT pretraining happens in two phases. In phase 1, the sequence length is 128. In phase 2, the sequence length increases to 512. This additional training phase will further reduce the pretraining loss and improve the metrics for the fine-tune tasks that usually follow. The setup is very similar to the phase 1, with some differences in training environment and command line arguments highlighted below.

### Training Environment

The following dataset and checkpoint are required:

- `~/examples_datasets/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512` is WikiCorpus training dataset that is preprocessed (tokenized and pre-masked) for phase 2.
- `~/examples/dp_bert_hf_pretrain/output/ckpt_<phase1_end_step>.pt` is the final checkpoint from phase 1. It's generated automatically at the end of phase 1 pretraining. For convenience, one can also download the example available at `s3://neuron-s3/training_checkpoints/pytorch/dp_bert_large_hf_pretrain/ckpt_28125.pt`, which is collected after 28125 training steps in phase 1. Phase 2 will continue training by loading this checkpoint. During its progression, phase 2 continues to generate its own checkpoints in output directory, following the naming convention `ckpt_<global_steps>.pt`

### Initiating a Training Job

To launch the phase 2 pretraining job with AdamW optimizer, run the same python script `dp_bert_large_hf_pretrain_hdf5.py` as before except with different options for phase 2. For phase 2, we are using global batch size of 32768, with worker device batch size of 2 and gradient accumulation microsteps of 512. The pretokenized dataset is expected to be at `~/examples_datasets/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512/` following the setup steps above and is set via `--data_dir` option.

```
cd ~/aws-neuron-samples/torch-neuronx/training/dp_bert_hf_pretrain
torchrun --nproc_per_node=32 dp_bert_large_hf_pretrain_hdf5.py \
  --data_dir ~/examples_datasets/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512/ \
  --lr 2.8e-4 \
  --phase2 \
  --resume_ckpt \
  --phase1_end_step 28125 \
  --batch_size 2 \
  --grad_accum_usteps 512 \
  --seq_len 512 \
  --max_pred_len 80 \
  --warmup_steps 781 \
  --max_steps 1563 \
  | tee run_pretrain_log_phase2.txt
```

The script `run_dp_bert_large_hf_pretrain_bf16_s512_phase2.sh` is provided in the same BERT tutorial directory for convenience and you can simply run the script to start the training with AdamW optimizer. Similarly, you can use LAMB optimizer using the script `run_dp_bert_large_hf_pretrain_bf16_s512_lamb_phase2.sh`.

The output below is expected as the job is initiated. Step 28125 is the `phase1_end_step` in this run, which could be different if phase1 training stops at a different global step.

```

Worker 21 resuming from checkpoint ./output/ckpt_28125.pt at step 28125
Worker 23 resuming from checkpoint ./output/ckpt_28125.pt at step 28125
Worker 27 resuming from checkpoint ./output/ckpt_28125.pt at step 28125
Worker 26 resuming from checkpoint ./output/ckpt_28125.pt at step 28125
Worker 20 resuming from checkpoint ./output/ckpt_28125.pt at step 28125
Worker 22 resuming from checkpoint ./output/ckpt_28125.pt at step 28125

-----TRAINING CONFIG-----
Namespace(batch_size=2, data_dir='/home/ec2-user/examples_datasets/
bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512/', debug=False,
enable_pt_autocast=False, grad_accum_usteps=512, local_rank=0, lr=0.0002,
max_pred_len=80, max_steps=28125, metrics_file='/tmp/test_dict.json',
minimal_ckpt=False, num_ckpts_to_keep=1, output_dir='./output',
phase1_end_step=28125, phase2=True, resume_ckpt=True, resume_step=-1,
seed=12349, seq_len=512, shards_per_ckpt=1, steps_this_run=32, warmup_steps=781)

-----MODEL CONFIG-----
BertConfig {
  "_name_or_path": "bert-large-uncased",
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "classifier_dropout": null,
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 1024,
  "initializer_range": 0.02,
  "intermediate_size": 4096,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 16,
  "num_hidden_layers": 24,
  "pad_token_id": 0,
  "position_embedding_type": "absolute",
  "transformers_version": "4.15.0",
  "type_vocab_size": 2,
  "use_cache": true,
  "vocab_size": 30522
}

```

As the phase 2 training proceeds, similar metrics to phase 1 will appear on the console, showing the loss, learning rate, and throughput:

```

LOG Tue Sep 27 20:56:35 2022 - (0, 26) step_loss : 4.3438 learning_rate : 6.66e-06 ↵
↵throughput : 494.55
LOG Tue Sep 27 20:57:40 2022 - (0, 27) step_loss : 4.0938 learning_rate : 6.91e-06 ↵
↵throughput : 495.67
LOG Tue Sep 27 20:58:46 2022 - (0, 28) step_loss : 4.1875 learning_rate : 7.17e-06 ↵
↵throughput : 496.18
LOG Tue Sep 27 20:59:53 2022 - (0, 29) step_loss : 4.0000 learning_rate : 7.43e-06 ↵

```

(continues on next page)

(continued from previous page)

```

↪throughput : 495.31
LOG Tue Sep 27 21:00:58 2022 - (0, 30) step_loss : 4.2500 learning_rate : 7.68e-06 ↪
↪throughput : 495.60
LOG Tue Sep 27 21:02:05 2022 - (0, 31) step_loss : 4.3125 learning_rate : 7.94e-06 ↪
↪throughput : 495.50
LOG Tue Sep 27 21:03:10 2022 - (0, 32) step_loss : 4.4688 learning_rate : 8.19e-06 ↪
↪throughput : 496.02

```

## Tools

While running the tutorial, try experimenting with the following Neuron tools, which help monitor and evaluate compute utilization in real-time:

### neuron-ls

The `neuron-ls` command describes the number of Neuron devices present in the system, along with the associated NeuronCore count, memory, and PCI device information:

```

[ec2-user@ip-10-0-1-68 ~]$ neuron-ls
Neuron Tools have changed with Neuron v1.16.0. Go here to learn more:
https://awsdocs-neuron.readthedocs-hosted.com/en/latest/neuron-guide/neuron-tools/

If you're still running Neuron Runtime 1.0, run this program as sudo to display neuron-rt
runtime information

+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | PCI   |
| DEVICE | CORES  | MEMORY | BDF   |
+-----+-----+-----+-----+
| 0       | 2       | 32 GB  | 10:1c.0 |
| 1       | 2       | 32 GB  | 20:1b.0 |
| 2       | 2       | 32 GB  | 10:1d.0 |
| 3       | 2       | 32 GB  | 20:1c.0 |
| 4       | 2       | 32 GB  | a0:1d.0 |
| 5       | 2       | 32 GB  | 90:1c.0 |
| 6       | 2       | 32 GB  | a0:1c.0 |
| 7       | 2       | 32 GB  | 90:1b.0 |
| 8       | 2       | 32 GB  | 10:1e.0 |
| 9       | 2       | 32 GB  | 20:1d.0 |
| 10      | 2       | 32 GB  | 10:1b.0 |
| 11      | 2       | 32 GB  | 20:1e.0 |
| 12      | 2       | 32 GB  | a0:1b.0 |
| 13      | 2       | 32 GB  | 90:1e.0 |
| 14      | 2       | 32 GB  | a0:1e.0 |
| 15      | 2       | 32 GB  | 90:1d.0 |
+-----+-----+-----+-----+
[ec2-user@ip-10-0-1-68 ~]$

```

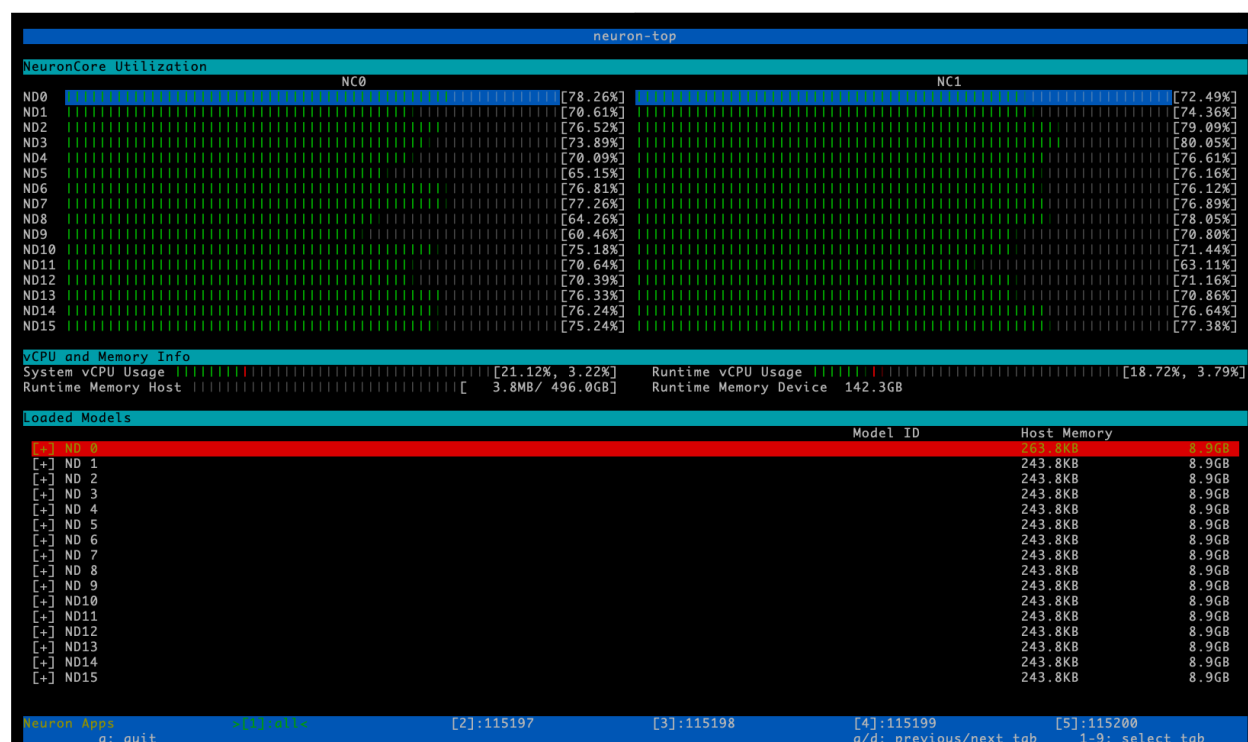
You will find that the Trn1 instance has 16 Neuron devices, each with 2 NeuronCores. This configuration allows you to train the model using a total of 32 workers, one per NeuronCore, within a single instance.

Additional information regarding `neuron-ls` can be found in the [neuron-ls user guide](#).

## neuron-top

The `neuron-top` command presents a high-level view of the Neuron environment, including the utilization of each of the NeuronCores, any models that are currently loaded onto one or more NeuronCores, process IDs for any processes that are leveraging the Neuron runtime, and basic system statistics relating to vCPU and memory usage.

Please note that `neuron-top` can either display aggregate NeuronCore utilization for ‘all’ processes (the default), or alternatively display the NeuronCore utilization for a particular process. You can toggle through the aggregate and per-process views using the `a` and `d` keys. The screenshot below illustrates the default aggregate view:



Please refer to the `neuron-top` user guide for additional details.

## Generating tokenized and sharded dataset files

This section is for generating tokenized and sharded dataset files from WikiCorpus dataset. If you just want the pre-generated dataset files, please see [Downloading tokenized and sharded dataset files](#) section above.

On a `c5n.18xlarge` instance launched with Deep Learning Conda AMI and 512GB disk space, you can generate the preprocessed datasets from WikiCorpus dataset using NVidia’s `DeepLearningExamples` for BERT pretraining. The preprocessing converts the WikiCorpus dataset to tokenized data and shard the data into multiple shards for parallel loading. The full flow takes about 8.7 hours:

```
source activate pytorch_latest_p37
cd ~/
git clone https://github.com/NVIDIA/DeepLearningExamples.git
cd DeepLearningExamples
git checkout 81b9010096b6f9812e3977b607669f6ec8b16561
sudo mkdir -m a=rwx /workspace
cp -rf PyTorch/LanguageModeling/BERT /workspace/bert
cd /workspace
```

(continues on next page)



(continued from previous page)

```
git clone https://github.com/attardi/wikiextractor.git
cd wikiextractor
git checkout 6408a430fc504a38b04d37ce5e7fc740191dee16
cd /workspace/bert
# increase num processes and shards
ex -s "+%s/(bertPrep\.py)\( --action create_hdf5_files\)/\1 --n_processes 32 --n_test_
→shards 1024 --n_training_shards 1024\2" "+wq" data/create_datasets_from_start.sh
export BERT_PREP_WORKING_DIR=/workspace/data/
time ./data/create_datasets_from_start.sh wiki_only |& tee log
```

After execution is finished, phase 1 pre-tokenized and sharded dataset is located at:

```
/workspace/data/hdf5_lower_case_1_seq_len_128_max_pred_20_masked_lm_prob_0.
15_random_seed_12345_dupe_factor_5/wikicorpus_en/
```

Copy this entire directory to `~/examples_datasets/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128` of the trn1.32xlarge machine.

Phase 2 pre-tokenized dataset is located at:

```
/workspace/data/hdf5_lower_case_1_seq_len_512_max_pred_80_masked_lm_prob_0.
15_random_seed_12345_dupe_factor_5/wikicorpus_en/
```

Copy this entire directory to `~/examples_datasets/bert_pretrain_wikicorpus_tokenized_hdf5_seqlen512` of the trn1.32xlarge machine.

## Known issues and limitations

### BERT-large compilation limitations

Optimal BERT-large phase 1 (sequence length 128) batch size is currently 8 for FP32 and 16 for full BF16 with stochastic rounding. Optimal BERT-large phase 2 (sequence length 512) batch size is currently 1 for FP32 and 2 for full BF16 with stochastic rounding.

### BERT-large pretraining with pretokenized dataset hangs when using `xm.save`

Currently, BERT-large pretraining with pretokenized dataset hangs when `xm.save` is used outside of the main training loop.

```
Loop through HDF5 sharded dataset files:
  Train on one HDF5 sharded dataset file
    Loop through batched samples:
      Training iteration
    Save checkpoint using xm.save
```

The reason is that `xm.save` has a synchronization point. However, the HDF5 shared data files do not have the same number of training samples so the workers cannot all reach `xm.save` in the same iteration.

The workaround is to use `xm._maybe_convert_to_cpu` to ensure tensors are moved to CPU followed by `torch.save` as done in the BERT-large pretraining tutorial:

```
cpu_data = xm._maybe_convert_to_cpu(data)
```

## BERT-large two worker pretraining hangs or run out of host memory during checkpointing on trn1.2xlarge

On trn1.2xlarge, where there's limited host memory and CPU resources, the BERT-large two worker pretraining may hang or run out of host memory during checkpointing. This problem can be worked around by not saving optimizer and LR scheduler states in the checkpoint. This is enabled by `--minimal_ckpt` option of the pretraining script.

## BERT precompilation using `neuron_parallel_compile` hangs when using `torchrun`

We use `neuron_parallel_compile` in front of the short run command to do precompilation. However, the following command hangs when running BERT parallel compilation with `torchrun`:

```
neuron_parallel_compile XLA_DOWNCAST_BF16=1 torchrun --nproc_per_node=32 --nnodes=1 dp_
↳ bert_large_hf_pretrain_hdf5.py --steps_this_run 5

...
Updating train metrics in provide results.json file
Current data: {'num_workers': 32, 'epoch': 0, 'steps': 5, 'microsteps': 320, 'loss': -
↳ 22172234.0, 'train_time_minutes': 0.7424166639645894, 'throughput_average': 1839.
↳ 0391805624324, 'throughput_peak': 1840.0107059878164, 'batch_size': 8, 'max_length': 1
↳ 28}
Updating with data: {'num_workers': 32, 'epoch': 0, 'steps': 5, 'microsteps': 320, 'loss
↳ ': -22172234.0, 'train_time_minutes': 0.7826640844345093, 'throughput_average': 1744.
↳ 4691285659471, 'throughput_peak': 1745.4964663587539, 'batch_size': 8, 'max_length': 1
↳ 28}
Checkpointing...
Checkpointing done...
(hangs)
```

The fix is to add `xm.rendezvous` at the end of training to ensure all workers sync up before exiting the script `dp_bert_large_pretrain_hdf5.py`.

```
def _mp_fn(index, flags):
    torch.set_default_tensor_type('torch.FloatTensor')
    train_bert_hdf5(flags)
    xm.rendezvous("_mp_fn finished")
```

## Troubleshooting

The following are troubleshooting tips related to this tutorial. See *PyTorch Neuron on Trainium Troubleshooting Guide* for additional troubleshooting tips.

### ModuleNotFoundError: No module named 'torch', 'torch\_xla', 'transformers', etc

If you encounter 'ModuleNotFoundError' messages while attempting to run the demo scripts, please ensure that you have activated the appropriate Python *virtualenv* which contains all of the demo dependencies:

```
cd ~
source <python virtual environment path>/bin/activate
```

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Multi-Layer Perceptron Training Tutorial

MNIST is a standard dataset for handwritten digit recognition. A multi-layer perceptron (MLP) model can be trained with MNIST dataset to recognize hand-written digits. This tutorial starts with a 3-layer MLP training example in PyTorch on CPU, then show how to modify it to run on Trainium using PyTorch Neuron. It also shows how to do multiple worker data parallel MLP training.

### Table of Contents

- *Setup environment and download examples*
- *Multi-layer perceptron MNIST model*
- *Single-worker MLP training script in PyTorch on CPU*
- *Single-worker MLP training on Trainium*
- *Multi-worker data-parallel MLP training using torchrun*
- *Single-worker MLP evaluation on Trainium*
- *Known issues and limitations*

**Note:** Logs used in tutorials do not present latest performance numbers

For latest performance numbers visit *Neuron performance*

## Setup environment and download examples

Before running the tutorial please follow the installation instructions at:

*Install PyTorch Neuron on Trn1*

Please set the storage of instance to *512GB* or more if you also want to run through the BERT pretraining and GPT pretraining tutorials.

For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

Install needed dependencies in your environment by running:

```
pip install pillow
```

Torchvision package is needed for MNIST dataset and has already been installed as part of Install PyTorch Neuron on Trn1. Installing torchvision together with torch-neuronx ensures that the compatible version of torchvision is selected. For example, torchvision==0.12 is compatible with torch==1.11 and torchvision==0.13 is compatible with torch==1.12.

To download the MNIST MLP examples, do:

```
git clone https://github.com/aws-neuron/aws-neuron-samples.git
cd aws-neuron-samples/torch-neuronx/training/mnist_mlp
```

## Multi-layer perceptron MNIST model

In `model.py`, we define the multi-layer perceptron (MLP) MNIST model with 3 linear layers and ReLU activations, followed by a log-softmax layer. This model will be used in multiple example scripts.

## Single-worker MLP training script in PyTorch on CPU

We will show how to modify a training script that runs on other platform to run on Trainium.

We begin with a single-worker MLP training script for running on the host CPUs of the Trainium instance. The training script imports the MLP model from `model.py`.

In this training script, we load the MNIST train dataset and, within the `main()` method, set the data loader to read batches of 32 training examples and corresponding labels.

Next we instantiate the MLP model and move it to the device. We use `device = 'cpu'` to illustrate the use of device in PyTorch. On GPU you would use `device = 'cuda'` instead.

We also instantiate the other two components of a neural network trainer: stochastic-gradient-descent (SGD) optimizer and negative-log-likelihood (NLL) loss function (also known as cross-entropy loss).

After the optimizer and loss function, we create a training loop to iterate over the training samples and labels, performing the following steps for each batch in each iteration:

- Zero gradients using:

```
optimizer.zero_grad()
```

- Move training samples and labels to device using the `'tensor.to'` method.

- Perform forward/prediction pass using

```
output = model(train_x)
```

- The prediction results are compared against the corresponding labels using the loss function to compute the loss

```
loss_fn(output, train_label)
```

- The loss is propagated back through the model using chain-rule to compute the weight gradients

```
loss.backward()
```

- The weights are updated with a change that is proportional to the computed weights gradients

```
optimizer.step()
```

At the end of training we compute the throughput, display the final loss and save the checkpoint.

Expected CPU output:

```
-----Training -----
Train throughput (iter/sec): 286.96994718801335
Final loss is 0.1040
-----End Training -----
```

Run the command below to execute this script:

```
python train_cpu.py
```

For a full tutorial on training in PyTorch, please see <https://pytorch.org/tutorials/beginner/introyt/trainingyt.html>.

Thus far we have used PyTorch without Trainium. Next, we will show how to change this script to run on Trainium.

## Single-worker MLP training on Trainium

To run on Trainium, first we modify the CPU training script `train_cpu.py` to run with PyTorch Neuron `torch_xla` as described in *PyTorch Neuron for Trainium Getting Started Guide* by changing the device:

```
import torch_xla.core.xla_model as xm
device = xm.xla_device()
# or
device = 'xla'
```

When the model is moved to the XLA device using `model.to(device)` method, subsequent operations on the model are recorded for later execution. This is XLA's lazy execution which is different from PyTorch's eager execution. Within the training loop, we must mark the graph to be optimized and run on XLA device (NeuronCore) using `xm.mark_step()` (unless `MpDeviceLoader` is used as you will see in the next section). Without this mark, XLA cannot determine where the graph ends. The collected computational graph also gets compiled and executed when you request the value of a tensor such as by calling `loss.item()` or `print(loss)`.

To save a checkpoint, it is recommended to use the `xm.save()` function instead of `torch.save()` to ensure states are moved to CPU. `xm.save()` also prevents the "XRT memory handle not found" warning at the end of evaluation script (if the checkpoint saved using `torch.save()` is used for evaluation).

The resulting script `train.py` can be executed as `python3 train.py`. Again, note that we import the MLP model from `model.py`. When you examine the script, the comments that begin with 'XLA' indicate the changes required to make the script compatible with `torch_xla`.

Run the command below to execute this script:

```
python train.py
```

Expected output on trn1.32xlarge (start from a fresh compilation cache, located at /var/tmp/neuron-compile-cache by default):

```
2022-04-12 16:15:00.000947: INFO ||NCC_WRAPPER||: No candidate found under /var/tmp/
↳neuron-compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_18200615679846498221.
2022-04-12 16:15:00.000949: INFO ||NCC_WRAPPER||: Cache dir for the neff: /var/tmp/
↳neuron-compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_18200615679846498221/
↳MODULE_0_SyncTensorsGraph.318_18200615679846498221_ip-172-31-69-14.ec2.internal-
↳8355221-28940-5dc775cd78aa2/83a0fd4a-b07e-4404-aa55-701ab3b2700c
.....
Compiler status PASS
2022-04-12 16:18:05.000843: INFO ||NCC_WRAPPER||: Exiting with a successfully compiled_
↳graph
2022-04-12 16:18:05.000957: INFO ||NCC_WRAPPER||: No candidate found under /var/tmp/
↳neuron-compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_5000680699473283909.
2022-04-12 16:18:05.000960: INFO ||NCC_WRAPPER||: Cache dir for the neff: /var/tmp/
↳neuron-compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_5000680699473283909/
↳MODULE_1_SyncTensorsGraph.390_5000680699473283909_ip-172-31-69-14.ec2.internal-8355221-
↳28940-5dc7767e5fc69/7d0a2955-11b4-42e6-b536-6f0f02cc68df
.
Compiler status PASS
2022-04-12 16:18:12.000912: INFO ||NCC_WRAPPER||: Exiting with a successfully compiled_
↳graph
-----Training -----
Train throughput (iter/sec): 95.06756661972014
Final loss is 0.1979
-----End Training -----
```

If you re-run the training script a second time, you will see messages indicating that the compiled graphs are cached in the persistent cache from the previous run and that the startup time is quicker:

```
(aws_neuron_venv_pytorch_p36) [ec2-user@ip-172-31-69-14 mnist_mlp]$ python train.py |&
↳tee log_trainium
2022-04-12 16:21:58.000241: INFO ||NCC_WRAPPER||: Using a cached neff at /var/tmp/neuron-
↳compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_18200615679846498221/MODULE_0_
↳SyncTensorsGraph.318_18200615679846498221_ip-172-31-69-14.ec2.internal-8355221-28940-
↳5dc775cd78aa2/83a0fd4a-b07e-4404-aa55-701ab3b2700c/MODULE_0_SyncTensorsGraph.318_
↳18200615679846498221_ip-172-31-69-14.ec2.internal-8355221-28940-5dc775cd78aa2.neff.
↳Exiting with a successfully compiled graph
2022-04-12 16:21:58.000342: INFO ||NCC_WRAPPER||: Using a cached neff at /var/tmp/neuron-
↳compile-cache/USER_neuroncc-1.0.47218.0+162039557/MODULE_5000680699473283909/MODULE_1_
↳SyncTensorsGraph.390_5000680699473283909_ip-172-31-69-14.ec2.internal-8355221-28940-
↳5dc7767e5fc69/7d0a2955-11b4-42e6-b536-6f0f02cc68df/MODULE_1_SyncTensorsGraph.390_
↳5000680699473283909_ip-172-31-69-14.ec2.internal-8355221-28940-5dc7767e5fc69.neff.
↳Exiting with a successfully compiled graph
-----Training -----
Train throughput (iter/sec): 93.16748895384832
Final loss is 0.1979
-----End Training -----
```

Multiple graphs can be created during execution since there are differences between some iterations (first, steady state,

last). After the first iteration, the graph for each iteration should remain the same from iteration to iteration. This allows XLA runtime to execute a previous compiled graph that has been cached in XLA runtime cache.

If the inner training loop has some control-flows, for example for gradient accumulation, the number of compiled graphs may increase due to the generation and consumption of intermediates as well as additional operations when the conditional path is taken.

## Multi-worker data-parallel MLP training using torchrun

Data parallel training allows you to replicate your script across multiple workers, each worker processing a proportional portion of the dataset, in order to train faster.

The PyTorch distributed utility `torchrun` can be used to launch multiple processes in a server node for multi-worker data parallel training.

To run multiple workers in data parallel configuration using `torchrun`, modify the single-worker training script `train.py` as follows (below we use `xm` as alias for `torch_xla.core.xla_model` and `xmp` as alias for `torch_xla.distributed.xla_multiprocessing`):

1. Import XLA backend for `torch.distributed` using `import torch_xla.distributed.xla_backend`.
2. Use `torch.distributed.init_process_group('xla')` to initialize PyTorch XLA runtime and Neuron runtime.
3. Use XLA multiprocessing device loader (`MpDeviceLoader`) from `torch_xla.distributed` to wrap PyTorch data loader.
4. Use `xm.optimizer_step(optimizer)` to perform allreduce and take optimizer step.

XLA `MpDeviceLoader` is optimized for XLA and is recommended for best performance. It also takes care of marking the step for execution (compile and execute the lazily collected operations for an iteration) so no separate `xm.mark_step()` is needed.

The following are general best-practice changes needed to scale up the training:

1. Set the random seed to be the same across workers.
2. Scale up the learning rate by the number of workers. Use `xm.xrt_world_size()` to get the global number of workers.
3. Add distributed sampler to allow different worker to sample different portions of dataset.

Also, the `xm.save()` function used to save checkpoint automatically saves only for the rank-0 worker's parameters.

The resulting script is `train_torchrun.py` (note again that we import the MLP model from `model.py`):

Next we use the `torchrun` utility that is included with torch installation to run multiple processes, each using one Logical NeuronCore. Use the option `nproc_per_node` to indicate the number of processes to launch. For example, to run on two Logical NeuronCores on one Trn1/Trn2 instance only, do:

Run the command below to execute this script:

```
torchrun --nproc_per_node=2 train_torchrun.py
```

**Note:** Currently we only support: - 1 and 2 worker configurations on `trn1.2xlarge` (default Logic NeuronCores size of 1) - 1, 2, 8, and 32-worker configurations on `trn1.32xlarge` (default Logic NeuronCores size of 1) - 1, 4, 16 and 64-worker configurations on `trn2.48xlarge` (default Logic NeuronCores size of 2)

Expected output on `trn1.32xlarge` (second run to avoid compilations):

```

-----Training -----
-----Training -----
... (Info messages truncated)
Train throughput (iter/sec): 163.25353269069706
Train throughput (iter/sec): 163.23261047441036
Final loss is 0.3469
Final loss is 0.1129
-----End Training -----
-----End Training -----

```

In another example, we run on two trn1.32xlarge instances launched with EFA-enabled interfaces, using [EFA-enabled security group](#), and setup using Install PyTorch Neuron on Trn1. NOTE: To run on multiple instances, you will need to use trn1.32xlarge instances and using all 32 NeuronCores on each instance.

On the rank-0 Trn1 host (root), run with `--node_rank=0` using `torchrun` utility, and `--master_addr` set to rank-0 host's IP address:

```

export FI_EFA_USE_DEVICE_RDMA=1
export FI_PROVIDER=efa
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=0 --master_addr=<root IP> --master_
↪port=2020 train_torchrun.py

```

On another Trn1 host, run with `--node_rank=1`, and `--master_addr` also set to rank-0 host's IP address:

```

export FI_EFA_USE_DEVICE_RDMA=1
export FI_PROVIDER=efa
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=1 --master_addr=<root IP> --master_
↪port=2020 train_torchrun.py

```

It is important to launch rank-0 worker with `--node_rank=0` to avoid hang.

For trn2.48xlarge, use `--nproc_per_node=64` for 64 Logical NeuronCores default (each Logical NeuronCores using two physical NeuronCores).

To train on multiple instances, it is recommended to use a `ParallelCluster`. For a `ParallelCluster` example, please see [Train a model on AWS Trn1 ParallelCluster](#).

## Single-worker MLP evaluation on Trainium

After training, the final checkpoint is saved in `checkpoints` directory. You can run the evaluation step by running the `eval.py` script in the same directory as the training script:

Run the command below to execute this script:

```

cd ~/aws-neuron-samples/torch-neuronx/training/mnist_mlp
python eval.py

```

This evaluation phase can be merged with the training script to check accuracy, for example at the end of every epoch. It is kept separate for illustration purpose.

The evaluation script follow similar flow as the training script with the following differences:

- The input data used is the validation subset of the MNIST dataset.
- Only need to loop through the dataset once (no epochs).
- There's only forward pass through the model, and no backward pass or optimizer update.



- Compute the accuracy across validation set instead of loss per batch.

Expected results (after a second execution to eliminate warmup compilation time during first execution):

```
-----Evaluating-----
Test throughput (iter/sec): 47.897945949832845
Accuracy: 0.9273833632469177
-----Done Evaluating-----
```

If you get a lower accuracy than above, please check that the training is done with at least 4 epochs.

You can also use *PyTorch NeuronX Tracing API for Inference* in the evaluation loop. This can be achieved by the following changes to the `eval.py`:

- Use `device = 'cpu'` instead of XLA device.
- Don't use `mark_step()`.
- Trace the model at the first iteration to freeze it and precompile for inference:

```
if idx == 0:
    import torch_neuronx
    model = torch_neuronx.trace(model, test_x)
```

However, note that the inference trace API fixed the input tensor shape, so that every input tensor will need to match the size used during the tracing step. To ensure every batch from `DataLoader` has the same tensor shape, pass `drop_last=True` option when instantiating `DataLoader`.

```
test_loader = DataLoader(test_dataset, batch_size=32, drop_last=True)
```

The script `eval_using_trace.py` can be compared against `eval.py` to show the above modifications. It can be executed using:

Run the command below to execute this script:

```
python eval_using_trace.py
```

Expected results (note the large increase in performance when using trace API for inference):

```
-----Evaluating-----
Test throughput (iter/sec): 409.0836291417652
Accuracy: 0.9288585186004639
-----Done Evaluating-----
```

## Known issues and limitations

MLP model is not optimized for performance. For the single-worker training, the performance can be improved by using `MpDeviceLoader` which exists in the multiprocessing example. For example, by setting `--nproc_per_node=1` in the `torchrun` example, you will see higher MLP performance.

```
(aws_neuron_venv_pytorch_p36) [ec2-user@ip-172-31-69-14 mnist_mlp]$ torchrun --nproc_per_
node=1 train_torchrun.py

-----Training -----
... (Info messages truncated)
Train throughput (iter/sec): 192.43508922834008
```

(continues on next page)

(continued from previous page)

```
Final loss is 0.2720
-----End Training -----
```

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## PyTorch Neuron for Trainium Hugging Face BERT MRPC task finetuning using Hugging Face Trainer API

**Note:** Please use Hugging Face *Optimum-Neuron* <<https://huggingface.co/docs/optimum-neuron/index>> for best coverage and support of Hugging Face models running on Trainium and Inferentia devices.

In this tutorial, we show how to run a Hugging Face script that uses Hugging Face Trainer API to do fine-tuning on Trainium. The example follows the [text-classification example](#) which fine-tunes BERT-base model for sequence classification on the GLUE benchmark.

### Table of Contents

- *Setup and compilation*
- *Single-worker training*
- *Multi-worker data-parallel training*
- *Converting BERT pretrained checkpoint to Hugging Face pretrained model format*
- *Older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0*
- *Known issues and limitations*

**Note:** Logs used in tutorials do not present latest performance numbers

For latest performance numbers visit [Neuron performance](#)

## Setup and compilation

Before running the tutorial please follow the installation instructions at:

*[Install PyTorch Neuron on Trn1](#)*

Please set the storage of instance to *512GB* or more if you also want to run through the BERT pretraining and GPT pretraining tutorials.

For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

First we install a recent version of HF transformers, scikit-learn and evaluate packages in our environment as well as download the source matching the installed version. In this example, we use the text classification example from HF transformers source:

```
export HF_VER=4.52.0
export ACC_VER=1.7.0
pip install -U transformers==$HF_VER accelerate==$ACC_VER datasets evaluate scikit-learn
cd ~/
git clone https://github.com/huggingface/transformers --branch v$HF_VER
```

## Single-worker training

We will run MRPC task fine-tuning following the example in README.md located in the path `~/transformers/examples/pytorch/text-classification`. In this part of the tutorial we will use the Hugging Face model hub's pretrained `bert-large-uncased` model.

**Note:** If you are using older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0, please see section *Older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0* for necessary workarounds.

We use BF16 mixed-precision casting using trainer API `--bf16` option and compiler flag `--model-type=transformer` to enable best performance. We also launch the `run_glue.py` script with `torchrun` using `--nproc_per_node=N` option to specify the number of workers. Here we start of with 1 worker.

**Note:** With transformers version 4.44 and up, please use `torchrun` even for one worker (`--nproc_per_node=1`) to avoid execution hang.

First, paste the following script into your terminal to create a “run.sh” file and change it to executable:

```
tee run.sh > /dev/null <<EOF
#!/usr/bin/env bash
set -eExuo
export TASK_NAME=mrpc
export NEURON_CC_FLAGS="--model-type=transformer"
NEURON_RT_STOCHASTIC_ROUNDING_EN=1 torchrun --nproc_per_node=1 ./run_glue.py \\\
--model_name_or_path bert-large-uncased \\\
--task_name $TASK_NAME \\\
--do_train \\\
--do_eval \\\
--bf16 \\\
--use_cpu True \\\
--max_seq_length 128 \\\
--per_device_train_batch_size 8 \\\
--learning_rate 2e-5 \\\
--num_train_epochs 5 \\\
--save_total_limit 1 \\\
--overwrite_output_dir \\\
--output_dir /tmp/$TASK_NAME/ |& tee log_run
EOF

chmod +x run.sh
```

We optionally precompile the model and training script using `neuron_parallel_compile` to warm up the persistent graph cache (Neuron Cache) such that the actual run has fewer compilations (faster run time):

```
neuron_parallel_compile ./run.sh
```

Please ignore the results from this precompile run as it is only for extracting and compiling the XLA graphs.

---

**Note:** With both train and evaluation options (`--do_train` and `--do_eval`), you will encounter harmless error `ValueError: Target is multiclass but average='binary' when using neuron_parallel_compile`.

---

Precompilation is optional and only needed to be done once unless hyperparameters such as batch size are modified. After the optional precompilation, the actual run will be faster with minimal additional compilations.

```
./run.sh
```

If precompilation was not done, the first execution of `./run.sh` will be slower due to serial compilations. Rerunning the same script a second time would show quicker execution as the compiled graphs will be already cached in persistent cache.

### Multi-worker data-parallel training

The above script would run one worker on one Logical NeuronCore. To run on multiple Logical NeuronCores in data-parallel configuration, launch the `run_glue.py` script with `torchrun` using `--nproc_per_node=N` option to specify the number of workers (N=2 for `trn1.2xlarge`, and N=2, 8, or 32 for `trn1.32xlarge`).

---

**Note:** If you are using older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0, please see section [Older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0](#) for necessary workarounds.

---

The following example runs 2 workers. Paste the following script into your terminal to create a “`run_2w.sh`” file and change it to executable:

```
tee run_2w.sh > /dev/null <<EOF
#!/usr/bin/env bash
set -eExuo
export TASK_NAME=mrpc
export NEURON_CC_FLAGS="--model-type=transformer"
NEURON_RT_STOCHASTIC_ROUNDING_EN=1 torchrun --nproc_per_node=2 ./run_glue.py \
--model_name_or_path bert-large-uncased \
--task_name \${TASK_NAME} \
--do_train \
--do_eval \
--bf16 \
--use_cpu True \
--max_seq_length 128 \
--per_device_train_batch_size 8 \
--learning_rate 2e-5 \
--num_train_epochs 5 \
--save_total_limit 1 \
--overwrite_output_dir \
--output_dir /tmp/\${TASK_NAME}/ |& tee log_run_2w
EOF

chmod +x run_2w.sh
```

Again, we optionally precompile the model and training script using `neuron_parallel_compile` to warm up the persistent graph cache (Neuron Cache), ignoring the results from this precompile run as it is only for extracting and compiling the XLA graphs:

```
neuron_parallel_compile ./run_2w.sh
```

Precompilation is optional and only needed to be done once unless hyperparameters such as batch size are modified. After the optional precompilation, the actual run will be faster with minimal additional compilations.

```
./run_2w.sh
```

During run, you will now notice that the “Total train batch size” is now 16 and the “Total optimization steps” is now half the number for one worker training.

## Converting BERT pretrained checkpoint to Hugging Face pretrained model format

If you have a pretrained checkpoint (i.e., from the BERT phase 2 pretraining tutorial), you can run the script below (saved as “convert.py”) to convert BERT pretrained saved checkpoint to Hugging Face pretrained model format. An example phase 2 pretrained checkpoint can be downloaded from `s3://neuron-s3/training_checkpoints/pytorch/dp_bert_large_hf_pretrain/ckpt_29688.pt`. Note that here we also use the `bert-large-uncased` model configuration to match the BERT-Large model trained following BERT phase 2 pretraining tutorial.

```
tee convert.py > /dev/null <<EOF
import os
import sys
import argparse
import torch
import transformers
from transformers import (
    BertForPreTraining,
)
import torch_xla.core.xla_model as xm
from transformers.utils import check_min_version
from transformers.utils.versions import require_version

if __name__ == '__main__':
    parser = argparse.ArgumentParser()
    parser.add_argument('--model_name', type=str, default='bert-large-uncased', help=
    ↪ "Path to model identifier from huggingface.co/models")
    parser.add_argument('--output_saved_model_path', type=str, default='./hf_saved_model
    ↪ ', help="Directory to save the HF pretrained model format.")
    parser.add_argument('--checkpoint_path', type=str, required=True, help="Path to
    ↪ pretrained checkpoint which needs to be converted to a HF pretrained model format")
    args = parser.parse_args(sys.argv[1:])

    model = BertForPreTraining.from_pretrained(args.model_name)
    check_point = torch.load(args.checkpoint_path, map_location='cpu')
    model.load_state_dict(check_point['model'], strict=False)
    model.save_pretrained(args.output_saved_model_path, save_config=True, save_
    ↪ function=xm.save)
    print("Done converting checkpoint {} to HuggingFace saved model in directory {}".
    ↪ format(args.checkpoint_path, args.output_saved_model_path))
EOF
```

Run the conversion script as:

```
python convert.py --checkpoint_path ckpt_29688.pt
```

After conversion, the new Hugging Face pretrained model is stored in the output directory specified by the `--output_saved_model_path` option which is `hf_saved_model` by default. You will use this directory in the next step.

Paste the following script into your terminal to create a “`run_converted.sh`” file and change it to executable: (note that it uses the converted Hugging Face pretrained model in `hf_saved_model` directory):

```
tee run_converted.sh > /dev/null <<EOF
#!/usr/bin/env bash
set -eExuo
export TASK_NAME=mrpc
export NEURON_CC_FLAGS="--model-type=transformer"
NEURON_RT_STOCHASTIC_ROUNDING_EN=1 torchrun --nproc_per_node=2 ./run_glue.py \
--model_name_or_path hf_saved_model \
--tokenizer_name bert-large-uncased \
--task_name $TASK_NAME \
--do_train \
--do_eval \
--bf16 \
--use_cpu True \
--max_seq_length 128 \
--per_device_train_batch_size 8 \
--learning_rate 2e-5 \
--num_train_epochs 5 \
--save_total_limit 1 \
--overwrite_output_dir \
--output_dir /tmp/$TASK_NAME/ |& tee log_run_converted
EOF

chmod +x run_converted.sh
```

If it is the first time running with `bert-large-uncased` model or if hyperparameters have changed, then the optional one-time precompilation step can save compilation time:

```
neuron_parallel_compile ./run_converted.sh
```

If you have run the single worker training in a previous section, then you can skip the precompilation step and just do:

```
./run_converted.sh
```

### Older versions of transformers <4.27.0 or PyTorch Neuron <1.13.0

If using older versions of transformers package before 4.27.0 or PyTorch Neuron before 1.13.0, please edit the python script `run_glue.py` and add the following lines after the Python imports. They set the compiler flag for transformer model type and enable data parallel training using `torchrun`:

```
# Enable torchrun
import os
import torch
```

(continues on next page)

(continued from previous page)

```

import torch_xla.distributed.xla_backend
from packaging import version
from transformers import __version__, Trainer
if version.parse(__version__) < version.parse("4.26.0") and os.environ.get("WORLD_SIZE"):
    torch.distributed.init_process_group('xla')

# Disable DDP for torchrun
import contextlib
if version.parse(__version__) < version.parse("4.20.0"):
    def _wrap_model(self, model, training=True):
        model.no_sync = lambda: contextlib.nullcontext()
        return model
else:
    def _wrap_model(self, model, training=True, dataloader=None):
        model.no_sync = lambda: contextlib.nullcontext()
        return model
Trainer._wrap_model = _wrap_model

# Workaround for NaNs seen with transformers version >= 4.21.0
# https://github.com/aws-neuron/aws-neuron-sdk/issues/593
import transformers
if os.environ.get("XLA_USE_BF16") or os.environ.get("XLA_DOWNCAST_BF16"):
    transformers.modeling_utils.get_parameter_dtype = lambda x: torch.bfloat16

```

## Known issues and limitations

The following are currently known issues:

- During model evaluation, there can be small compilations for every evaluation step due to a [known transformers issue](#). The work-around is to set training arguments `eval_do_concat_batches=False` and apply the changes in [the PR](#) which will be in a future release of transformers package (version 4.52 or later).
- With transformers==4.44.0, running one worker fine-tuning without torchrun would result in a hang. To workaround and run one worker fine-tuning, use `torchrun --nproc_per_node=1 <script>`.
- With torch-neuronx 2.1, HF Trainer API's use of XLA function `xm.mesh_reduce` causes "EOFError: Ran out of input" or "`_pickle.UnpicklingError: invalid load key, '!'`" errors during Neuron Parallel Compile. This is an issue with the trial execution of empty NEFFs and should not affect the normal execution of the training script.
- Multi-worker training using Trainer API resulted in too many graph compilations for HF transformers>=4.35: This is resolved with HF transformers>=4.37 with the additional workarounds as shown in [`the ticket<https://github.com/aws-neuron/aws-neuron-sdk/issues/813>`](#).
- Long compilation times: this can be alleviated with `neuron_parallel_compile` tool to extract graphs from a short trial run and compile them in parallel ahead of the actual run, as shown above.
- When precompiling using batch size of 16 on trn1.2xlarge, you will see `ERROR ||PARALLEL_COMPILE||: parallel compilation with neuronx-cc exited with error.Received error code: -9.` To workaround this error, please set `NEURON_PARALLEL_COMPILE_MAX_RETRIES=1` in the environment.
- With release 2.6 and transformers==4.25.1, using `neuron_parallel_compile` tool to run `run_glue.py` script with both train and evaluation options (`--do_train` and `--do_eval`), you will encounter harmless error `ValueError: Target is multiclass but average='binary'`

- Reduced accuracy for RoBERTa-Large is seen with Neuron PyTorch 1.12 (release 2.6) in FP32 mode with compiler BF16 autocast. The workaround is to set `NEURON_CC_FLAGS="-auto-cast none"` or set `NEURON_RT_STOCHASTIC_ROUNDING_EN=1`.
- When using DDP in PT 1.13, compilation of one graph will fail with “Killed” error message for `bert-large-uncased`. For `bert-base-cased`, the final MRPC evaluation accuracy is 31% which is lower than expected. These issues are being investigated and will be fixed in an upcoming release. For now, DDP is disabled with the workaround shown above in *Multi-worker Training*.
- When using DDP in PT 1.13 with `neuron_parallel_compile` precompilation, you will hit an error `Rank 1 has 393 params, while rank 0 has inconsistent 0 params..` To workaround this error, add the follow code snippet at the top of `run_glue.py` to skip the problematic shape verification code during precompilation:

```
import os
if os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY", None):
    import torch.distributed as dist
    _verify_param_shape_across_processes = lambda process_group, tensors, logger=None:
    ↪ True
```

- Variable input sizes: When fine-tune models such as `dslim/bert-base-NER` using the *token-classification example*, you may encounter timeouts (lots of “socket.h:524 CCOM WARN Timeout waiting for RX” messages) and execution hang. This occurs because NER dataset has different sample sizes, which causes many recompilations and compiled graph (NEFF) reloads. Furthermore, different data parallel workers can execute different compiled graph. This multiple-program multiple-data behavior is currently unsupported. To workaround this issue, please pad to maximum length using the Trainer API option `--pad_to_max_length`.
- When running HuggingFace GPT fine-tuning with transformers version `>= 4.21.0` and using `XLA_USE_BF16=1` or `XLA_DOWNCAST_BF16=1`, you might see NaNs in the loss immediately at the first step. This issue occurs due to large negative constants used to implement attention masking (<https://github.com/huggingface/transformers/pull/17306>). To workaround this issue, please use transformers version `<= 4.20.0`.
- When using Trainer API option `-bf16`, you will see “RuntimeError: No CUDA GPUs are available”. To workaround this error, please add “`import torch; torch.cuda.is_bf16_supported = lambda: True`” to the Python script (i.e. `run_glue.py`). (Trainer API option `-fp16` is not yet supported).
- When using latest HuggingFace transformers version, you may see “ValueError: Your setup doesn’t support bf16/gpu.” To fix this, please use `--use_cpu True` in your scripts.

The following are resolved issues:

- Using `neuron_parallel_compile` tool to run `run_glue.py` script with both train and evaluation options (`--do_train` and `--do_eval`), you will encounter `INVALID_ARGUMENT` error. To avoid this, only enable train for parallel compile (`--do_train`). This will cause compilations during evaluation step. The `INVALID_ARGUMENT` error is fixed in release 2.6 together with latest transformers package version 4.25.1.
- When running HuggingFace BERT (any size) fine-tuning tutorial or pretraining tutorial with transformers version `>= 4.21.0` and `< 4.25.1` and using `XLA_USE_BF16=1` or `XLA_DOWNCAST_BF16=1`, you will see NaNs in the loss immediately at the first step. More details on the issue can be found at [pytorch/xla#4152](https://pytorch/xla#4152). The workaround is to use transformers version `< 4.21.0` or `>= 4.25.1`, or add `transformers.modeling_utils.get_parameter_dtype = lambda x: torch.bfloat16` to your Python script (i.e. `run_glue.py`).
- Some recompilation is seen at the epoch boundary even after `neuron_parallel_compile` is used. This can be fixed by using the same number of epochs both during precompilation and the actual run.
- When running multi-worker training, you may see the process getting killed at the time of model saving on `trn1.2xlarge`. This happens because the transformers `trainer.save_model` api uses `xm.save` for saving models. This api is known to cause high host memory usage in multi-worker setting see *Saving and Loading XLA Tensors in* . Coupled with a compilation at the same time results in a host OOM. To avoid this issue, we can: Pre-compile all the graphs in multi-worker training. This can be done by running the multi-worker training first with



`neuron_parallel_compile <script>` followed by the actual training. This would avoid the compilation at model save during actual training.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Fine-tune T5 model on Trn1

**Note:** Update 01/03/24: This tutorial is currently broken and the AWS Neuron team is working on the fix.

In this tutorial, we show how to fine-tune a Hugging Face (HF) T5 model using HF trainer API. This example fine-tunes a [T5 model for a text-summarization](#) task on CNN/DailyMail dataset.

### Table of Contents

- [Setup and compilation](#)
- [Single-worker training](#)
- [Multi-worker Training](#)
- [Known issues and limitations](#)

**Note:** Logs used in tutorials do not present latest performance numbers

For latest performance numbers visit [Neuron performance](#)

## Setup and compilation

Before running the tutorial please follow the installation instructions at:

[Install PyTorch Neuron on Trn1](#)

Please set the storage of instance to *512GB* or more if you also want to run through the BERT pretraining and GPT pretraining tutorials.

For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

First we install a recent version of HF transformers, scikit-learn and evaluate packages in our environment as well as download the source matching the installed version. In this example, we chose version 4.26.0 and the text summarization example from HF transformers source:

```
export HF_VER=4.26.0
pip install -U transformers==$HF_VER datasets evaluate scikit-learn rouge_score_
↪ pandas==1.4.0
cd ~/
git clone https://github.com/huggingface/transformers --branch v$HF_VER
cd ~/transformers/examples/pytorch/summarization
```

## Single-worker training

We will run text-summarization fine-tuning task following the example in README.md located in the path `~/transformers/examples/pytorch/summarization`.

We use full BF16 casting using `XLA_USE_BF16=1` to enable best performance. First, paste the following script into your terminal to create a “run.sh” file and change it to executable:

```
tee run.sh > /dev/null <<EOF
#!/bin/bash
set -eExuo
if [ \${NEURON_PARALLEL_COMPILE} == "1" ]
then
    XLA_USE_BF16=1 python3 ./run_summarization.py \
        --model_name_or_path t5-small \
        --dataset_name cnn_dailymail \
        --dataset_config "3.0.0" \
        --do_train \
        --do_eval \
        --source_prefix "summarize: " \
        --max_source_length 512 \
        --per_device_train_batch_size 32 \
        --per_device_eval_batch_size 4 \
        --overwrite_output_dir \
        --pad_to_max_length \
        --max_steps 100 \
        --max_eval_samples 100 \
        --gradient_accumulation_steps=32 \
        --output_dir /tmp/tst-summarization |& tee log_run
else
    XLA_USE_BF16=1 python3 ./run_summarization.py \
        --model_name_or_path t5-small \
        --dataset_name cnn_dailymail \
        --dataset_config "3.0.0" \
        --do_train \
        --do_eval \
        --source_prefix "summarize: " \
        --max_source_length 512 \
        --per_device_train_batch_size 32 \
        --per_device_eval_batch_size 4 \
        --overwrite_output_dir \
        --pad_to_max_length \
        --gradient_accumulation_steps=32 \
        --output_dir /tmp/tst-summarization |& tee log_run
fi
EOF

chmod +x run.sh
```

We optionally precompile the model and training script using `neuron_parallel_compile` to warm up the persistent graph cache (Neuron Cache) such that the actual run has fewer compilations (faster run time):

```
neuron_parallel_compile ./run.sh
```

Note: For these auto-regressive models, do not run the `predict_with_generate` method when doing the precompile

step. This is because the `neuron_parallel_compile` utility will run the training script in graph extraction mode and no actual execution of the graph will be done. Hence, the outputs at each step are invalid. Since the auto-regressive generation at each step is dependent on output of previous step, the generate step would fail since the outputs from previous steps are invalid.

Precompilation is optional and only needs to be done once unless hyperparameters such as batch size are modified. After the optional precompilation, the actual run will be faster with minimal additional compilations.

```
./run.sh
```

If precompilation was not done, the first execution of `./run.sh` will be slower due to serial compilations. Rerunning the same script a second time would show quicker execution as the compiled graphs will be already cached in persistent cache.

Running the above script will run the T5-small fine-tuning on a single process.

**Note:** As you may have noticed, we are not running the `predict_with_generate` as part of training. This is because, `predict_with_generate` requires auto-regressive sampling where the inputs to the decoder are created by appending outputs of previous steps. This causes the inputs to the decoder to change shape and thereby resulting in a new graph. In other words, the current `generate` api provided by HF transformers leads to repeated compilations. We are working on building a Neuron friendly version of `generate` api and it will be made available as part of future release. This will enable us to run `predict_with_generate` as part of training script.

As a workaround, we can run the `predict_with_generate` on CPU after the model is trained. Once training is completed, a trained checkpoint would be saved. We can load the trained model and run the `predict_with_generate` to compute the final accuracy.

To do so, in `run_summarization.py`, add the following before `transformers` get imported. This can be done by adding the below lines before all the imports:

```
import libneuronxla
# Disable configuring xla env
def _configure_env():
    pass
libneuronxla.configure_environment = _configure_env
```

You can now run the following and it should run the predict method on CPU device.

```
NEURON_NUM_DEVICES=0 python3 ./run_summarization.py \
  --model_name_or_path <CHECKPOINT_DIR> \
  --dataset_name cnn_dailymail \
  --dataset_config "3.0.0" \
  --do_predict \
  --predict_with_generate \
  --source_prefix "summarize: " \
  --per_device_eval_batch_size 4 \
  --max_source_length 512 \
  --pad_to_max_length \
  --no_cuda \
  --output_dir /tmp/tst-summarization |& tee log_run
```

Note: To run on CPU, we need to make sure that `NEURON_NUM_DEVICES` is set to 0. This will make sure no `xla_devices` are created and the trainer would use the default device (CPU).

## Multi-worker Training

The above script will run one worker on one NeuronCore. To run on multiple cores, first add these lines to top of `run_summarization.py` to disable Distributed Data Parallel (DDP) when using `torchrun` (see Known issues and limitations section below):

```
# Disable DDP for torchrun
from transformers import __version__, Trainer
Trainer._wrap_model = lambda self, model, training=True, dataloader=None: model
```

Then launch the `run_summarization.py` script with `torchrun` using `--nproc_per_node=N` option to specify the number of workers (N=2 for `trn1.2xlarge`, and N=2, 8, or 32 for `trn1.32xlarge`). The following example runs 2 workers. Paste the following script into your terminal to create a “`run_2w.sh`” file and change it to executable:

```
tee run_2w.sh > /dev/null <<EOF
#!/bin/bash
set -eExuo
if [ \${NEURON_PARALLEL_COMPILE} == "1" ]
then
    XLA_USE_BF16=1 torchrun --nproc_per_node=2 ./run_summarization.py \
    --model_name_or_path t5-small \
    --dataset_name cnn_dailymail \
    --dataset_config "3.0.0" \
    --do_train \
    --do_eval \
    --source_prefix "summarize: " \
    --max_source_length 512 \
    --per_device_train_batch_size 32 \
    --per_device_eval_batch_size 4 \
    --overwrite_output_dir \
    --pad_to_max_length \
    --max_steps 100 \
    --max_eval_samples 100 \
    --gradient_accumulation_steps=32 \
    --output_dir /tmp/tst-summarization |& tee log_run
else
    XLA_USE_BF16=1 torchrun --nproc_per_node=2 ./run_summarization.py \
    --model_name_or_path t5-small \
    --dataset_name cnn_dailymail \
    --dataset_config "3.0.0" \
    --do_train \
    --do_eval \
    --source_prefix "summarize: " \
    --max_source_length 512 \
    --per_device_train_batch_size 32 \
    --per_device_eval_batch_size 4 \
    --overwrite_output_dir \
    --pad_to_max_length \
    --gradient_accumulation_steps=32 \
    --output_dir /tmp/tst-summarization |& tee log_run
fi
EOF

chmod +x run_2w.sh
```

Again, we optionally precompile the model and training script using `neuron_parallel_compile` to warm up the persistent graph cache (Neuron Cache), ignoring the results from this precompile run as it is only for extracting and compiling the XLA graphs:

```
neuron_parallel_compile ./run_2w.sh
```

Precompilation is optional and only needs to be done once unless hyperparameters such as batch size are modified. After the optional precompilation, the actual run will be faster with minimal additional compilations.

```
./run_2w.sh
```

During run, you will notice that the “Total train batch size” is now 8 and the “Total optimization steps” is now half the number for one worker training. Also, if you open `neuron-top` in a separate terminal, you should see 2 cores been utilized.

To train T5-large model, you can set the `model_name_or_path` argument to `t5-large`. Please note, currently running `t5-large` on `trn1-2xl` machine can result in `HOST OOM` during compilation. Hence, it is recommended that you run a `t5-large` model training on a `trn1-32xl` machine.

On a `trn1-32xl` machine, you can create a `run_32w.sh` on the terminal using the following commands:

```
tee run_32w.sh > /dev/null <<EOF
#!/bin/bash
set -eExuo
if [ \${NEURON_PARALLEL_COMPILE} == "1" ]
then
    XLA_USE_BF16=1 torchrun --nproc_per_node=32 ./run_summarization.py \
    --model_name_or_path t5-large \
    --dataset_name cnn_dailymail \
    --dataset_config "3.0.0" \
    --do_train \
    --do_eval \
    --source_prefix "summarize: " \
    --max_source_length 512 \
    --per_device_train_batch_size 4 \
    --per_device_eval_batch_size 4 \
    --overwrite_output_dir \
    --pad_to_max_length \
    --max_steps 100 \
    --max_eval_samples 100 \
    --gradient_accumulation_steps=11 \
    --output_dir /tmp/tst-summarization |& tee log_run
else
    XLA_USE_BF16=1 torchrun --nproc_per_node=32 ./run_summarization.py \
    --model_name_or_path t5-large \
    --dataset_name cnn_dailymail \
    --dataset_config "3.0.0" \
    --do_train \
    --do_eval \
    --source_prefix "summarize: " \
    --max_source_length 512 \
    --per_device_train_batch_size 4 \
    --per_device_eval_batch_size 4 \
    --overwrite_output_dir \
    --pad_to_max_length \
```

(continues on next page)

(continued from previous page)

```
--gradient_accumulation_steps=11 \  
--output_dir /tmp/tst-summarization |& tee log_run  
fi  
EOF  
  
chmod +x run_32w.sh
```

You can now follow the same steps as listed above. This script would run a t5-large model by launching a training script using 32 data-parallel workers.

## Known issues and limitations

The following are currently known issues:

- Long compilation times: this can be alleviated with `neuron_parallel_compile` tool to extract graphs from a short trial run and compile them in parallel ahead of the actual run, as shown above.
- T5-Large compilation causing processes to get killed on trn1-2xl: It is recommended to `t5-large` model training on a trn1-32xl machine, as it avoids CPU OOM and also provides faster training by making use of 32 data-parallel workers.

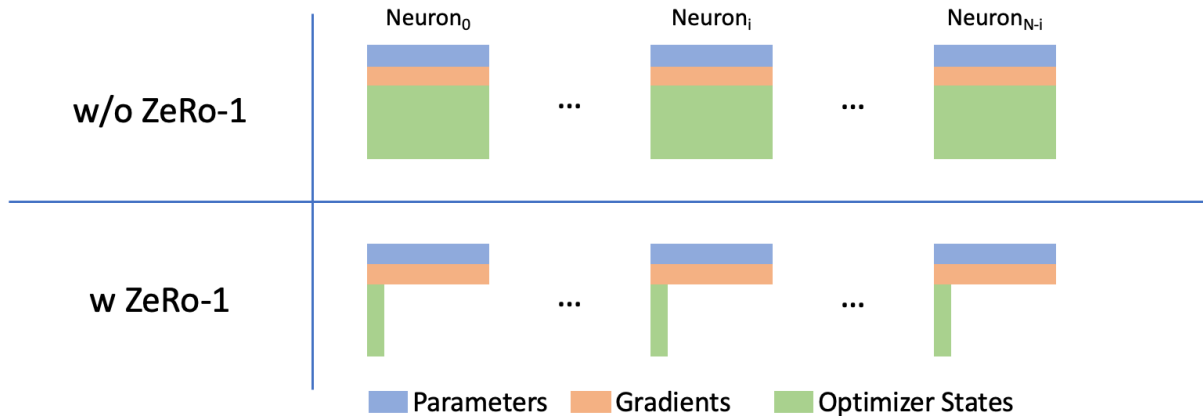
*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## ZeRO-1 Tutorial

### What is ZeRO-1?

ZeRO-1 (Zero Redundancy Optimizer Stage 1, <https://arxiv.org/abs/1910.02054>) is an optimization technique for large-scale deep learning models. It is a memory efficient variation of data parallelism. ZeRO leverages the aggregate computation and memory resources of data parallelism to reduce the memory and compute requirements of each accelerator used for model training. ZeRO reduces the memory consumption of each accelerator by partitioning the various model training states (weights, gradients, and optimizer states) across the available devices in the distributed training hardware. ZeRO is being implemented as incremental stages of optimizations. In stage 1, the optimizer states (e.g., for Adam optimizer, 32-bit weights, and the first, and second moment estimates) are partitioned across the processes, so that each process updates only its partition.



Assume we use mixed precision training with Adam optimizer. Comparing the total memory usage with and with out ZeRO-1.  $\psi$  denotes model size (number of parameters), and  $N_d$  denotes DP degree. Then, without ZeRO-1 memory consumption is  $(2 + 2 + 3 * 4) * \psi = 16\psi$ , with ZeRO-1 is  $2\psi + 2\psi + 12\psi/N_d$ .

We implemented an XLA-friendly version of ZeRO-1 and it has been merged in open-source PyTorch/XLA project. Users can use it to enable ZeRO-1 algorithm by simply wrapping the origin optimizer as shown below.

```
# Before:
optimizer = torch.optim.Adam(model.parameters(), lr=0.0001)

# After
optimizer = ZeroRedundancyOptimizer(model.parameters(), torch.optim.Adam, lr=0.0001)
```

Then just call `optimizer.step()` directly, the wrapped optimizer will handle the distributed operations automatically.

The above code snippet illustrates the basic usage. Generally, users can use ZeRO-1 optimizer like a normal optimizer. In addition, `ZeroRedundancyOptimizer` also provides other features: enable gradient clipping or use other data type for wrapped optimizer. Note that though the most of optimizers can be used with ZeRO-1, optimizers that compute norm for parameters (e.g. LAMB) might lead to accuracy disparities compared to using original local optimizer when using ZeRO-1, because these optimizers cannot get full parameters but shards.

## Usage

To enable ZeRO-1 optimizer, just import it and replace origin optimizer with ZeRO-1 wrapped version

```
from torch_xla.distributed.zero_redundancy_optimizer import ZeroRedundancyOptimizer
...
...

device = xm.xla_device()
model = model.to(device)

optimizer = ZeroRedundancyOptimizer(model.parameters(), AdamW, lr=0.001)
```

Then in training loop, just call `optimizer.step()`, note that we should not use `xm.reduce_gradients()` or `xm.optimizer_step()` as gradient reduction will be handle by ZeRO-1.

```
...
loss.backward()
xm.mark_step()
optimizer.step()
xm.mark_step()
```

ZeRO-1 optimizer also provides some additional features, user can pass these arguments to the wrapper constructor:

- Change `optimizer_dtype` to choose data dtype used by optimizer, default is `torch.float32`. For example, when parameter data type is `bfloat16`, set `optimizer_dtype` to be `float32` to enable ‘master weight’.
- Change `grad_clipping` to enable grad clipping, default is `True`.
- Change `max_norm` to determine the maximum norm value used by grad clipping, default is `1.0`.
- Change `use_grad_acc_hook` to enable using buffers to store gradients, it will use the same data type as `optimizer_dtype` to accumulate gradients. (Added in neuron 2.19.0 release).
- Change `higher_cc_precision` to force reduce-scatter operator to use the same data type as `optimizer_dtype`, default is `False`. When `use_grad_acc_hook` is `True`, it has no effects. (Added in neuron 2.19.0 release).

Note: ZeRO-1 optimizer now forces to use the same data type as parameters for all-gather operator. (Changed in neuron 2.19.0 release)

## GPT2-XL Pretraining Tutorial

### Table of contents

- [Setup](#)
- [Dataset](#)
- [Training](#)
- [Known Issues, Work-arounds and Limitations](#)

## Setup

We use single Trn1.32xlarge instance. Follow [Install PyTorch Neuron on Trn1](#) to setup the environment first. For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

**requirements.txt:** We pin the following Hugging Face Library versions necessary for the tutorial

```
transformers==4.27.3
accelerate==0.17
tensorboard==2.12.2
```

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

```
git clone https://github.com/aws-neuron/aws-neuron-samples.git
cd aws-neuron-samples/torch-neuronx/training/zero1_gpt2
python3 -m pip install -r requirements.txt
```



The specific files you need for this tutorial:

- `config_1p5B_gpt2.json`: The model configuration used in the tutorial for GPT 2.7B Neo
- `neuron_utils.py`: includes utility functions and the logging tools
- `run_clm_no_trainer.py`: the main training script that runs the actual training
- `run_clm.sh`: the shell script to launch the training job

## Dataset

For the dataset, we use the wikitext dataset, specifically `wikitext-103-raw-v1`, provided by the HuggingFace <https://huggingface.co/datasets/wikitext>. The data will be preprocessed the first time running through the training script and then preprocessed data will be cached in the HuggingFace cache directory for any future training runs.

If the main process downloads the dataset, tokenizes the data and groups them together successfully, the expected output would be as below at the beginning of the training.

```
***** Running training *****
Num examples = 114248
Num Epochs = 29
Instantaneous batch size per device = 1
Total train batch size (w. parallel, distributed & accumulation) = 32
Gradient Accumulation steps = 1
Total optimization steps = 100000
```

## Training

The GPT2 python fine-tuning script is adapted from the example `run_clm_no_trainer.py` in <https://github.com/huggingface/transformers/tree/main/examples/pytorch/language-modeling>. It incorporates the Accelerate <https://github.com/huggingface/accelerate>. Given its beta stage, some modifications are needed, along with the bridge code to XLA. Particularly, some workarounds to support Accelerate for the training script are listed in “Known Issues Workarounds and Limitations” below.

In this example, we use GPT2-xl as example, and show the training steps with mixed precision (bfloat16 and float32)

- single node training:

```
# Run precompile and training
neuron_parallel_compile bash run_clm.sh MIXED wikitext-103-raw-v1
bash run_clm.sh MIXED wikitext-103-raw-v1
```

- multi-node training, run:

```
sbatch run_clm_compile.slurm
```

then

```
sbatch run_clm.slurm
```

## Known Issues, Work-arounds and Limitations

1. Error message: `ValueError: invalid literal for int() with base 10: ''`. Simply re-run the script can solve this issue. This issue is already solved in the newer versions of transformers, see <https://github.com/huggingface/transformers/pull/22427>.
2. Accelerator API workarounds:
  - Error message: “Gradient accumulation is not supported on TPU. Please set `gradient_accumulation_steps` to 1 and don’t pass in a `GradientAccumulationPlugin` object.” More context here: <https://github.com/huggingface/accelerate/pull/479>. The training still works by commenting out the assertion and avoid using the accumulation wrapper with `accelerator.accumulate(model)`
  - `Accelerator.prepare` call: We have noticed that using the optimizer returned by this API are not directly reusable. It is due to gaps in configuring accelerate API for XLA devices.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Analyze for Training Tutorial

This tutorial explains how to analyze a model for training support using via `torch-neuronx`.

---

**Note:** For analyzing models for inference support via `torch-neuronx`, please refer to `torch_neuronx.analyze()`

---

## Setup

For this tutorial we’ll be using two scripts: `supported.py` and `unsupported.py`. Create these files by copy pasting the below code to their respective files.

`supported.py`

```
import torch
import torch_xla.core.xla_model as xm

class NN(torch.nn.Module):
    def __init__(self):
        super().__init__()

        self.layer1 = torch.nn.Linear(4,4)
        self.nl1 = torch.nn.ReLU()
        self.layer2 = torch.nn.Linear(4,2)
        self.nl2 = torch.nn.Tanh()

    def forward(self, x):
        x = self.nl1(self.layer1(x))
        return self.nl2(self.layer2(x))

def main():
    device = xm.xla_device()
```

(continues on next page)

(continued from previous page)

```

model = NN().to(device)
optimizer = torch.optim.SGD(model.parameters(), lr=0.01)
loss_fn = torch.nn.MSELoss()

inp = torch.rand(4)
target = torch.tensor([1,0])

model.train()
for epoch in range(2):
    optimizer.zero_grad()
    inp = inp.to(device)
    target = target.to(device)
    output = model(inp)
    loss = loss_fn(output,target)
    loss.backward()
    optimizer.step()
    xm.mark_step()

if __name__ == '__main__':
    main()

```

unsupported.py

```

import torch
import torch_xla.core.xla_model as xm

class UnsupportedModel(torch.nn.Module):
    def __init__(self):
        super().__init__()

    def forward(self, x):
        y = torch.fft.fft(x)
        x = x + 10
        return x * y

def main():
    device = xm.xla_device()

    model = UnsupportedModel().to(device)

    inp = torch.rand(4)

    model.train()
    for epoch in range(1):
        inp = inp.to(device)
        output = model(inp)

        xm.mark_step()

if __name__ == '__main__':

```

(continues on next page)

(continued from previous page)

```
main()
```

## Running analyze via neuron\_parallel\_compile

To analyze a model, we supply the training script to the analyze command, which is shipped with neuron\_parallel\_compile. The command is:

```
neuron_parallel_compile --command analyze python supported.py
```

This will generate a lot of output showing a lot of compilation statuses. Here's a snippet of the output when running the above command.

```
.2023-05-25 00:43:43.000394: 776642 INFO ||ANALYZE||: Compiling /tmp/model_analysis_
→graphs/compare_7841189860629745939_23.hlo.pb using following command: neuronx-cc_
→compile --target=trn1 --framework XLA /tmp/model_analysis_graphs/compare_
→7841189860629745939_23.hlo.pb --verbose=35 --query-compute-placement
2023-05-25 00:43:43.000418: 776642 INFO ||ANALYZE||: Compiling /tmp/model_analysis_
→graphs/multiply_15640857564712679356_53.hlo.pb using following command: neuronx-cc_
→compile --target=trn1 --framework XLA /tmp/model_analysis_graphs/multiply_
→15640857564712679356_53.hlo.pb --verbose=35 --query-compute-placement
.
Compiler status PASS
2023-05-25 00:43:43.000549: 776642 INFO ||ANALYZE||: Compiling /tmp/model_analysis_
→graphs/subtract_1927104012014828209_49.hlo.pb using following command: neuronx-cc_
→compile --target=trn1 --framework XLA /tmp/model_analysis_graphs/subtract_
→1927104012014828209_49.hlo.pb --verbose=35 --query-compute-placement
...
Compiler status PASS
```

The analysis report will be generated as a JSON file. The location of the report is shown as the last log entry:

```
2023-05-25 00:43:49.000252: 776642 INFO ||ANALYZE||: Removing existing report /home/
→ubuntu/analyze_for_training/model_analysis_result/result.json
2023-05-25 00:43:49.000252: 776642 INFO ||ANALYZE||: Model analysis completed. Report -
→ /home/ubuntu/analyze_for_training/model_analysis_result/result.json
```

**Note:** Note that if a report is already present in the specified path, analyze will remove/overwrite it.

The report generated running the above command looks like:

```
{
  "torch_neuronx_version": "1.13.0.1.6.1",
  "neuronx_cc_version": "2.5.0.28+1be23f232",
  "support_percentage": "100.00%",
  "supported_operators": {
    "aten": {
      "aten::permute": 8,
      "aten::add": 8,
      "aten::mul": 8,
      "aten::expand": 18,
```

(continues on next page)

(continued from previous page)

```

        "aten::mm": 10,
        "aten::mse_loss_backward": 12,
        "aten::relu": 3,
        "aten::threshold_backward": 4,
        "aten::squeeze": 4,
        "aten::view": 4,
        "aten::pow": 2,
        "aten::mse_loss": 2,
        "aten::tanh": 2
    }
},
"unsupported_operators": {
    "aten": []
}
}

```

---

**Note:** Note that the torch\_neuronx and neuronx\_cc versions may be different from this example

---

## Understanding analyze report for Unsupported Models

### Default Verbosity

Let's run analyze for unsupported.py

```
neuron_parallel_compile --command analyze python unsupported.py
```

Here is the report generated by the above command:

```

{
  "torch_neuronx_version": "1.13.0.1.6.1",
  "neuronx_cc_version": "2.5.0.28+1be23f232",
  "support_percentage": "60.00%",
  "supported_operators": {
    "aten": {
      "aten::add": 2,
      "aten::mul": 1
    }
  },
  "unsupported_operators": {
    "aten": [
      {
        "kind": "aten::mul",
        "failureAt": "neuronx-cc",
        "call": "test2_unsup.py 24"
      }
    ]
  }
}

```

In the list of unsupported operators we are provided the specific aten op that failed, and where that operator is in the training script.

One thing to notice is that the `support_percentage` doesn't exactly add up. This is because the `support_percentage` is calculated based on the supported number of XLA/HLO instructions (explained more in the next section). To see the specific XLA/HLO op lowerings, use the flag `--analyze-verbosity 1`, as the default is 2.

The last thing is that a specific aten operator can be supported and unsupported simultaneously. In our example, this can be seen with `aten::mul`. This is due to the configuration of the aten op. The below section will describe what went wrong with the `aten::mul` op.

## Lower Level Verbosity

Let's run again with lower verbosity level:

```
neuron_parallel_compile --command analyze --analyze-verbosity 1 python unsupported.py
```

The report looks like:

```
{
  "torch_neuronx_version": "1.13.0.1.6.1",
  "neuronx_cc_version": "2.5.0.28+1be23f232",
  "support_percentage": "60.00%",
  "supported_operators": {
    "aten": {
      "aten::mul": 1,
      "aten::add": 2
    },
    "xla": [
      "f32[] multiply(f32[], f32[])",
      "f32[4]{0} broadcast(f32[]), dimensions={}",
      "f32[4]{0} add(f32[4]{0}, f32[4]{0})"
    ]
  },
  "unsupported_operators": {
    "aten": [
      {
        "kind": "aten::mul",
        "failureAt": "neuronx-cc",
        "call": "test2_unsup.py 24"
      }
    ],
    "xla": [
      {
        "hlo_instruction": "c64[4]{0} convert(f32[4]{0})",
        "aten_op": "aten::mul"
      },
      {
        "hlo_instruction": "c64[4]{0} multiply(c64[4]{0}, c64[4]{0})",
        "aten_op": "aten::mul"
      }
    ]
  }
}
```

(continues on next page)

(continued from previous page)

```
}
}
```

This report provides both the aten operator and the failed XLA/HLO instructions. There will be more HLO instructions than aten ops since an aten op generally lowers to multiple HLO instructions. As a result, the `support_percentage` field doesn't exactly line up with the aten operator count, but does line up the XLA/HLO instruction count. This level of verbosity is intended for use when you have the ability to modify the model's HLO lowering, or generally have insight into the HLO lowering.

As mentioned before, the `aten::mul` op appears to be both supported and unsupported. This is because the compiler does not support a specific configuration of `aten::mul`, which can be seen more clearly with the HLO lowering. In the above example, the `aten::mul` operator is unsupported since at least one parameter provided was a complex type (C64), which is unsupported by `neuronx-cc`.

This concludes the tutorial. The API for `analyze` can be found within [neuron\\_parallel\\_compile](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Neuron Custom C++ Operators in MLP Training

In this tutorial we'll demonstrate how to prepare a PyTorch model that contains a custom operator (ie. `CppExtension`) for Neuron compilation to run on Trainium EC2 instances. To learn more about Neuron CustomOps see [Neuron Custom C++ Operators \[Beta\]](#). For a deeper dive on MNIST or Multi-Layer Perceptron models, see the [Multi-Layer Perceptron Training Tutorial](#). This tutorial assumes the reader is familiar with [PyTorch Custom Extensions](#).

### Table of Contents

- [Setup Environment and Download Examples](#)
- [Basic PyTorch Custom Relu Operator](#)
- [Multi-layer perceptron MNIST model](#)
- [Training the MLP model on CPU](#)
- [Neuron Relu CustomOp](#)
- [Training the MLP model on Trainium](#)

## Setup Environment and Download Examples

Before running the tutorial please follow the installation instructions at:

- `pytorch-neuronx-install` on Trn1

**Note:** The name of `aws-neuronx-gpsimd-customop` has been changed to `aws-neuronx-gpsimd-customop-lib` as of the neuron 2.10 release.

**Note:** Custom C++ Operators are supported as of Neuron SDK Version 2.7 as a beta feature. As such this feature is not installed by default, additional tooling and library packages (RPM and DEB) are required.

For AL2023 only, the following packages need be installed as dependencies:

```
sudo yum install libns1
sudo yum install libxcrypt-compat
```

On AL2 and AL2023, they can be installed with the following commands:

```
sudo yum remove python3-devel -y
sudo yum remove aws-neuronx-gpsimd-tools-0.* -y
sudo yum remove aws-neuronx-gpsimd-customop-lib-0.* -y

sudo yum install python3-devel -y
sudo yum install aws-neuronx-gpsimd-tools-0.* -y
sudo yum install aws-neuronx-gpsimd-customop-lib-0.* -y
```

On Ubuntu, they can be installed with the following commands:

```
sudo apt-get remove python3-dev -y
sudo apt-get remove aws-neuronx-gpsimd-tools=0.* -y
sudo apt-get remove aws-neuronx-gpsimd-customop-lib=0.* -y

sudo apt-get install python3-dev -y
sudo apt-get install aws-neuronx-gpsimd-tools=0.* -y
sudo apt-get install aws-neuronx-gpsimd-customop-lib=0.* -y
```

---

For all the commands below, make sure you are in the virtual environment that you have created above before you run the commands:

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

Install dependencies for PyTorch Custom Extensions in your environment by running:

```
pip install regex
pip install ninja
```

The `ninja` package is only needed for the reference CPU example. It is not needed by Neuron to run on Trainium instances.

To download the source code for this tutorial, do:

```
git clone https://github.com/aws-neuron/aws-neuron-samples.git
cd aws-neuron-samples/torch-neuronx/training/customop_mlp
```

In the `customop_mlp` directory there are two subdirectories. The `pytorch` directory contains an example model and training script using a custom operator that runs using the `cpu` device with standard PyTorch APIs and libraries (ie. not specific to AWS/Neuron). The `neuron` directory contains a version of the same model and training script with the custom operator ported to Neuron to run on `trn1` using the XLA device.



## Basic PyTorch Custom Relu Operator

For the next few sections we'll review the example model in the `pytorch` directory. This is a condensed and simplified explanation of PyTorch C++ Extensions, for more details see the [PyTorch documentation](#). In `my_ops.py` we implement a custom relu activation op as a torch autograd function so that we can use it in a training loop:

```
import torch

torch.ops.load_library('librelu.so')

class Relu(torch.autograd.Function):
    @staticmethod
    def forward(ctx, input):
        ctx.save_for_backward(input)
        return torch.ops.my_ops.relu_forward(input)

    @staticmethod
    def backward(ctx, grad):
        input, = ctx.saved_tensors
        return torch.ops.my_ops.relu_backward(grad, input), None
```

Notice that here we first load `librelu.so` using the `load_library` API. And then call the `relu_forward` and `relu_backward` functions from our library within the relevant static methods.

We implemented these two library functions in the `relu.cpp` file:

```
torch::Tensor relu_forward(const torch::Tensor& t_in) {
    ...
    t_out_acc[i][j] = t_in_acc[i][j] > 0.0 ? t_in_acc[i][j] : 0.0;
    ...
}

torch::Tensor relu_backward(const torch::Tensor& t_grad, const torch::Tensor& t_in) {
    ...
    t_out_acc[i][j] = t_in_acc[i][j] > 0.0 ? t_grad_acc[i][j] : 0.0;
    ...
}

TORCH_LIBRARY(my_ops, m) {
    m.def("relu_forward", &relu_forward);
    m.def("relu_backward", &relu_backward);
}
```

And then built them into a library using the PyTorch Cpp Extension APIs in the `build.py` script:

```
torch.utils.cpp_extension.load(
    name='librelu',
    sources=['relu.cpp'],
    is_python_module=False,
    build_directory=os.getcwd()
)
```

Run `python build.py` to produce the `librelu.so` library.

## Multi-layer perceptron MNIST model

In `model.py`, we define the multi-layer perceptron (MLP) MNIST model with 3 linear layers and a custom ReLU activation, followed by a log-softmax layer. Highlighted below are the relevant custom changes in the `model.py` file:

```
import torch
import torch.nn as nn
from torch.nn import functional as F
import my_ops

# Declare 3-layer MLP for MNIST dataset
class MLP(nn.Module):
    def __init__(self, input_size = 28 * 28, output_size = 10, layers = [120, 84]):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = my_ops.Relu.apply(f1)
        f2 = self.fc2(r1)
        r2 = my_ops.Relu.apply(f2)
        f3 = self.fc3(r2)
        return torch.log_softmax(f3, dim=1)
```

## Training the MLP model on CPU

In the `train_cpu.py` script we load the MNIST train dataset, instantiate the MLP model, and use `device='cpu'` to execute on the host CPU. Expected CPU output:

```
-----Training -----
Train throughput>(*iter/sec*): *286*.96994718801335
Final loss is *0*.1040
-----End Training -----
```

## Neuron Relu CustomOp

Now switch over into the `neuron` directory. To migrate our PyTorch customOp to Neuron, we have to make a few small changes. First, we create a new `shape.cpp` file to implement our shape function as required by XLA (see [Neuron Custom C++ Operators Developer Guide \[Beta\]](#) for details). We also replace the `TORCH_LIBRARY` API with `NEURON_LIBRARY`.

```
torch::Tensor relu_fwd_shape(torch::Tensor t_in) {
    torch::Tensor t_out = torch::zeros(t_in.sizes(), torch::kFloat);
    return t_out;
}

torch::Tensor relu_bwd_shape(torch::Tensor t_grad, torch::Tensor t_in) {
    torch::Tensor t_out = torch::zeros(t_in.sizes(), torch::kFloat);
```

(continues on next page)

(continued from previous page)

```

    return t_out;
}

NEURON_LIBRARY(my_ops, m) {
    m.def("relu_forward", &relu_fwd_shape, "relu_forward");
    m.def("relu_backward", &relu_bwd_shape, "relu_backward");
}

```

And then we build it using the `torch_neuronx` package in `build.py`:

```

from torch_neuronx.xla_impl import custom_op

custom_op.load(
    name='relu',
    compute_srcs=['relu.cpp'],
    shape_srcs=['shape.cpp'],
    build_directory=os.getcwd()
)

```

Notice that here we specify both the `relu.cpp` and `shape.cpp` files separately. This is because the shape functions will be compiled with an x86 compiler and run on the host during the XLA compilation, and the compute functions will be compiled for the NeuronCore accelerator and executed during the training loop. Running `build.py` produces the same `librelu.so` as in the CPU example, but compiles the source code to execute on the NeuronCore.

In our `my_ops.py` file we just use the `torch_neuronx` API to load our new library and execute our `customOp` exactly the same way we did before:

```

import torch
import torch_neuronx
from torch_neuronx.xla_impl import custom_op

custom_op.load_library('librelu.so')

class Relu(torch.autograd.Function):
    @staticmethod
    def forward(ctx, input):
        ctx.save_for_backward(input)
        return torch.ops.my_ops.relu_forward(input)

    @staticmethod
    def backward(ctx, grad):
        input, = ctx.saved_tensors
        return torch.ops.my_ops.relu_backward(grad, input), None

```

## Training the MLP model on Trainium

In the `train.py` script we modify the CPU training script `train_cpu.py` to run with PyTorch Neuron `torch_xla`. Expected output on a `trn1` instance:

```
-----Training -----
2023-02-02 22 (tel:2023020222):46:58.000299: INFO ||NCC_WRAPPER||: Using a cached neff_
→at /var/tmp/neuron-compile-cache/USER_neuroncc-2.0.0.8683a0+c94c3936c/MODULE_
→4447837791278761679/MODULE_0_SyncTensorsGraph.329_4447837791278761679_ip-172-31-38-167.
→us-west-2.compute.internal-49ad7ade-14011-5f3bf523d8788/1650ba41-bcfd-4d15-9038-
→16d391c4a57c/MODULE_0_SyncTensorsGraph.329_4447837791278761679_ip-172-31-38-167.us-
→west-2.compute.internal-49ad7ade-14011-5f3bf523d8788.neff. Exiting with a successfully_
→compiled graph
2023-02-02 22 (tel:2023020222):46:58.000433: INFO ||NCC_WRAPPER||: Using a cached neff_
→at /var/tmp/neuron-compile-cache/USER_neuroncc-2.0.0.8683a0+c94c3936c/MODULE_
→16964505026440903899/MODULE_1_SyncTensorsGraph.401_16964505026440903899_ip-172-31-38-
→167.us-west-2.compute.internal-4d0cabba-14011-5f3bf529794a3/23d74230-59dd-4347-b247-
→fa98aed416bd/MODULE_1_SyncTensorsGraph.401_16964505026440903899_ip-172-31-38-167.us-
→west-2.compute.internal-4d0cabba-14011-5f3bf529794a3.neff. Exiting with a successfully_
→compiled graph
Train throughput (iter/sec): 117.47151142662648
Final loss is 0.1970
-----End Training -----
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Neuron Custom C++ Operators Performance Optimization

In this tutorial, we will build on the small MLP model shown in [Neuron Custom C++ Operators in MLP Training](#) and demonstrate methods to optimize the performance of a custom C++ operator. We will be taking advantage of the TCM accessor as well as the usage of multiple GPSIMD cores to enhance performance.

This tutorial assumes the reader has read and set up an environment described in [Neuron Custom C++ Operators in MLP Training](#).

### Table of Contents

- [Download Examples](#)
- [Model Configuration Adjustment](#)
- [Performance with Element-wise Accessor](#)
- [Performance with TCM Accessor](#)
- [Extending the example to utilize multiple GPSIMD cores](#)

## Download Examples

To download the source code for this tutorial, do:

```
git clone https://github.com/aws-neuron/aws-neuron-samples.git
cd aws-neuron-samples/torch-neuronx/inference/customop_mlp
```

**Note:** We will be using an inference example in this tutorial in order to adhere to certain Custom C++ operator restrictions when using multiple GPUSIMD cores (see [Custom Operators API Reference Guide \[Beta\]](#) for details on current restrictions).

**Note:** Custom C++ Operators are supported as of Neuron SDK Version 2.7 as a beta feature. As such this feature is not installed by default, additional tooling and library packages (RPM and DEB) are required.

For AL2023 only, the following packages need be installed as dependencies:

```
sudo yum install libnsl
sudo yum install libxcrypt-compat
```

On AL2 and AL2023, they can be installed with the following commands:

```
sudo yum remove python3-devel -y
sudo yum remove aws-neuronx-gpsimd-tools=0.* -y
sudo yum remove aws-neuronx-gpsimd-customop-lib=0.* -y

sudo yum install python3-devel -y
sudo yum install aws-neuronx-gpsimd-tools=0.* -y
sudo yum install aws-neuronx-gpsimd-customop-lib=0.* -y
```

On Ubuntu, they can be installed with the following commands:

```
sudo apt-get remove python3-dev -y
sudo apt-get remove aws-neuronx-gpsimd-tools=0.* -y
sudo apt-get remove aws-neuronx-gpsimd-customop-lib=0.* -y

sudo apt-get install python3-dev -y
sudo apt-get install aws-neuronx-gpsimd-tools=0.* -y
sudo apt-get install aws-neuronx-gpsimd-customop-lib=0.* -y
```

Activate the virtual environment created in [Neuron Custom C++ Operators in MLP Training](#),

```
source ~/aws_neuron_venv_pytorch/bin/activate
```

As a reminder, ninja should be already installed in the virtual environment. If not, install it for PyTorch Custom Extensions in your environment by running:

```
pip install regex
pip install ninja
```

## Model Configuration Adjustment

For this tutorial, we will enlarge the size of the hidden layer from [120, 84] to [4096, 2048] in `model.py`.

```
import torch
import torch.nn as nn
from torch.nn import functional as F
import my_ops

# Declare 3-layer MLP for MNIST dataset
class MLP(nn.Module):
    def __init__(self, input_size = 28 * 28, output_size = 10, layers = [4096, 2048]):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = my_ops.Relu.apply(f1)
        f2 = self.fc2(r1)
        r2 = my_ops.Relu.apply(f2)
        f3 = self.fc3(r2)
        return torch.log_softmax(f3, dim=1)
```

## Performance with Element-wise Accessor

The neuron directory contains the same code shown in *Neuron Custom C++ Operators in MLP Training*, where the `relu_forward` is implemented with element-wise accessor. Go to neuron directory, run `build.py` then inference. `py`, the expected output on a `trn1` instance is,

```
Inf throughput (iter/sec): 8.098649744235592
-----End Inference -----
```

## Performance with TCM Accessor

Now we switch to neuron-tcm folder. As mentioned in *Custom Operators API Reference Guide [Beta]*, TCM accessors provide faster read and write performance. We implement the `relu_forward` using TCM accessor in `relu.cpp`:

```
torch::Tensor relu_forward(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = torch::zeros(t_in.sizes(), torch::kFloat);

    static constexpr size_t buffer_size = 1024;
    float *tcm_buffer = (float*)torch::neuron::tcm_malloc(sizeof(float) * buffer_size);

    if (tcm_buffer != nullptr) {
        auto t_in_tcm_acc = t_in.tcm_accessor();
        auto t_out_tcm_acc = t_out.tcm_accessor();
```

(continues on next page)

(continued from previous page)

```

    for (size_t i = 0; i < num_elem; i += buffer_size) {
        size_t remaining_elem = num_elem - i;
        size_t copy_size = (remaining_elem > buffer_size) ? buffer_size : remaining_elem;

        t_in_tcm_acc.tensor_to_tcm<float>(tcm_buffer, i, copy_size);
        for (size_t j = 0; j < copy_size; j++) {
            tcm_buffer[j] = tcm_buffer[j] > 0.0 ? tcm_buffer[j] : 0.0;
        }
        t_out_tcm_acc.tcm_to_tensor<float>(tcm_buffer, i, copy_size);
    }
    torch::neuron::tcm_free(tcm_buffer);
    return t_out;
}

```

Run build.py then inference.py, the expected output on a trn1 instance is:

```

Inf throughput (iter/sec): 220.73800131604054
-----End Inference -----

```

## Extending the example to utilize multiple GPSIMD cores

Now we switch to the neuron-multicore folder. We first enable the usage of multiple GPSIMD cores by multicore=True in the build.py.

```

custom_op.load(
    name='relu',
    compute_srcs=['relu.cpp'],
    shape_srcs=['shape.cpp'],
    build_directory=os.getcwd(),
    multicore=True,
    verbose=True
)

```

After passing the flag, the kernel function relu\_forward defined in relu.cpp will execute on all GPSIMD cores. Thus we need to use cpu\_id to partition the workload among all cores.

```

torch::Tensor relu_forward(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = get_dst_tensor();

    uint32_t cpu_id = get_cpu_id();
    uint32_t cpu_count = get_cpu_count();
    uint32_t partition = num_elem / cpu_count;
    if (cpu_id == cpu_count - 1) {
        partition = num_elem - partition * (cpu_count - 1);
    }

    static constexpr size_t buffer_size = 1024;
    float *tcm_buffer = (float*)torch::neuron::tcm_malloc(sizeof(float) * buffer_size);
}

```

(continues on next page)

(continued from previous page)

```

if (tcm_buffer != nullptr) {
    auto t_in_tcm_acc = t_in.tcm_accessor();
    auto t_out_tcm_acc = t_out.tcm_accessor();

    for (size_t i = 0; i < partition; i += buffer_size) {
        size_t remaining_elem = partition - i;
        size_t copy_size = (remaining_elem > buffer_size) ? buffer_size : remaining_elem;

        t_in_tcm_acc.tensor_to_tcm<float>(tcm_buffer, partition *cpu_id + i, copy_size);
        for (size_t j = 0; j < copy_size; j++) {
            tcm_buffer[j] = tcm_buffer[j] > 0.0 ? tcm_buffer[j] : 0.0;
        }
        t_out_tcm_acc.tcm_to_tensor<float>(tcm_buffer, partition *cpu_id + i, copy_size);
    }
    torch::neuron::tcm_free(tcm_buffer);
    return t_out;
}

```

There are two things noteworthy in the code:

1. We use `cpu_id` and `cpu_count` to distribute the workload among all cores. Particularly, each cores performs relu on a partition of the tensor, the offset is computed based on `cpu_id`.
2. The output of the operator is directly written to the tensor from `get_dst_tensor()`. The `return t_out;` statement is ignored during execution.

Run `build.py` then `inference.py`, the expected output on a `trn1` instance is:

```

Inf throughput (iter/sec): 269.936119707143
-----End Inference -----

```

Details of the API used in the sample here can be found in [Custom Operators API Reference Guide \[Beta\]](#).

This document is relevant for: `Inf2`, `Trn1`, `Trn2`

- [Hugging Face BERT Pretraining Tutorial \(Data-Parallel\)](#)
- [Multi-Layer Perceptron Training Tutorial](#)
- [PyTorch Neuron for Trainium Hugging Face BERT MRPC task finetuning using Hugging Face Trainer API](#)
- [Fine-tune T5 model on Trn1](#)
- [ZeRO-1 Tutorial](#)
- [Analyze for Training Tutorial](#)
- [Neuron Custom C++ Operators in MLP Training](#)
- [Neuron Custom C++ Operators Performance Optimization](#)

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
- [running-jupyter-notebook-as-script](#)

This document is relevant for: `Trn1`, `Trn2`



*This document is relevant for: Inf2, Trn1, Trn2*

### Additional Examples (torch-neuronx)

- [AWS Neuron Reference for Nemo Megatron GitHub Repository](#)
- [AWS Neuron Samples for EKS](#)
- [AWS Neuron Samples for AWS ParallelCluster](#)
- [AWS Neuron Samples GitHub Repository](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### API Reference Guide for Training (torch-neuronx)

*This document is relevant for: Trn1, Trn2*

### PyTorch NeuronX neuron\_parallel\_compile CLI

PyTorch NeuronX performs just-in-time compilation of graphs during execution. At every step, a graph is traced. If the traced graph varies from the previous executions, it is compiled by the neuron compiler. For large models, the compilation time for each graph can be high. Moreover, because of JIT, we would compile all these graphs sequentially, hence incurring huge compilation penalty.

To reduce this compilation time during execution, the `neuron_parallel_compile` utility is provided as part of PyTorch Neuron installation. The `neuron_parallel_compile` will extract graphs from a trial run of your script, perform parallel pre-compilation of the graphs, and populate the *Neuron Persistent Cache* on disk or in AWS S3 bucket with compiled graphs. Your trial run should be limited to a few steps (eg. 10-15), enough for the utility to extract the different graphs needed for full execution. To run the utility:

```
neuron_parallel_compile <run commands>
```

Where `<run commands>` are the commands to run a short run (i.e. 10 steps) to trace training loops for pre-compilation. The example for the run command is `torchrun --nproc_per_node=2 <train script>`, where train script accepts `--steps_this_run` option to limit number of run steps:

```
neuron_parallel_compile torchrun --nproc_per_node=2 <train script> --steps_this_run=10
```

You may notice that the output from the model is invalid when you use `neuron_parallel_compile`. This is because when you initiate your training run command with `neuron_parallel_compile`, the utility will run your command with environment variables that puts your training script into graph extraction mode. In this mode, no real execution is performed and the outputs are invalid. You will also see outputs similar to the following about the compile cache path and the extracted graphs:

```
INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/neuron-compile-cache
INFO ||NEURON_CC_WRAPPER||: Extracting graphs (/var/tmp/neuron-compile-cache/neuronxcc-2.
→0.0.22266a0+a69f71e55/MODULE_9219523464496887986+abb26765/model.hlo.pb) for ahead-of-
→time parallel compilation. No compilation was done.
```

After the trial execution ends and the graphs are extracted, `neuron_parallel_compile` would launch multiple compilation processes in parallel to compile all these graphs. Compiled graphs (NEFFs) are inserted into the Neuron Persistent Cache. You will also see outputs similar to the following about the compile cache path, the list of graphs (HLOs) to be compiled, and the running statistics of compiled graphs (count of remaining graphs, locked graphs, failed graphs, done compiled graphs).

```
INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/neuron-compile-cache
INFO ||NEURON_CACHE||: Current remaining items are 5, locked are 0, failed are 0, done_
→are 0, total is 5
INFO ||NEURON_PARALLEL_COMPILE||: master grab hlos to compile: ['/var/tmp/neuron-compile-
→cache/neuronxcc-2.0.0.22266a0+a69f71e55/MODULE_8068656800389078395+abb26765/model.hlo.
→pb', '/var/tmp/neuron-compile-cache/neuronxcc-2.0.0.22266a0+a69f71e55/MODULE_
→17109392703413819652+abb26765/model.hlo.pb', '/var/tmp/neuron-compile-cache/neuronxcc-
→2.0.0.22266a0+a69f71e55/MODULE_9219523464496887986+abb26765/model.hlo.pb', '/var/tmp/
→neuron-compile-cache/neuronxcc-2.0.0.22266a0+a69f71e55/MODULE_
→16969875447143373016+abb26765/model.hlo.pb', '/var/tmp/neuron-compile-cache/neuronxcc-
→2.0.0.22266a0+a69f71e55/MODULE_3000743782456078279+abb26765/model.hlo.pb']
...
INFO ||NEURON_CACHE||: Current remaining items are 0, locked are 0, failed are 0, done_
→are 5, total is 5
```

After all compilations are completed, a compilation summary is shown:

```
INFO: 2023-08-24 20:21:11.000895: 161136 INFO ||NEURON_PARALLEL_COMPILE||: {
INFO:   "compilation_summary": {
INFO:     "true": 2
INFO:   },
INFO:   "compilation_report": {
INFO:     "/var/tmp/neuron-compile-cache/neuronxcc-2.0.0.22266a0+a69f71e55/MODULE_
→1970132581169579119+abb26765/model.hlo.pb": {
INFO:       "status": true,
INFO:       "retry": 0
INFO:     },
INFO:     "/var/tmp/neuron-compile-cache/neuronxcc-2.0.0.22266a0+a69f71e55/MODULE_
→16141953836240613513+abb26765/model.hlo.pb": {
INFO:       "status": true,
INFO:       "retry": 0
INFO:     }
INFO:   }
INFO: }
INFO: 2023-08-24 20:21:11.000895: 161136 INFO ||NEURON_PARALLEL_COMPILE||: Total_
→graphs: 2
INFO: 2023-08-24 20:21:11.000895: 161136 INFO ||NEURON_PARALLEL_COMPILE||: Total_
→successful compilations: 2
INFO: 2023-08-24 20:21:11.000895: 161136 INFO ||NEURON_PARALLEL_COMPILE||: Total_
→failed compilations: 0
```

Now if you run your script (without `neuron_parallel_compile`), it will be faster since the compiled graphs are already cached.

```
torchrun --nproc_per_node=2 <train script>
```

Note: Except for the option to limit number of run steps (such as `--steps_this_run`), the other options of `<run commands>` must match between the pre-compilation and actual run. If this is not the case, you may see additional compilations during training run because of new graphs getting generated, resulting in cache miss.

There may be additional compilations due to unreachable execution paths (in case the execution path is not reached in the first few steps of graph extraction), or changes in parameters such as number of data parallel workers.

Each precompilation command or actual script execution command above can be prefixed with `NEURON_COMPILE_CACHE_URL=<cache URL>` or `NEURON_CC_FLAGS="--cache_dir=<cache URL>"` to specify a

different cache location than the default (with `--cache_dir` taking precedence over `NEURON_COMPILE_CACHE_URL` if both are specified). Alternatively, the cache URL can also be specify in Python code using:

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + "--cache_dir=
↪<cache URL>"
```

You need to specify the same cache URL for both the precompilation command (using `neuron_parallel_compile`) and the actual script execution command if you want the previously compiled and cached graphs to be used for actual script execution.

The environment variables below are available to help modify `neuron_parallel_compile` behavior:

**NEURON\_PARALLEL\_COMPILE\_MAX\_RETRIES :**

- Set the maximum number of retries when using *Neuron Persistent Cache* or *neuron\_parallel\_compile*. If set to N, the tool will try compilation N more time(s) if the first graph compilation failed. Example: Set `NEURON_PARALLEL_COMPILE_MAX_RETRIES=1` when precompiling on `trn1.2xlarge` where there's limited host memory and CPU resources. Default is 0.

**NEURON\_IGNORE\_TRAINING\_SCRIPT\_ERROR\_AND\_COMPILE :**

- When using *Neuron Persistent Cache* or *neuron\_parallel\_compile*, if you want to ignore the error in training script and compile the accumulated HLO graphs, you can do so by setting this environment variable. Example: If `NEURON_IGNORE_TRAINING_SCRIPT_ERROR_AND_COMPILE=1` is set when using `neuron_parallel_compile`, a crash in the training script would be ignored and the graphs collected up to the crash would be compiled.

**NEURON\_COMPILE\_CACHE\_URL:**

- Set the *Neuron Persistent Cache* URL or *neuron\_parallel\_compile*. If starts with `s3://`, it will use AWS S3 as cache backend. Otherwise it will use local disk cache. Default is `/var/tmp/neuron-compile-cache`. If this is specified together with `cache_dir=<cache_url>` option via `NEURON_CC_FLAGS`, the `--cache_dir` option takes precedence.

## Debugging with Neuron Persistent Cache

A graph compilation can fail because of a compilation error or an environment issue (for example, compilation is interrupted by ctrl-C). The graph would be marked as failed and subsequent rerun would encounter message like below:

```
INFO ||NCC_WRAPPER||: Got a cached failed neff at /var/tmp/neuron-compile-cache/
↪neuronxcc-2.8.0.25+a3ad0f342/MODULE_12486829708343293975+d41d8cd9/model.neff. Will_
↪skip compilation, please set --retry_failed_compilation for recompilation.
```

To retry compilation, add `--retry_failed_compilation` in `NEURON_CC_FLAGS` environment variable. This will retry the compilation even if the graph was previously marked as failed compilation.

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --retry_failed_
↪compilation'
```

See *Neuron Persistent Cache* for more information.

## Separate collection and compilation commands

For cases like finetuning, there could be multiple independent training tasks running on different nodes and sharing many compilation graphs in common. `neuron_parallel_compile` provides commands to separate the graph collection and compilation phases, so users can collect all graphs across different training sessions in advance to avoid duplicate compilations.

To only collect the graphs from trial executions of training scripts into Neuron Persistent Cache:

```
neuron_parallel_compile --command collect <run_script>
```

To compile the graph previously collected using `collect` command and store compiled result (NEFFs) back into Neuron Persistent Cache (make sure to use the same `neuronx-cc` compiler version as during the graph collection step):

```
``neuron_parallel_compile --command compile <run_script>``
```

Note: if `--command` is not specified, `neuron_parallel_compile` will do both collection and compilation phases by default.

## Cache maintenance commands

The following commands are available to help maintain the cache.

**Warning:** Make sure no running process is using the cache when you use `clean` or `clear-locks` command because it can cause cache errors.

To clean cached files:

```
# WARNING: Make sure no running process is using the cache
neuron_parallel_compile --command clean
```

To clear file locks left behind when a `neuron_parallel_compile` execution was interrupted:

```
# WARNING: Make sure no running process is using the cache
neuron_parallel_compile --command clear-locks
```

Each command above can be prefixed with `NEURON_COMPILE_CACHE_URL=<cache URL>` or `NEURON_CC_FLAGS="--cache_dir=<cache URL>"` to specify a different cache location than the default.

---

**Note:** Currently there's no automatic maintenance of cache size either on disk or in S3. Please delete files (i.e. older compiler versions) as necessary to keep cache size within your limit.

---

## Analyze operations support

The analyze command checks the support of operations within the training script by checking each operator against neuronx-cc. It is only supported for PyTorch models. The output of the tool will be available as result.json within the output location.

```
neuron_parallel_compile --command analyze python3 training_script.py
```

Optional Arguments:

- analyze-output ANALYZE\_OUTPUT\_LOCATION Only supported for --command analyze. Path to location where output will be persisted. Default: cwd/model\_analysis\_result
- analyze-verbosity {1,2} Only supported for --command analyze. Level of information to be included within the output. 1: add XLA operator information into the results. 2: add aten metadata into results. Default: 2

The tutorial for analyze can be found [here](#)

*This document is relevant for:* Trn1, Trn2

*This document is relevant for:* Trn1, Trn2

## PyTorch NeuronX Environment Variables

Environment variables allow modifications to PyTorch NeuronX behavior without requiring code change to user script. It is recommended to set them in code or just before invoking the python process, such as `NEURON_FRAMEWORK_DEBUG=1 python3 <script>` to avoid inadvertently changing behavior for other scripts. Environment variables specific to PyTorch Neuron are (beta ones are noted):

### NEURON\_CC\_FLAGS

- Compiler options. Full compiler options are described in the [mixed-precision-casting-options](#). Additional options for the Neuron Persistent Cache can be found in the [Neuron Persistent Cache](#) guide.

### NEURON\_FRAMEWORK\_DEBUG

- Enable dumping of XLA graphs in both HLO format (intermediate representation) and text form for debugging.

### NEURON\_EXTRACT\_GRAPHS\_ONLY

- Dump the XLA graphs in HLO format (intermediate representation) and execute empty stubs with zero outputs in order to allow multiple XLA graphs to be traced through a trial execution. Used automatically for ahead-of-time graph extraction for parallel compilation in [neuron\\_parallel\\_compile](#) tool. This environment variable can be checked in the training script to prevent checking of bad outputs during trial run.

### NEURON\_NUM\_RECENT\_MODELS\_TO\_KEEP

- Keep only N number of graphs loaded in Neuron runtime for each process, where N is the value this environment variable is set to. Default is to keep all graphs loaded by a process.

### NEURON\_COMPILE\_CACHE\_URL

- Set the [Neuron Persistent Cache](#) URL or [neuron\\_parallel\\_compile](#). If starts with `s3://`, it will use AWS S3 as cache backend. Otherwise it will use local disk cache. Default is `/var/tmp/neuron-compile-cache`. If this is specified together with `cache_dir=<cache_url>` option via `NEURON_CC_FLAGS`, the `--cache_dir` option takes precedence.

### NEURON\_PARALLEL\_COMPILE\_MAX\_RETRIES

- Set the maximum number of retries when using *Neuron Persistent Cache* or *neuron\_parallel\_compile*. If set to N, the tool will try compilation N more time(s) if the first graph compilation failed. Example: Set NEURON\_PARALLEL\_COMPILE\_MAX\_RETRIES=1 when precompiling on trn1.2xlarge where there's limited host memory and CPU resources. Default is 0.

#### NEURON\_IGNORE\_TRAINING\_SCRIPT\_ERROR\_AND\_COMPILE

- When using *Neuron Persistent Cache* or *neuron\_parallel\_compile*, if you want to ignore the error in training script and compile the accumulated HLO graphs, you can do so by setting this environment variable. Example: If NEURON\_IGNORE\_TRAINING\_SCRIPT\_ERROR\_AND\_COMPILE=1 is set when using *neuron\_parallel\_compile*, a crash in the training script would be ignored and the graphs collected up to the crash would be compiled.

#### NEURON\_PARALLEL\_COMPILE\_DUMP\_RESULTS

- When set to 1, *neuron\_parallel\_compile* would report compilation time results in the final JSON output.

#### NEURON\_FUSE\_SOFTMAX

- Enable custom lowering for Softmax operation to enable compiler optimizations.

#### NEURON\_CUSTOM\_SILU

- Enable custom lowering for SILU operation to enable compiler optimizations.

#### NEURON\_TRANSFER\_WITH\_STATIC\_RING\_OPS

- The list of torch.nn.Modules that will have all parameter input buffers marked as static to enable runtime optimizations. The default is “Embedding,LayerNorm,Linear,Conv2d,BatchNorm2d” for torch-neuronx 1.13/2.1, and “Embedding” for torch-neuronx 2.1 in SDK release 2.20, and empty for torch-neuronx 2.1+ in SDK release 2.21.

#### NEURONCORE\_NUM\_DEVICES [Use only with xmp.spawn]

- Number of NeuronCores for setting up distributed data parallel training when using torch\_xla.distributed.xla\_multiprocessing.spawn (xmp.spawn) utility only. See [MNIST MLP training with xmp.spawn](#) for example. NOTE: Do not use this environment variable when using torchrun, which has --nproc\_per\_node option instead for this purpose. torchrun is recommended for consistent experience on one instance as well as across multiple instances.

#### NEURON\_DUMP\_HLO\_SNAPSHOT [Beta] [Torch-NeuronX 1.13 only]

- Dump the inputs, outputs, and graph in HLO format of a graph execution in a snapshot file. This variable can be set to 1, ON\_NRT\_ERROR, ON\_NRT\_ERROR\_CPU, ON\_NRT\_ERROR\_HYBRID to dump snapshots at every iteration using CPU memory, or dump only on errors automatically using device, host, and both device and host memory respectively.

#### NEURON\_NC0\_ONLY\_SNAPSHOT [Beta] [Torch-NeuronX 1.13 only]

- Dump only the snapshot associated with Neuron Core 0 when NEURON\_NC0\_ONLY\_SNAPSHOT=1 and the NEURON\_DUMP\_HLO\_SNAPSHOT flag is set.

#### NEURON\_TRANSFER\_ALL\_PARAMETERS\_WITH\_STATIC\_RING [Beta]

- When set to 1, mark all parameter transfers as static to enable runtime optimizations for torch.nn modules that are wrapped as done in Megatron-LM. This setting is not needed if torch.nn modules are not wrapped.

#### BUCKET\_CAP\_MB [PyTorch XLA <=2.1]

- If there are many small gradient tensors, such as in BERT training, small allreduce sizes can limit performance. To improve performance, you can try increasing the bucket size using BUCKET\_CAP\_MB environment variable, which is set to 50MB by default. For example, BERT pretraining on multiple instances can see improved performance with BUCKET\_CAP\_MB=512. NOTE: While this is supported in PyTorch Neuron 2.5, it is recommended for users to switch to ALLREDUCE\_GRADIENTS\_BUCKET\_SIZE\_MB.

**ALLREDUCE\_GRADIENTS\_BUCKET\_SIZE\_MB [PyTorch XLA 2.5+]**

- If there are many small gradient tensors, such as in BERT training, small allreduce sizes can limit performance. To improve performance, you can try increasing the bucket size using `ALLREDUCE_GRADIENTS_BUCKET_SIZE_MB` environment variable, which is set to 50MB by default. For example, BERT pretraining on multiple instances can see improved performance with `ALLREDUCE_GRADIENTS_BUCKET_SIZE_MB=512`.

**XLA\_FLAGS [PyTorch XLA] [Torch-NeuronX 2.1+]**

- When set to `--xla_dump_hlo_snapshots --xla_dump_to=<dir>`, this environmental variable enables dumping snapshots in `<dir>` directory. See [Snapshotting With Torch-Neuronx 2.1](#) section for more information.

**XLA\_USE\_DUMMY\_STORE [PyTorch XLA]**

- When set to 1 along with `TORCH_DIST_INIT_BARRIER=0`, PJRT process group initialization will use Dummy-Store instead of TCPStore. This reduces the number of open file descriptors and enables scaling training up to a large number of nodes.

**XLA\_USE\_BF16 [PyTorch XLA <=2.1]**

- When `XLA_USE_BF16=1`, PyTorch Neuron will automatically map both `torch.float` and `torch.double` tensors to `bfloat16` tensors and turn on Stochastic Rounding mode. This can both reduce memory footprint and improve performance. Example: to enable `bfloat16` autocasting and stochastic rounding, set `XLA_USE_BF16=1` only, as stochastic rounding mode is on by default when `XLA_USE_BF16=1`. If you would like to preserve some tensors in `float32`, see `XLA_DOWNCAST_BF16` below. NOTE: This is deprecated in PyTorch Neuron 2.5. See `migration_from_xla_downcast_bf16`.

**XLA\_DOWNCAST\_BF16 [PyTorch XLA <=2.1]**

- When `XLA_DOWNCAST_BF16=1`, PyTorch Neuron will automatically map `torch.float` tensors to `bfloat16` tensors, `torch.double` tensors to `float32` tensors and turn on Stochastic Rounding mode. This can both reduce memory footprint and improve performance, while preserving some tensors in `float32`. Example: to enable `float` to `bfloat16` and `double` to `float` autocasting and stochastic rounding, set `XLA_DOWNCAST_BF16=1` only, as stochastic rounding mode is on by default when `XLA_DOWNCAST_BF16=1`. If you want to cast both `torch.float` and `torch.double` to `bfloat16`, please see `XLA_USE_BF16` above. NOTE: This is deprecated in PyTorch Neuron 2.5. See `migration_from_xla_downcast_bf16`.

**XLA\_DISABLE\_FUNCTIONALIZATION [PyTorch XLA 2.1+]**

- When `XLA_DISABLE_FUNCTIONALIZATION=0`, PyTorch XLA will enable the functionalization feature which makes graphs more compilable by removing mutations from functions. In PyTorch XLA 2.1 functionalization causes 15% performance degradations for BERT due to missing aliasing for gradient accumulation <https://github.com/pytorch/xla/issues/7174> so it is off by default (`XLA_DISABLE_FUNCTIONALIZATION=1`). Enabling functionalization can improve convergence for LLaMA 70B with ZeRO1 (when used with release 2.19 compiler).

**XLA\_ENABLE\_PARAM\_ALIASING [PyTorch XLA]**

- When `XLA_ENABLE_PARAM_ALIASING=0`, PyTorch Neuron will disable parameter aliasing in HLO graphs. This can be useful for debug. However, it would lead to increased device memory usage due to extra allocation of buffers (so higher chance of out-of-device memory errors) and decreased performance. When not set, parameter aliasing is enabled by default.

**NEURON\_RT\_STOCHASTIC\_ROUNDING\_EN [Neuron Runtime]**

- When `NEURON_RT_STOCHASTIC_ROUNDING_EN=1`, PyTorch Neuron will use stochastic rounding instead of round-nearest-even for all internal rounding operations when casting from FP32 to a reduced precision data type (FP16, BF16, FP8, TF32). This feature has been shown to improve training convergence for reduced precision training jobs, such as when `bfloat16` autocasting is enabled. This is set to 1 by default by PyTorch Neuron



when `XLA_USE_BF16=1` or `XLA_DOWNCAST_BF16=1`. To switch to round-nearest-even mode, please set `NEURON_RT_STOCHASTIC_ROUNDING_EN=0`.

#### `NEURON_RT_STOCHASTIC_ROUNDING_SEED` [Neuron Runtime]

- Sets the seed for the random number generator used in stochastic rounding (see previous section). If this environment variable is not set, the seed is set to 0 by default. Please set `NEURON_RT_STOCHASTIC_ROUNDING_SEED` to a fixed value to ensure reproducibility between runs.

#### `NEURON_RT_VISIBLE_CORES` [Neuron Runtime]

Integer range of specific NeuronCores needed by the process (for example, 0-3 specifies NeuronCores 0, 1, 2, and 3). You this environment variable when using `torchrun` to limit the launched processes to specific consecutive NeuronCores. To ensure best performance, the multi-core jobs requiring N NeuronCores for collective communication must be placed at the NeuronCore ID that starts at a multiple of N, where N is the world size limited to 1, 2, 8, 32. For example, a process using 2 NeuronCores can be mapped to 2 free NeuronCores starting at NeuronCore id 0, 2, 4, 6, etc, and a process using 8 NeuronCores can be mapped to 8 free NeuronCores starting at NeuronCore id 0, 8, 16, 24.

Additional Neuron runtime environment variables are described in [runtime configuration documentation](#).

Additional XLA runtime environment variables are described in [PyTorch-XLA troubleshooting guide](#).

*This document is relevant for:* Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Neuron Persistent Cache

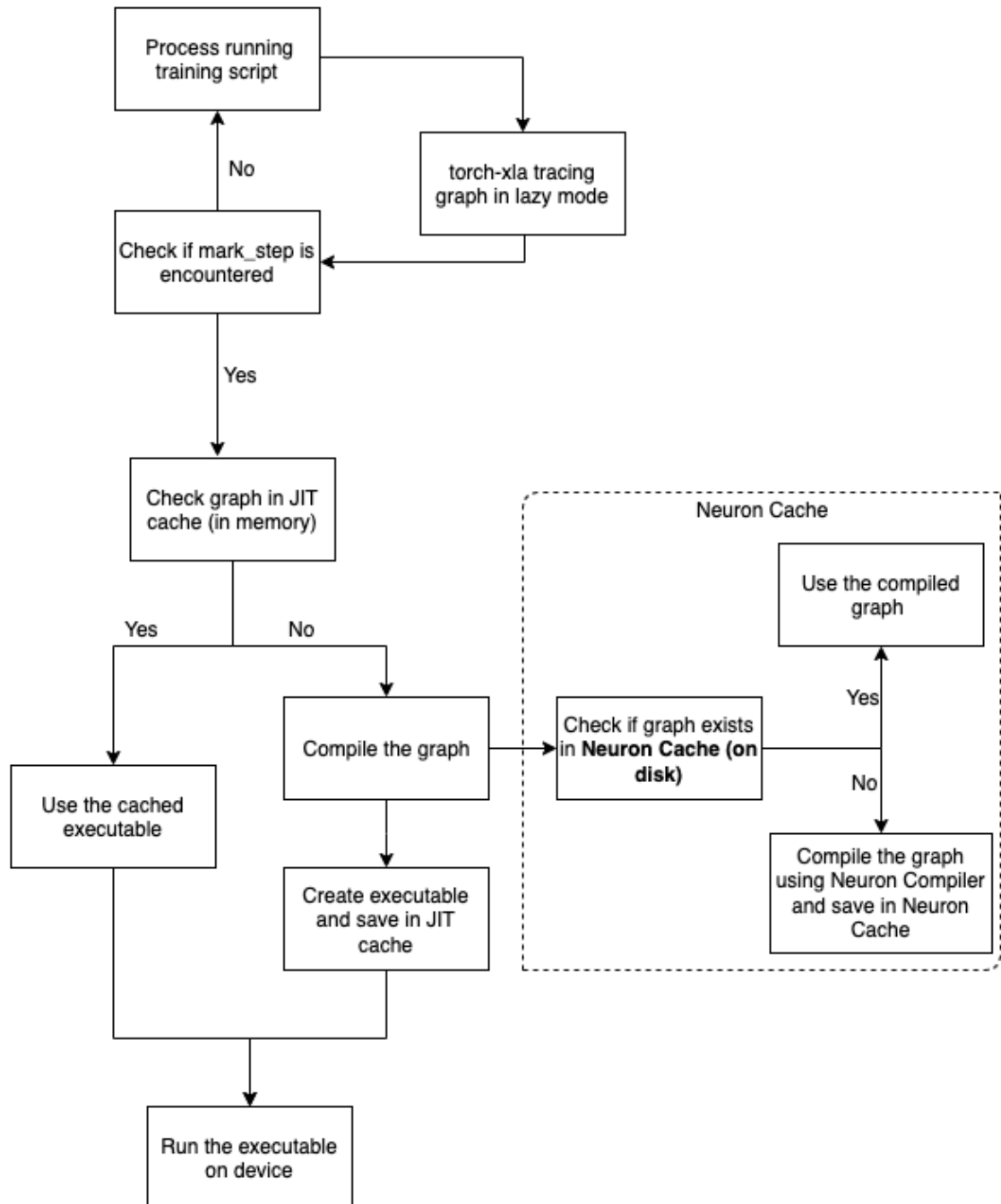
PyTorch Neuron (`torch-neuronx`) uses `torch-xla`, and `torch-xla` operates in lazy mode. In other words, every operation in training script is recorded in a graph. The graph is executed only when the results are requested by the user when they use `print` or `xm.mark_step`. Requesting results tells `torch-xla` that the recorded graph needs to be executed.

Before executing the graph on a Neuron device, `torch-xla` would call Neuron Compiler (`neuronx-cc`) to compile the graph into Neuron specific graph. Then the graph is executed on the *NeuronCore/s*. Compiling the graph involves running optimizations that can make use of the *NeuronCore/s* efficiently. Running these optimizations can be expensive and can result in long compile times. To save the users from compiling these graphs at every iteration, `torch-xla` maintains an in-memory cache called Just in Time (JIT) cache. When the user re-runs the same graph (eg. 2nd iteration of the training run), `torch-xla` would check in this JIT cache and re-use the cached compilation result, thereby avoiding the wait times.

Since the JIT cache is an in-memory cache, it needs to be constructed every time the training script is run. Hence, if the user re-runs the training script, a new JIT cache is created. This causes a compilation for the first training graph. To avoid such compilations across training runs, PyTorch Neuron (`torch-neuronx`) has built an on-disk **Neuron Persistent Cache**. Since this cache is on-disk, its persistent across training runs. So now, when a graph is compiled for the first time, the compilation result is saved in **Neuron Persistent Cache**. When the user re-runs the training script, since the JIT cache is not ready, it would send the graph for compilation. PyTorch Neuron (`torch-neuronx`) would then check if the compiled result is present in the **Neuron Persistent Cache**, if yes, it would return with the compiled result. This on-disk cache thereby avoids compilations across training runs. This cache is enabled by default for Neuron's PyTorch/XLA flow (training) as well as transformers-neuronx LLM inference package. The default cache path is the directory `/var/tmp/neuron-compile-cache`.

Look at the diagram below on the end to end flow:





As seen from the diagram, the operations are recorded in a graph in lazy mode and only when a `mark_step` is hit, the graph is executed. Before execution, the graph passes through two caches to check if we have compiled the graph sometime in the past. If yes, we reuse the compilation result and execute with it. This avoids duplicate compilations. One thing to note, both JIT cache and Neuron Cache are complementary to each other. JIT cache prevents duplicate compilation within a run and Neuron Cache prevents duplicate compilations across training runs. For example, within a training script, we have a training loop that iterates through the dataset. The first iteration would trace a unique graph and the following iteration would trace a graph that is similar to the first one. In this case, the subsequent iterations would hit the JIT cache and reuse the result. However, to save users from compiling for the first iteration graph, Neuron Persistent Cache would be used. In this case, the very first time when the script is run, the Neuron Persistent Cache would be updated. Going forward when we re-run the training script, compilation results from

Neuron Persistent Cache would be used.

To better understand how Neuron Persistent Cache works, consider the example below:

```
import torch
import torch_xla
import torch_xla.core.xla_model as xm
device = xm.xla_device()
t1 = torch.randn(3, 3).to(device)
t2 = t1 / 0.5
x = t2.cpu()
```

Running the above example produces the following logs:

```
2023-08-25 21:51:36.000433: INFO ||NCC_WRAPPER||: Compile cache path: /var/tmp/neuron-
↪compile-cache
.
Compiler status PASS
```

Re-running the above script would fetch the graph from the neuron cache and you would see logs as follows:

```
2023-08-25 21:52:23.000451: INFO ||NCC_WRAPPER||: Compile cache path: /var/tmp/neuron-
↪compile-cache
2023-08-25 21:52:23.000453: INFO ||NCC_WRAPPER||: Using a cached neff at /var/tmp/neuron-
↪compile-cache/neuronxcc-2.8.0.25+a3ad0f342/MODULE_198775565831884870+d41d8cd9/model.
↪neff. Exiting with a successfully compiled graph.
```

As you can see, the next run picks the compiled graph from cache, thereby saving the compilation time. The cache uses hash of the Neuron compiler flags and XLA graph as the key. If the Neuron compiler version or XLA graph changes, you will see recompilation. Examples of changes that would cause XLA graph change include:

- Model type and size
- Batch size
- Optimizer and optimizer hyperparameters
- Location of `xm.mark_step()`

To keep cache size small and to enable weights/parameters updates without recompilation, only the compute graphs are cached when using transformers-neuronx (weights/parameters are inputs to the compute graphs) and training flow using torch-neuronx's XLA (weights/parameters are inputs and outputs of the compute graphs). Note that this caching mechanism doesn't apply to the torch-neuronx trace API where the weights/parameters are frozen and converted to constants, then compiled together with the compute operations (traced graphs with frozen weights/parameters are not cached).

All compilation results are saved in the cache. To disable the cache, you can pass `--no_cache` option via `NEURON_CC_FLAGS`:

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --no_cache'
```

The default cache path is the directory `/var/tmp/neuron-compile-cache`. To change the cache's location, pass `cache_dir=<cache_url>` option via `NEURON_CC_FLAGS` or `NEURON_COMPILE_CACHE_URL=<cache_url>` environment variables:

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --cache_dir=
↪<cache URL>'
```

```
os.environ['NEURON_COMPILE_CACHE_URL'] = '<cache_url>'
```

The cache URL specified using `--cache_dir` is prioritized over that specified using `NEURON_COMPILE_CACHE_URL` if both are set. If `<cache_url>` starts with `s3://`, it will use the AWS S3 URL as the cache location, provided that the corresponding S3 bucket exists and is both readable and writeable.

You can change the verbose level of the compiler by adding `log_level` to either `WARNING`, `INFO` or `ERROR`. This can be done as follows:

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --log_
↪level=INFO'
```

A graph compilation can fail because of a compilation error or an environment issue (for example, compilation is interrupted by `ctrl-C`). The graph would be marked as failed and subsequent rerun would encounter message like below:

```
INFO ||NCC_WRAPPER||: Got a cached failed neff at /var/tmp/neuron-compile-cache/
↪neuronxcc-2.8.0.25+a3ad0f342/MODULE_12486829708343293975+d41d8cd9/model.neff. Will_
↪skip compilation, please set --retry_failed_compilation for recompilation.
```

To retry compilation, add `--retry_failed_compilation` in `NEURON_CC_FLAGS` environment variable. When the script is rerun, all the previously failed compilations are recompiled and fresh results are saved in the cache.

```
os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --retry_failed_
↪compilation'
```

Note that all flags demonstrated above will be parsed by a tool called `neuron_cc_wrapper`, which is a wrapper over Neuron Compiler CLI to provide caching mechanism. All these flags will not be passed into Neuron Compiler CLI.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch NeuronX Profiling API

The profiler provides a method to generate a context manager to capture trace events at the operator or runtime level.

```
torch_neuronx.experimental.profiler.profile(port=9012, ms_duration=60000,
                                             neuron_tensorboard_plugin_dir='logs/plugins/neuron',
                                             profile_type='operator', auto_start=True,
                                             delete_working=True)
```

The `torch_neuronx.experimental.profiler.profile()` method returns a profile context manager object. This object doesn't need to be used directly, as default options are set to auto capture events based on the `profile_type`.

The context manager will wrap around the entire model and training/inference loop. The context-manager is backwards-compatible with the `torch_xla.debug.profiler`

### Required Arguments

None

### Optional Keyword Arguments

#### Keyword Arguments

- **port** (*int*) – Port to run the profiling GRPC server on. Default is 9012.

- **ms\_duration** (*int*) – This defines how long the profiler will capture the HLO artifacts from the model to view in the profiler. The unit is in milliseconds. The default value is 60000 ms, or 1 minute.
- **neuron\_tensorboard\_plugin\_dir** (*str*) – The directory the neuron tensorboard plugin will file write to. This will be logs/plugins/neuron by default/
- **profile\_type** (*str*) – There is “trace” and “operator”. “trace” is the Torch Runtime Trace Level, while “operator” is the Model Operator Trace Level. Default is “operator”
- **auto\_start** (*bool*) – If set to true, the profiler will start profiling immediately. If set to false, the profiler can be set to start at a later condition. Refer to `profile.start()` for more details. Default is True.
- **delete\_working** (*bool*) – If set to False turns off the deletion of temporary files. Default True.
- **traced\_only** (*str*) – This should be set to True if profiling a model that has been traced with `torch_neuronx.trace()`. Default is False.

**Returns**

The traced *profile*

**Return type**

~profile

`torch_neuronx.experimental.profiler.profile.start()`

The `torch_neuronx.experimental.profiler.profile.start()` method starts the profiler if not started (i.e when `auto_start=False`). This function does not take in any parameters, nor return anything.

*Required Arguments*

None

*Optional Keyword Arguments*

None

**Returns**

None

*This document is relevant for:* Inf2, Trn1, Trn2

**API Reference Guide for Training (torch-neuronx)**

- *PyTorch NeuronX neuron\_parallel\_compile CLI*
- *Neuron Persistent Cache*
- *PyTorch NeuronX Environment Variables*
- *PyTorch NeuronX Profiling API*

*This document is relevant for:* Trn1, Trn2

*This document is relevant for:* Trn1, Trn2

## Developer Guide (torch-neuronx)

This document is relevant for: Trn1, Trn2

## Developer Guide for Training with PyTorch NeuronX

### Table of Contents

- *PyTorch NeuronX*
  - *Neuron XLA device*
- *PyTorch NeuronX single-worker training/evaluation quick-start*
- *PyTorch NeuronX multi-worker data parallel training using torchrun*
- *Conversion from Distributed Data Parallel (DDP) application*
- *PyTorch NeuronX environment variables*
- *Neuron Persistent Cache for compiled graphs*
- *Number of graphs*
- *Full BF16 with stochastic rounding enabled*
- *BF16 in GPU-compatible mode without stochastic rounding enabled*
- *BF16 automatic mixed precision using PyTorch Autocast*
- *Tips and Best Practices*
  - *Understand the lazy mode in PyTorch NeuronX*
  - *Minimize the number of compilation-and-executions*
  - *Aggregate the data transfers between host CPUs and devices*
  - *Ensure common initial weights across workers*
  - *Use PyTorch/XLA's model save function*
- *FAQ*
- *Debugging and troubleshooting*

Trainium is designed to speed up model training and reduce training cost. It is available on the Trn1 and Trn2 instances. On Trn1, each Trainium accelerator has two NeuronCores (default two Logical NeuronCores), which are the main neural network compute units. On Trn2, each Trainium accelerator has 8 NeuronCores (default 4 Logical NeuronCores). The examples in this guide applies to Trn1 and can be extended to run Trn2.

PyTorch NeuronX enables PyTorch users to train their models on Trainium's NeuronCores with little code change to their training code. It is based on the [PyTorch/XLA software package](#).

This guide helps you get started with single-worker training and distributed training using PyTorch Neuron.

## PyTorch NeuronX

### Neuron XLA device

With PyTorch NeuronX the default XLA device is mapped to a *Logical NeuronCore*. By default, one Logical NeuronCore is configured by a process. To use the Neuron XLA device, specify the device as `xm.xla_device()` or `'xla'`:

```
import torch_xla.core.xla_model as xm
device = xm.xla_device()
```

or

```
device = 'xla'
```

PyTorch models and tensors can be mapped to the device as usual:

```
model.to(device)
tensor.to(device)
```

To move tensor back to CPU, do :

```
tensor.cpu()
```

or

```
tensor.to('cpu')
```

### PyTorch NeuronX single-worker training/evaluation quick-start

PyTorch NeuronX uses XLA to enable conversion of PyTorch operations to Trainium instructions. To get started on PyTorch NeuronX, first modify your *training script* to use XLA in the same manner as described in [PyTorch/XLA documentation](#) and use XLA device:

```
import torch_xla.core.xla_model as xm

device = xm.xla_device()
# or
device = 'xla'
```

The Logical NeuronCore is mapped to an XLA device. On Trainium instance, the XLA device is automatically mapped to the first available Logical NeuronCore. You can use `NEURON_RT_VISIBLE_CORES` to select specific Logical NeuronCore to use.

By default the above steps will enable the training or evaluation script to run on one Logical NeuronCore. NOTE: Each process is mapped to one NeuronCore.

Finally, add `mark_step` at the end of the training or evaluation step to compile and execute the training or evaluation step:

```
xm.mark_step()
```

These changes can be placed in control-flows in order to keep the script the same between PyTorch Neuron and CPU/GPU. For example, you can use an environment variable to disable XLA which would cause the script to run in PyTorch native mode (using CPU on Trainium instances and GPU on GPU instances):

```

device = 'cpu'
if not os.environ.get("DISABLE_XLA", None):
    device = 'xla'

...

# end of training step
if not os.environ.get("DISABLE_XLA", None):
    xm.mark_step()

```

More on the need for `mark_step` is at [Understand the lazy mode in PyTorch Neuron](#).

For a full runnable example, please see the [Single-worker MLP training on Trainium tutorial](#).

## PyTorch NeuronX multi-worker data parallel training using torchrun

Data parallel training allows you to replicate your script across multiple workers, each worker processing a proportional portion of the dataset, in order to train faster.

To run multiple workers in data parallel configuration, with each worker using one NeuronCore, first add additional imports for parallel dataloader and multi-processing utilities:

```
import torch_xla.distributed.parallel_loader as pl
```

Next we initialize the Neuron distributed context using the XLA backend for `torch.distributed`:

```
import torch_xla.distributed.xla_backend
torch.distributed.init_process_group('xla')
```

Next, replace `optimizer.step()` function call with `xm.optimizer_step(optimizer)` which adds gradient synchronization across workers before taking the optimizer step:

```
xm.optimizer_step(optimizer)
```

If you're using a distributed dataloader, wrap your dataloader in the PyTorch/XLA's `MpDeviceLoader` class which provides buffering to hide CPU to device data load latency:

```
parallel_loader = pl.MpDeviceLoader(dataloader, device)
```

Within the training code, use `xm.xrt_world_size()` to get the world size, and `xm.get_ordinal` to get the global rank of the current process.

Then run use [PyTorch torchrun](#) utility to run the script. For example, to run 32 worker data parallel training on `trn1.32xlarge`:

```
torchrun --nproc_per_node=32 <script and options>
```

To run on multiple instances, make sure to use `trn1.32xlarge` instances and use all 32 NeuronCores on each instance. For example, with two instances, on the rank-0 Trn1 host, run with `--node_rank=0` using `torchrun` utility:

```
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=0 --master_addr=<root IP> --master_
↪port=<root port> <script and options>
```

On another Trn1 host, run with `--node_rank=1` :

```
torchrun --nproc_per_node=32 --nnodes=2 --node_rank=1 --master_addr=<root IP> --master_
↪port=<root port> <script and options>
```

It is important to launch rank-0 worker with `--node_rank=0` to avoid hang.

For trn2.48xlarge, use `--nproc_per_node=64` for 64 Logical NeuronCores default (each Logical NeuronCores using two physical NeuronCores).

To train on multiple instances, it is recommended to use a `ParallelCluster`. For a `ParallelCluster` example, please see [Train a model on AWS Trn1 ParallelCluster](#).

More information about `torchrun` can be found PyTorch documentation at <https://pytorch.org/docs/stable/elastic/run.html#launcher-api>.

See the Multi-worker data-parallel MLP training using `torchrun` tutorial for a full example.

## Conversion from Distributed Data Parallel (DDP) application

Distributed Data Parallel (DDP) in `torch.distributed` module is a wrapper to help convert a single-worker training to distributed training. To convert from `torch.distributed` Distributed Data Parallel (DDP) application to PyTorch Neuron, first convert the application back to single-worker training, which simply involves removing the DDP wrapper, for example `model = DDP(model, device_ids=[rank])`. After this, follow the previous section to change to multi-worker training.

## PyTorch NeuronX environment variables

Environment variables allow modifications to PyTorch Neuron behavior without requiring code change to user script. See *PyTorch Neuron environment variables* for more details.

## Neuron Persistent Cache for compiled graphs

See *Neuron Persistent Cache for compiled graphs*

## Number of graphs

PyTorch/XLA converts PyTorch's eager mode execution to lazy-mode graph-based execution. During this process, there can be multiple graphs compiled and executed if there are extra mark-steps or functions with implicit mark-steps. Additionally, more graphs can be generated if there are different execution paths taken due to control-flows.

## Full BF16 with stochastic rounding enabled

Previously, on `torch-neuronx` 2.1 and earlier, the environmental variables `XLA_USE_BF16` or `XLA_DOWNCAST_BF16` provided full casting to BF16 with stochastic rounding enabled by default. These environmental variables are deprecated in `torch-neuronx` 2.5, although still functional with warnings. To replace `XLA_USE_BF16` or `XLA_DOWNCAST_BF16` with stochastic rounding on Neuron, set `NEURON_RT_STOCHASTIC_ROUNDING_EN=1` and use the `torch.nn.Module.to` method to cast model floating-point parameters and buffers to data-type BF16 as follows:



```
os.environ["NEURON_RT_STOCHASTIC_ROUNDING_EN"] = "1"
```

```
# model is created
model.to(torch.bfloat16)
```

Stochastic rounding is needed to enable faster convergence for full BF16 model.

If the loss is to be kept in FP32, initialize it with `dtype=torch.float` as follows:

```
running_loss = torch.zeros(1, dtype=torch.float).to(device)
```

Similarly, if the optimizer states are to be kept in FP32, convert the gradients to FP32 before optimizer computations:

```
grad = p.grad.data.float()
```

For a full example, please see the *PyTorch Neuron BERT Pretraining Tutorial (Data-Parallel)*, which has been updated to use `torch.nn.Module.to` instead of `XLA_DOWNCAST_BF16`.

## BF16 in GPU-compatible mode without stochastic rounding enabled

Full BF16 training in GPU-compatible mode would enable faster convergence without the need for stochastic rounding, but would require a FP32 copy of weights/parameters to be saved and used in the optimizer. To enable BF16 in GPU-compatible mode without stochastic rounding enabled, use the `torch.nn.Module.to` method to cast model floating-point parameters and buffers to data-type `bfloat16` as follows without setting `NEURON_RT_STOCHASTIC_ROUNDING_EN=1`:

```
# model is created
model.to(torch.bfloat16)
```

In the initializer of the optimizer, for example AdamW, you can add code like the following code snippet to make a FP32 copy of weights:

```
# keep a copy of weights in highprec
self.param_groups_highprec = []
for group in self.param_groups:
    params = group['params']
    param_groups_highprec = [p.data.float() for p in params]
    self.param_groups_highprec.append({'params': param_groups_highprec})
```

In the *PyTorch Neuron BERT Pretraining Tutorial (Data-Parallel)*, this mode can be enabled by passing `--optimizer=AdamW_FP32ParamsCopy` option to `dp_bert_large_hf_pretrain_hdf5.py` and setting `NEURON_RT_STOCHASTIC_ROUNDING_EN=0` (or leave it unset).

## BF16 automatic mixed precision using PyTorch Autocast

By default, the compiler automatically casts internal FP32 operations to BF16. You can disable this and allow PyTorch's BF16 automatic mixed precision function (`torch.autocast`) to do the casting of certain operations to operate in BF16.

To enable PyTorch's BF16 mixed-precision, first turn off the Neuron compiler auto-cast:

```
os.environ["NEURON_CC_FLAGS"] = "--auto-cast=none"
```

Next, per recommendation from official PyTorch [torch.autocast documentation](#), place only the forward-pass of the training step in the `torch.autocast` scope with `xla` device type:

```
with torch.autocast(dtype=torch.bfloat16, device_type='xla'):
    # forward pass
```

The device type is XLA because we are using PyTorch-XLA’s autocast backend. The PyTorch-XLA [autocast mode source code](#) lists which operations are casted to lower precision BF16 (“lower precision fp cast policy” section), which are maintained in FP32 (“fp32 cast policy”), and which are promoted to the widest input types (“promote” section).

Example showing the original training code snippet:

```
def train_loop_fn(train_loader):
    for i, data in enumerate(train_loader):
        inputs = data[0]
        labels = data[3]
        outputs = model(inputs, labels=labels)
        loss = outputs.loss / flags.grad_acc_steps
        loss.backward()
        optimizer.step()
        xm.mark_step()
```

The following shows the training loop modified to use BF16 autocast:

```
os.environ["NEURON_CC_FLAGS"] = "--auto-cast=none"

def train_loop_fn(train_loader):
    for i, data in enumerate(train_loader):
        torch.cuda.is_bf16_supported = lambda: True
        with torch.autocast(dtype=torch.bfloat16, device_type='xla'):
            inputs = data[0]
            labels = data[3]
            outputs = model(inputs, labels=labels)
            loss = outputs.loss / flags.grad_acc_steps
            loss.backward()
            optimizer.step()
            xm.mark_step()
```

For a full example of BF16 mixed-precision, see [PyTorch Neuron BERT Pretraining Tutorial \(Data-Parallel\)](#).

See official PyTorch documentation for more details about `torch.autocast` .

## Tips and Best Practices

### Understand the lazy mode in PyTorch NeuronX

One significant difference between PyTorch NeuronX and native PyTorch is that the PyTorch NeuronX system runs in lazy mode while the native PyTorch runs in eager mode. Tensors in lazy mode are placeholders for building the computational graph until they are materialized after the compilation and evaluation are complete. The PyTorch NeuronX system builds the computational graph on the fly when you call PyTorch APIs to build the computation using tensors and operators. The computational graph gets compiled and executed when `xm.mark_step()` is called explicitly or implicitly by `pl.MpDeviceLoader/pl.ParallelLoader`, or when you explicitly request the value of a tensor such as by calling `loss.item()` or `print(loss)`.

## Minimize the number of compilation-and-executions

For best performance, you should keep in mind the possible ways to initiate compilation-and-executions as described in Understand the lazy mode in PyTorch/XLA and should try to minimize the number of compilation-and-executions. Ideally, only one compilation-and-execution is necessary per training iteration and is initiated automatically by `pl.MpDeviceLoader/pl.ParallelLoader`. The `MpDeviceLoader` is optimized for XLA and should always be used if possible for best performance. During training, you might want to examine some intermediate results such as loss values. In such case, the printing of lazy tensors should be wrapped using `xm.add_step_closure()` to avoid unnecessary compilation-and-executions.

## Aggregate the data transfers between host CPUs and devices

For best performance, you may try to aggregate the data transfers between host CPUs and devices. For example, increasing the value for `batches_per_execution` argument when instantiating `MpDeviceLoader` can help increase performance for certain where there's frequent host-device traffic like ViT as described in a [blog](#). NOTE: Increasing `batches_per_execution` value would delay the mark-step for multiple batches specified by this value, increasing graph size and could lead to out-of-memory (device OOM) error.

## Ensure common initial weights across workers

To achieve best accuracy during data parallel training, all workers need to have the same initial parameter states. This can be achieved by using the same seed across the workers. In the case of HuggingFace library, the `set_seed` function can be used. (<https://github.com/pytorch/xla/issues/3216>).

## Use PyTorch/XLA's model save function

To avoid problems with saving and loading checkpoints, make sure you use PyTorch/XLA's model save function to properly checkpoint your model. For more information about the function, see `torch_xla.core.xla_model.save` in the *PyTorch on XLA Devices* documentation.

When training using multiple devices, `xla_model.save` can result in high host memory usage. If you see such high usage causing the host to run out of memory, please use `torch_xla.utils.serialization.save`. This would save the model in a serialized manner. When saved using the `serialization.save` api, the model should be loaded using `serialization.load` api. More information on this here: [Saving and Loading Tensors](#)

## FAQ

### Debugging and troubleshooting

To debug on PyTorch Neuron, please follow the debug guide.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## How to debug models in PyTorch NeuronX

### Table of Contents

- [Printing metrics](#)
- [Printing tensors](#)
- [Use mark\\_step](#)
- [Using Eager Debug Mode](#)
- [Profiling model run](#)
- [Snapshotting With Torch-Neuronx 2.1](#)
- [Snapshotting with Torch-Neuronx 1.13](#)
  - [Snapshot FAQs:](#)

Torch-XLA evaluates operations lazily, which means it builds a symbolic graph in the background and the graph is executed in hardware only when the users request (print) for the output or a mark\_step is encountered. To effectively debug training scripts with torch-xla, please use one of the approaches mentioned below:

### Printing metrics

Torch-xla provides a utility that records metrics of different sections of the code. These metrics can help figure out things like: How much time is spent in compilation? How much time is spent in execution? To check the metrics:

1. Import metrics: `import torch_xla.debug.metrics as met`
2. Print metrics at the end of the step: `print(met.metrics_report())`

Printing metrics should produce an output that looks like this:

```
Metric: CompileTime
  TotalSamples: 1
  Accumulator: 09s969ms486.408us
  Percentiles: 1%=09s969ms486.408us; 5%=09s969ms486.408us; 10%=09s969ms486.408us; 20
↪%=09s969ms486.408us; 50%=09s969ms486.408us; 80%=09s969ms486.408us; 90%=09s969ms486.
↪408us; 95%=09s969ms486.408us; 99%=09s969ms486.408us
.....
Metric: ExecuteTime
  TotalSamples: 1
  Accumulator: 186ms062.970us
  Percentiles: 1%=186ms062.970us; 5%=186ms062.970us; 10%=186ms062.970us; 20%=186ms062.
↪970us; 50%=186ms062.970us; 80%=186ms062.970us; 90%=186ms062.970us; 95%=186ms062.970us;
↪99%=186ms062.970us
.....
Metric: TensorsGraphSize
  TotalSamples: 1
  Accumulator: 9.00
  Percentiles: 1%=9.00; 5%=9.00; 10%=9.00; 20%=9.00; 50%=9.00; 80%=9.00; 90%=9.00; 95%=9.
↪00; 99%=9.00
Metric: TransferFromServerTime
  TotalSamples: 2
```

(continues on next page)

(continued from previous page)

```

Accumulator: 010ms130.597us
ValueRate: 549ms937.108us / second
Rate: 108.372 / second
Percentiles: 1%=004ms948.602us; 5%=004ms948.602us; 10%=004ms948.602us; 20%=004ms948.
→602us; 50%=006ms181.995us; 80%=006ms181.995us; 90%=006ms181.995us; 95%=006ms181.995us; 9
→99%=006ms181.995us
Metric: TransferToServerTime
TotalSamples: 6
Accumulator: 061ms698.791us
ValueRate: 007ms731.182us / second
Rate: 0.665369 / second
Percentiles: 1%=006ms848.579us; 5%=006ms848.579us; 10%=006ms848.579us; 20%=007ms129.
→666us; 50%=008ms940.718us; 80%=008ms496.166us; 90%=024ms636.413us; 95%=024ms636.413us; 9
→99%=024ms636.413us
Metric: TransferToServerTransformTime
TotalSamples: 6
Accumulator: 011ms835.717us
ValueRate: 001ms200.844us / second
Rate: 0.664936 / second
Percentiles: 1%=108.403us; 5%=108.403us; 10%=108.403us; 20%=115.676us; 50%=167.399us; 8
→80%=516.659us; 90%=010ms790.400us; 95%=010ms790.400us; 99%=010ms790.400us
.....
Counter: xla::_copy_from
Value: 7
Counter: xla::addmm
Value: 2
Counter: xla::empty
Value: 5
Counter: xla::t
Value: 2
....
Metric: XrtCompile
TotalSamples: 1
Accumulator: 09s946ms607.609us
Mean: 09s946ms607.609us
StdDev: 000.000us
Percentiles: 25%=09s946ms607.609us; 50%=09s946ms607.609us; 80%=09s946ms607.609us; 90
→%=09s946ms607.609us; 95%=09s946ms607.609us; 99%=09s946ms607.609us
Metric: XrtExecute
TotalSamples: 1
Accumulator: 176ms932.067us
Mean: 176ms932.067us
StdDev: 000.000us
Percentiles: 25%=176ms932.067us; 50%=176ms932.067us; 80%=176ms932.067us; 90%=176ms932.
→067us; 95%=176ms932.067us; 99%=176ms932.067us
Metric: XrtReadLiteral
TotalSamples: 2
Accumulator: 608.578us
Mean: 304.289us
StdDev: 067.464us
Rate: 106.899 / second
Percentiles: 25%=236.825us; 50%=371.753us; 80%=371.753us; 90%=371.753us; 95%=371.753us;

```

(continues on next page)

(continued from previous page)

```
→ 99%=371.753us
```

As seen, you can get useful information about graph compile times/execution times. You can also know which operators are present in the graph, which operators are run on the CPU and which operators are run on an XLA device. For example, operators that have a prefix `aten::` would run on the CPU, since they do not have xla lowering. All operators with prefix `xla::` would run on an XLA device. Note: `aten` operators that do not have xla lowering would result in a graph fragmentation and might end up slowing down the entire execution. If you encounter such operators, create a request for operator support.

## Printing tensors

Users can print tensors in their script as below:

```
import os
import torch
import torch_xla
import torch_xla.core.xla_model as xm

device = xm.xla_device()
input1 = torch.randn(2,10).to(device)
# Defining 2 linear layers
linear1 = torch.nn.Linear(10,30).to(device)
linear2 = torch.nn.Linear(30,20).to(device)

# Running forward
output1 = linear1(input1)
output2 = linear2(output1)
print(output2)
```

Since `torch-xla` evaluates operations lazily, when you try to print `output2`, the graph associated with the tensor would be evaluated. When a graph is evaluated, it is first compiled for the device and executed on the selected device. Note: Each tensor would have a graph associated with it and can result in graph compilations and executions. For example, in the above script, if you try to print `output1`, a new graph is cut and you would see another evaluation. To avoid multiple evaluations, you can make use of `mark_step` (next section).

## Use mark\_step

Torch-XLA provides an api called `mark_step` which evaluates a graph collected up to that point. While this is similar to printing of an output tensor wherein a graph is also evaluated, there is a difference. When an output tensor is printed, only the graph associated with that specific tensor is evaluated, whereas `mark_step` enables all the output tensors up to `mark_step` call to be evaluated in a single graph. Hence, any tensor print after `mark_step` would be effectively free of cost as the tensor values are already evaluated. Consider the example below:

```
import os
import torch
import torch_xla
import torch_xla.core.xla_model as xm
import torch_xla.debug.metrics as met

device = xm.xla_device()
```

(continues on next page)

(continued from previous page)

```

input1 = torch.randn(2,10).to(device)
# Defining 2 linear layers
linear1 = torch.nn.Linear(10,30).to(device)
linear2 = torch.nn.Linear(30,20).to(device)

# Running forward
output1 = linear1(input1)
output2 = linear2(output1)
xm.mark_step()
print(output2)
print(output1)
# Printing the metrics to check if compilation and execution occurred
print(met.metrics_report())

```

In the printed metrics, the number of compiles and executions is only 1, even though 2 tensors are printed. Hence, to avoid multiple graph evaluations, it is recommended that you visualize tensors after a `mark_step`. You can also make use of the `add_step_closure` api for this purpose. With this api, you pass in the tensors that needs to be visualized/printed. The added tensors would then be preserved in the graph and can be printed as part of the callback function passed to the api. Here is a sample usage: [https://github.com/pytorch/xla/blob/master/test/test\\_train\\_mp\\_mnist.py#L133](https://github.com/pytorch/xla/blob/master/test/test_train_mp_mnist.py#L133)

**Note:** Graph compilations can take time as the compiler optimizes the graph to run on device.

## Using Eager Debug Mode

Eager debug mode provides a convenient utility to step through the code and evaluate operators one by one for correctness. Eager debug mode is useful to inspect your models the way you would do in eager-mode frameworks like PyTorch and Tensorflow. With Eager Debug Mode operations are executed eagerly. As soon as an operation is registered with torch-xla, it would be sent for compilation and execution. Since compiling a single operation, the time spent would be minimal. Moreover, the chances of hitting the framework compilation cache increases as models would have repeated operations throughout. Consider example 1 below:

```

# Example 1

import os
# You need to set this env variable before importing torch-xla
# to run in eager debug mode.
os.environ["NEURON_USE_EAGER_DEBUG_MODE"] = "1"

import torch
import torch_xla
import torch_xla.core.xla_model as xm
import torch_xla.debug.metrics as met

device = xm.xla_device()
input1 = torch.randn(2,10).to(device)
# Defining 2 linear layers
linear1 = torch.nn.Linear(10,30).to(device)
linear2 = torch.nn.Linear(30,20).to(device)

# Running forward
output1 = linear1(input1)
output2 = linear2(output1)

```

(continues on next page)

(continued from previous page)

```

# Printing the metrics to check if compilation and execution occurred
# Here, in the metrics you should notice that the XRTCompile and XRTExecute
# value is non-zero, even though no tensor is printed. This is because, each
# operation is executed eagerly.
print(met.metrics_report())

print(output2)
print(output1)
# Printing the metrics to check if compilation and execution occurred.
# Here the XRTCompile count should be same as the previous count.
# In other words, printing tensors did not incur any extra compile
# and execution of the graph
print(met.metrics_report())

```

As seen from the above scripts, each operator is evaluated eagerly and there is no extra compilation when output tensors are printed. Moreover, together with the on-disk Neuron persistent cache, eager debug mode only incurs one time compilation cost when the ops is first run. When the script is run again, the compiled ops will be pulled from the persistent cache. Any changes you make to the training script would result in the re-compilation of only the newly inserted operations. This is because each operation is compiled independently. Consider example 2 below:

```

# Example 2

import os
# You need to set this env variable before importing torch-xla
# to run in eager debug mode.
os.environ["NEURON_USE_EAGER_DEBUG_MODE"] = "1"

import torch
import torch_xla
import torch_xla.core.xla_model as xm
import torch_xla.debug.metrics as met

os.environ['NEURON_CC_FLAGS'] = "--log_level=INFO"

device = xm.xla_device()
input1 = torch.randn(2,10).to(device)
# Defining 2 linear layers
linear1 = torch.nn.Linear(10,30).to(device)
linear2 = torch.nn.Linear(30,20).to(device)
linear3 = torch.nn.Linear(20,30).to(device)
linear4 = torch.nn.Linear(30,20).to(device)

# Running forward
output1 = linear1(input1)
output2 = linear2(output1)
output3 = linear3(output2)

# Note the number of compiles at this point and compare
# with the compiles in the next metrics print
print(met.metrics_report())

```

(continues on next page)



(continued from previous page)

```
output4 = linear4(output3)
print(met.metrics_report())
```

Running the above example 2 script after running example 1 script, you may notice that from the start until the statement `output2 = linear2(output1)`, all the graphs would hit the persistent cache. Executing the line `output3 = linear3(output2)` would result in a new compilation for `linear3` layer only because the layer configuration is new. Now, when we run `output4 = linear4(output3)`, you would observe no new compilation happens. This is because the graph for `linear4` is same as the graph for `linear2` and hence the compiled graph for `linear2` is reused for `linear4` by the framework's internal cache.

Eager debug mode avoids the wait times involved with tensor printing because of larger graph compilation. It is designed only for debugging purposes, so when the training script is ready, please remove the `NEURON_USE_EAGER_DEBUG_MODE` environment variable from the script in order to obtain optimal performance.

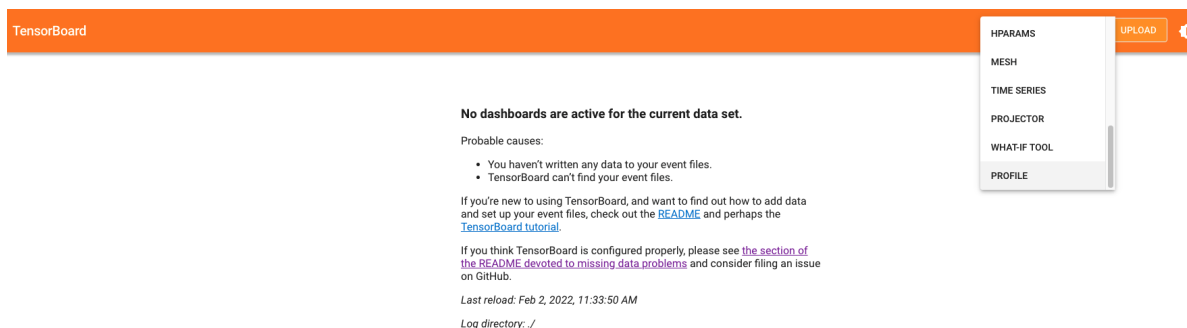
By default, in eager debug mode the logging level in the Neuron compiler is set to error mode. Hence, no logs would be generated unless there is an error. Before your first print, if there are many operations that needs to be compiled, there might be a small delay. In case you want to check the logs, switch on the INFO logs for compiler using:

```
os.environ["NEURON_CC_FLAGS"] = "--log_level=INFO"
```

## Profiling model run

Profiling model run can help to identify different bottlenecks and resolve issues faster. You can profile different sections of the code to see which block is the slowest. To profile model run, you can follow the steps below:

1. Add: `import torch_xla.debug.profiler as xp`
2. Start server. This can be done by adding the following line after creating xla device: `server = xp.start_server(9012)`
3. In a separate terminal, start tensorboard. The logdir should be in the same directory from which you run the script.



Open the tensorboard on a browser. Go to profile section in the top right. Note: you may have to install the profile plugin using: `pip install tensorboard-plugin-profile`

4. When you click on the profile, it should give an option to capture profile. Clicking on capture profile produces the following pop-up.

### Profile Service URL(s) or TPU name \*

Address Type: ☒ IP Address ☐ TPU Name

Profiling Duration (milliseconds)

1000

Host Trace (TraceMe) Level

info

Device Trace Level

enable

Python Trace Level

disable

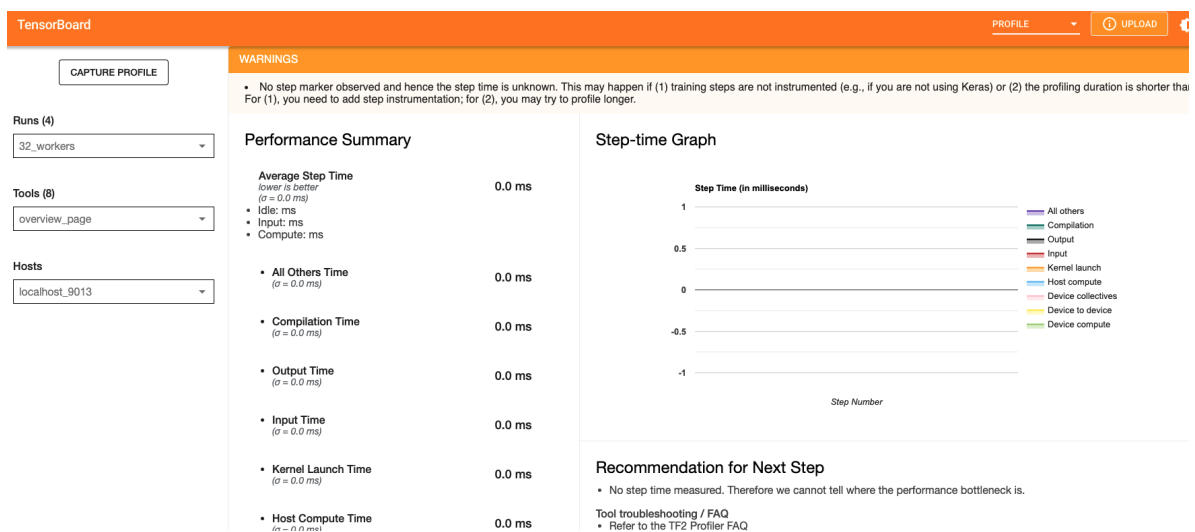
Advanced options

CAPTURE

CLOSE

In the URL enter: localhost:9012 . Port in this URL should be same as the one you gave when starting the server in the script.

- Once done, click capture and it should automatically load the following page:



- To check the profile for different blocks of code, head to `trace_viewer` under Tools (on the left column).

## Runs (4)

32\_workers ▾

## Tools (8)

memory\_profile

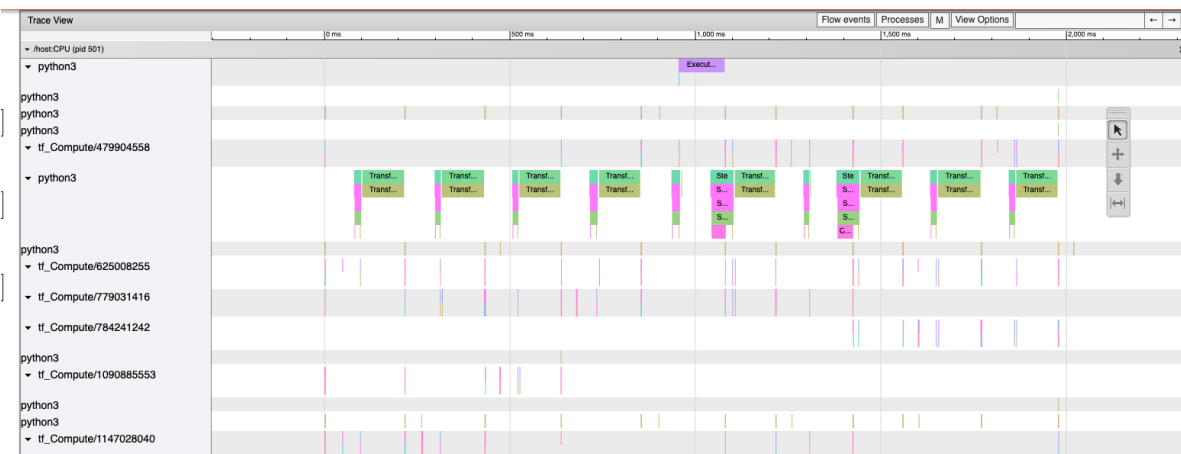
pod\_viewer

tensorflow\_stats

tf\_data\_bottleneck\_analysis

trace\_viewer

7. It should show a profile that looks like this:



Note: By default, torch-xla would time different blocks of code inside the library. However, you can also profile block of code in your scripts. This can be done by adding the code within a `xp.Trace` context as follows:

```

...
for epoch in range(total_epochs):
    inputs = torch.randn(1,10).to(device)
    labels = torch.tensor([1]).to(device)
    with xp.Trace("model_build"):
        loss = model(inputs, labels)

```

(continues on next page)

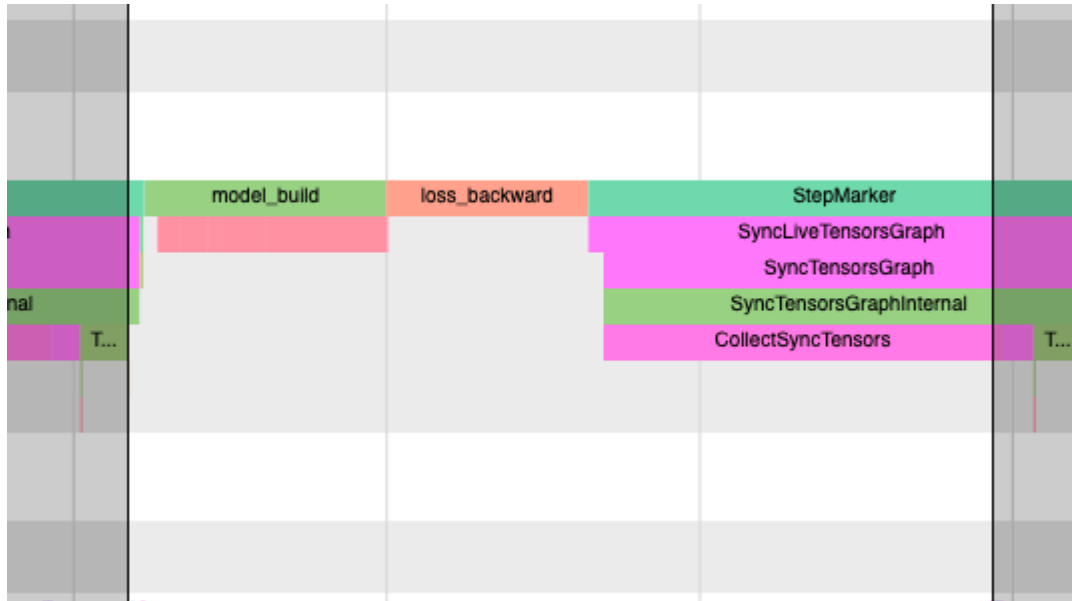
(continued from previous page)

```

with xp.Trace("loss_backward"):
    loss.backward()
....

```

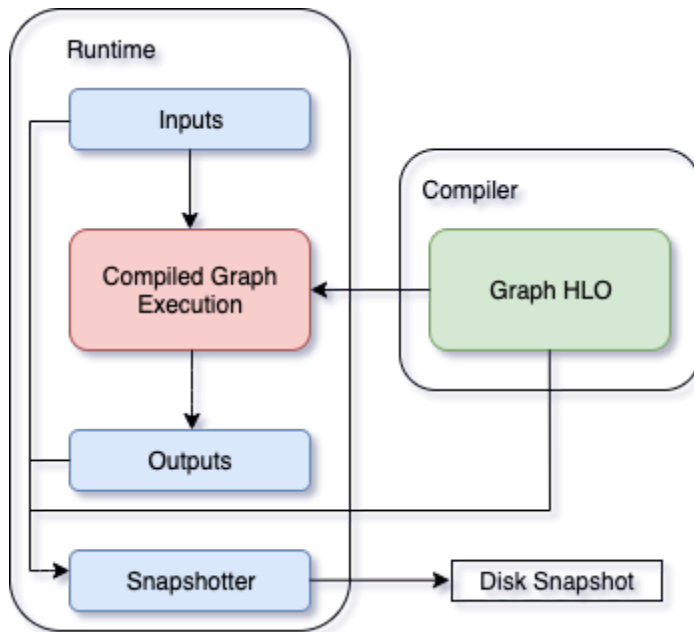
It should produce a profile that has the `model_build` and `loss_backward` section timed. This way you can time any block of script for debugging.



Note: If you are running your training script in a docker container, to view the tensorboard, you should launch the docker container using flag: `--network host` eg. `docker run --network host my_image:my_tag`

### Snapshotting With Torch-Neuronx 2.1

Snapshotting models can be used to dump debug information that can then be sent to the Neuron team. Neuron execution relies on a series of compiled graphs. Internally, graph HLOs are used as an intermediate representation which is then compiled. Then, during execution, the graph inputs are passed to the Neuron runtime, which produces outputs using the compiled graph. Snapshotting saves the inputs to a graph execution, executes the graphs, saves the outputs of the execution, and then bundles and dumps the inputs, outputs and graph HLO in one file. This is illustrated here:



This feature can be enabled using the following environment variables, which can be set at the beginning of your script as follows (./dump is the snapshot dump directory that will be created):

```
....
os.environ["XLA_FLAGS"] = "--xla_dump_hlo_snapshots --xla_dump_to=./dump"
....
```

This environment variable will produce snapshots in the ./dump folder with the extension .decomposed\_hlo\_snapshot at every iteration for every process. For example, files that look like the following would be produced.

```
SyncTensorsGraph.27737-process000000-executable000003-device000000-execution000496.
↪inputs.decomposed_hlo_snapshot
```

Note that NEURON\_FRAMEWORK\_DEBUG does not need to be set, as in torch-neuronx 1.13. Also note that NEURON\_DUMP\_HLO\_SNAPSHOT and NEURON\_NC0\_ONLY\_SNAPSHOT environment variables used in torch-neuronx 1.13 are now no longer used to control snapshot dumping.

Snapshots can take up a large amount of disk space. To avoid running out of disk space, you can limit the snapshotting for a certain rank, such as rank 0. The following example code would work with `torchrun` utility which sets the RANK environment variable for each process:

```
if os.environ.get("RANK", "0") == "0":
    os.environ["XLA_FLAGS"]="--xla_dump_hlo_snapshots --xla_dump_to=./dump"
```

or if not using `torchrun`:

```
import torch_xla.core.xla_model as xm

....
if xm.is_master_ordinal():
    os.environ["XLA_FLAGS"]="--xla_dump_hlo_snapshots --xla_dump_to=./dump"
....
```

Torch-NeuronX 2.1+ provides a `register_hlo_snapshot_callback` API to allow more control over when to dump the snapshot. By default, Torch-NeuronX 2.1+ includes the following callback function:

```
def _dump_hlo_snapshot_callback(name: str, addressable_device_index: int, execution_
    count: int) -> str:
    return 'inputs'
```

As the return value is always 'inputs', the backend will always dump snapshot files containing HLO and input data only. Recognized return value keywords are 'inputs' and 'outputs'. If the return value is an empty string '', then the backend will skip this dump. If the return value is 'inputs outputs', then the backend will dump two snapshot files for each execution, one holding inputs, and another one holding outputs.

To implement selective dumping, we can make use of the callback function's parameters `name`, `addressable_device_index`, `execution_count`, where:

- `name` is a string that stands for the HLO graph's name.
- `addressable_device_index` is an integer that refers to the index of the addressable Neuron device as one NEFF can load onto multiple addressable Neuron devices (NeuronCores) for SPMD executions. Note that this is not the same as the worker process rank in multi-process execution, in which `RANK/xm.get_ordinal()` or `LOCAL_RANK/xm.get_local_ordinal()` should be used. See examples above.
- `execution_count` is an integer that indicates the value of an internal execution counter that increments by one for each execution of a compiled graph when `HloSnapshot` dumping is requested. Note that each compiled graph maintains multiple execution counters, one for each addressable device that it loads onto.

For example, the following will dump snapshot files containing outputs at execution #2 (Note that this is graph execution number, not the iteration or step; for iteration or step, see the next example):

```
def callback(name, addressable_device_index, execution_count):
    if execution_count == 2:
        return 'outputs'
    else:
        return ''

import libneuronxla
old_callback = libneuronxla.register_hlo_snapshot_callback(callback)
```

Callback functions can be used to dump at a certain condition, such as when the global step count equals a value:

```
step = 0
def callback(name, addressable_device_index, execution_count):
    if step == 5:
        return 'inputs'
    else:
        return ''

import libneuronxla
old_callback = libneuronxla.register_hlo_snapshot_callback(callback)

...
for epoch in range(EPOCHS):
    for idx, (train_x, train_label) in enumerate(train_loader):
        step += 1
    ...
```

---

**Note:** Snapshot dumping triggered by a runtime error such as NaN is not yet available. It will be available in a feature release.

---

## Snapshotting with Torch-Neuronx 1.13

---

**Note:** If you are using Torch-NeuronX 2.1, please see *Snapshotting With Torch-Neuronx 2.1*

---

With Torch-Neuronx 1.13, the snapshotting feature can be enabled using the following environment variables, which can be set at the beginning of your script as follows.

```
....
os.environ["XLA_FLAGS"] = " --xla_dump_to=dump"
os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
os.environ["NEURON_DUMP_HLO_SNAPSHOT"] = "1"
....
```

This set of environment variables will produce snapshots under the dump folder with the extensions `.hlo.snapshot.pb` or `.decomposed_hlo_snapshot` at every iteration. For example a file that looks like the following would be produced.

```
dump/module_SyncTensorsGraph.387.pid_12643.execution_7496.hlo_snapshot.pb
```

The dumping environment variable can be set and unset at specific iterations as shown in the following example.

```
....
for step in range(STEPS):
    if step == 20:
        os.environ["NEURON_DUMP_HLO_SNAPSHOT"] = "1"
    else:
        os.environ.pop('NEURON_DUMP_HLO_SNAPSHOT', None)
    train_x = torch.randn(BATCH_SIZE, 28, 28)
    train_x = train_x.to(device)
    loss = model(train_x)
    loss.backward()
    optimizer.step()
    xm.mark_step()
....
```

Additionally, we provide capabilities to snapshot graphs automatically. The environment variables above can be set as follows:

```
....
os.environ["XLA_FLAGS"] = " --xla_dump_to=dump"
os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
os.environ["NEURON_DUMP_HLO_SNAPSHOT"] = "ON_NRT_ERROR"
....
```

When unexpected errors such as a graph execution producing NaNs occurs, snapshots will be automatically produced and execution will be terminated. Occasionally, for larger models, automatic snapshotting may not capture

snapshots due to the device memory being exhausted. In this case, the above flag can be set to `os.environ["NEURON_DUMP_HLO_SNAPSHOT"] = "ON_NRT_ERROR_HYBRID"`, this will allocate memory for inputs on both the device and host memory. In some additional cases, this may still go out of memory and may need to be set to `os.environ["NEURON_DUMP_HLO_SNAPSHOT"] = "ON_NRT_ERROR_CPU"` to avoid allocating any memory on the device at all for automatic snapshotting.

## Snapshot FAQs:

### When should I use this features?

This feature should be used when debugging errors that requires interfacing with and providing debug data to the Neuron team. Snapshotting may be redundant and unnecessary in some situations. For example, when only the model weights are necessary for debugging, methods such as checkpointing may be more convenient to use.

### What sort of data is captured with these snapshots?

The type of data captured by these snapshots may include model graphs in HLO form, weights/parameters, optimizer states, intermediate tensors and gradients. This data may be considered sensitive and this should be taken into account before sending the data to the Neuron team.

### What is the size of these snapshots?

The size of snapshots can be significant for larger models such as GPT or BERT with several GBs worth of data for larger graphs, so it is recommended to check that sufficient disk space exists before using snapshotting. In addition, limiting the amount of snapshots taken in a run will help to preserve disk space.

### Will snapshotting add overhead to my execution?

Snapshotting does add a small overhead to the execution in most cases. This overhead can be significant if snapshots are dumped at every iteration. In order to alleviate some of this overhead, in the case that snapshotting is not necessary on all cores the following environment variable can be set to collect snapshots only on the first core in torch-neuronx 1.13:

```
....
os.environ["NEURON_NC0_ONLY_SNAPSHOT"] = "1"
....
```

In torch-neuronx 2.1, use RANK environmental variable when using torchrun or `xm.is_master_ordinal()` to limit dumping to the first process (see above):

```
....
if os.environ.get("RANK", "0") == "0":
    os.environ["XLA_FLAGS"]="--xla_dump_hlo_snapshots --xla_dump_to=./dump"
....
```

or (not using torchrun):

```
import torch_xla.core.xla_model as xm

....
if xm.is_master_ordinal():
    os.environ["XLA_FLAGS"]="--xla_dump_hlo_snapshots --xla_dump_to=./dump"
....
```

In addition, checkpointing in tandem with snapshotting can be useful to reduce overhead. A checkpoint close to the problem iteration can be captured, then execution resumed with snapshots enabled.

### How can I share snapshots with the Neuron team?



These snapshots can be shared with the Neuron team via S3 bucket.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer Guide for Profiling with PyTorch NeuronX

### Table of Contents

- *Introduction*
- *Example used in this guide*
  - *Prerequisites*
  - *Environment*
  - *Setup*
- *Viewing the Trace on TensorBoard*
- *Using Named Blocks for the Trace*

## Introduction

The Neuron PyTorch profiler is a context manager wrapping around the entire model and training loop. Specifically this is the context manager: `torch_neuronx.experimental.profiler.profile`. This is a wrapper of the XLA Debug Profiler which we imported earlier as `import torch_xla.debug.profiler as xp`, and is backwards-compatible. Here are the parameters of the profiler context manager:

1. `port`: Port to run the profiling GRPC server on. Default is 9012.
2. `profile_type`: There is “trace” and “operator”. “trace” is the Torch Runtime Trace Level, while “operator” is the Model Operator Trace Level.
3. `ms_duration`: This defines how long the profiler will capture the HLO artifacts from the model to view in the profiler. The unit is in milliseconds.
4. `neuron_tensorboard_plugin_dir`: The directory the neuron tensorboard plugin will file write to (NB: Assumes that the tensorboard logdir=“log/”)
5. `delete_working`: If set to False turns off the deletion of temporary files (default True)

We move the model to the xla device *inside the context manager*. This is important, as this allows the profiler to collect the operations and processes from the `neuronx-cc` compiler artifacts. If the model is moved to the xla device outside of the context manager, the profiling won’t work.

---

**Note:** The warnings about the `XLA_IR_DEBUG` and `XLA_HLO_DEBUG` env vars not being set can be ignored for the most part. This warning only comes into play when compiling the model for Neuron outside of the profiler context manager.

---

After running this script, notice a `./logs` directory has been created. It contains the TensorBoard logs including the profiler views.

## Example used in this guide

We will use the following code sample to describe in detail how to use the Neuron PyTorch profiling API.

## Prerequisites

1. Initial Trn1 setup for PyTorch (torch-neuronx) has been done

## Environment

```
#activate python virtual environment and install tensorboard_plugin_neuron
source ~/aws_neuron_venv_pytorch_p38/bin/activate
pip install tensorboard_plugin_neuronx

#create work directory for the Neuron Profiling tutorials
mkdir -p ~/neuron_profiling_tensorboard_examples
cd ~/neuron_profiling_tensorboard_examples
```

## Setup

Create a new working directory:

```
mkdir simple_demo
cd simple_demo
```

Save the following code as demo.py:

```
import os

import torch
import torch.nn as nn
import torch.nn.functional as F

# XLA imports
import torch_xla
import torch_xla.core.xla_model as xm
import torch_xla.debug.profiler as xp

import torch_neuronx
from torch_neuronx.experimental import profiler

os.environ["NEURON_CC_FLAGS"] = "--cache_dir=./compiler_cache"

# Global constants
EPOCHS = 10

# Declare 3-layer MLP Model
class MLP(nn.Module):
    def __init__(self, input_size = 10, output_size = 2, layers = [5, 5]):
```

(continues on next page)

(continued from previous page)

```

    super(MLP, self).__init__()
    self.fc1 = nn.Linear(input_size, layers[0])
    self.fc2 = nn.Linear(layers[0], layers[1])
    self.fc3 = nn.Linear(layers[1], output_size)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return F.log_softmax(x, dim=1)

def main():
    # Fix the random number generator seeds for reproducibility
    torch.manual_seed(0)

    # XLA: Specify XLA device (defaults to a NeuronCore on Trn1 instance)
    device = xm.xla_device()

    # Start the profiler context-manager
    with torch_neuronx.experimental.profiler.profile(
        port=9012,
        profile_type='trace',
        ms_duration=15000 ) as profiler:

        # IMPORTANT: the model has to be transferred to XLA within
        # the context manager, otherwise profiling won't work
        model = MLP().to(device)
        optimizer = torch.optim.SGD(model.parameters(), lr=0.01)
        loss_fn = torch.nn.NLLLoss()

        # start training loop
        print('-----Training -----')
        model.train()
        for epoch in range(EPOCHS):
            optimizer.zero_grad()
            train_x = torch.randn(1,10).to(device)
            train_label = torch.tensor([1]).to(device)

            #forward
            loss = loss_fn(model(train_x), train_label)

            #back
            loss.backward()
            optimizer.step()

            # XLA: collect ops and run them in XLA runtime
            xm.mark_step()

        print('-----End Training -----')

if __name__ == '__main__':

```

(continues on next page)

(continued from previous page)

```
main()
```

Then run it!

```
python demo.py
```

## Viewing the Trace on TensorBoard

To view the TensorBoard logs, run `tensorboard --logdir=./logs`

**Note:** Depending on TensorBoard version `--load_fast=false` might be an additional parameter to add to view the trace.

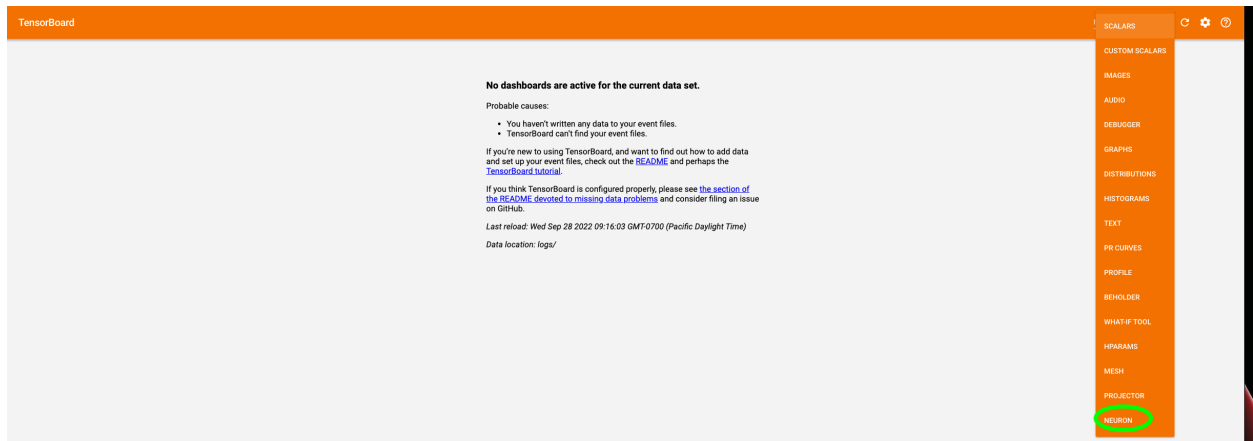
Take note of the port (usually 6006) and enter `localhost:<port>` into the local browser (assuming port forwarding is set up properly):

```
TensorBoard 1.15.0 at http://ip-172-31-33-242:6006/ (Press CTRL+C to quit)
```

Once `localhost:<port>` is entered, verify that the “NEURON” view is shown:

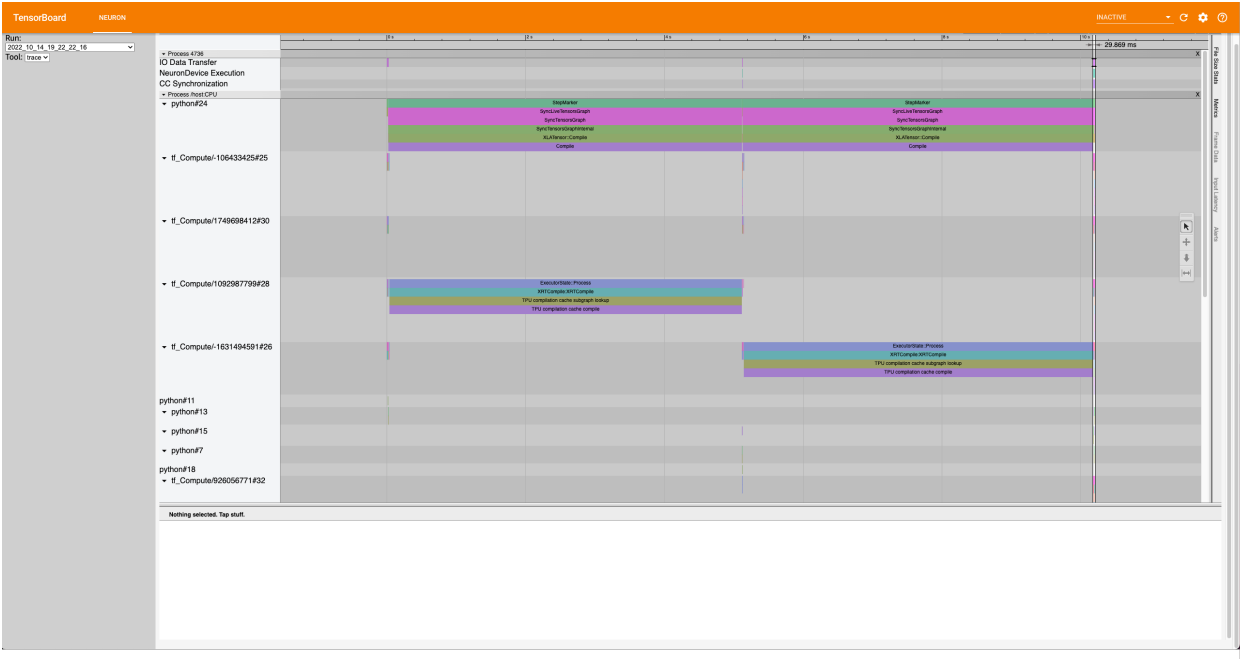


If “NEURON” isn’t shown on the top left hand side, select “NEURON” from the drop down on the top right hand side

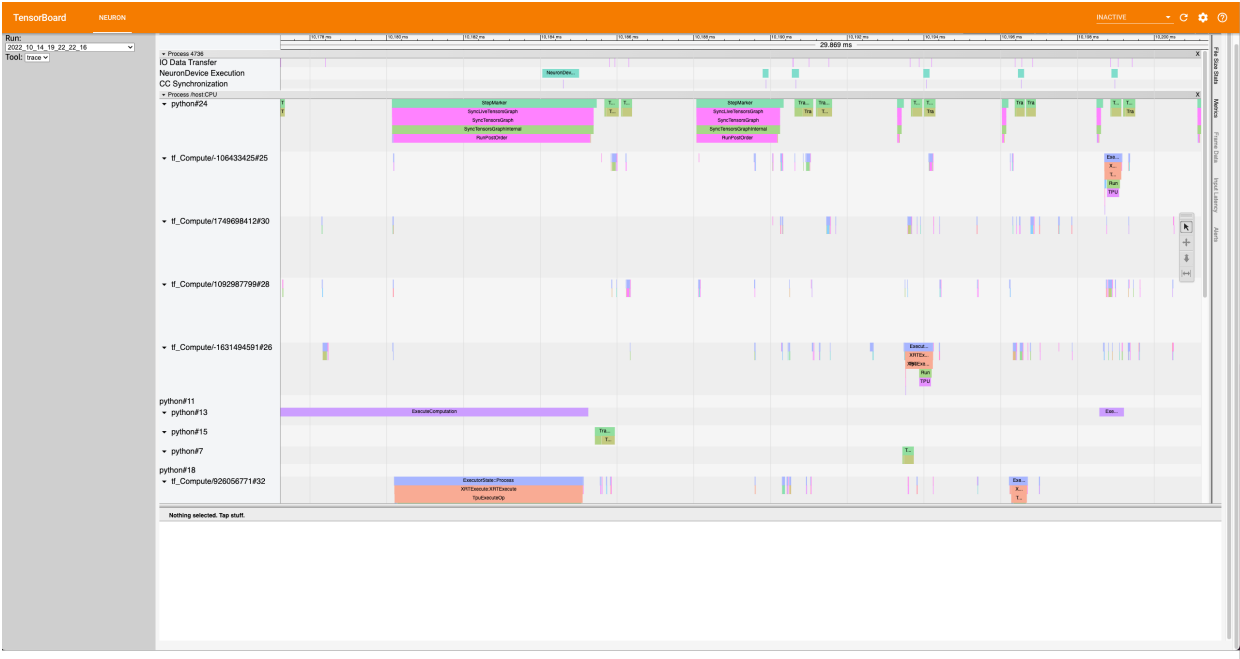


On the Left Hand Side, there are two dropdown menus: Run & Tool.





After zooming in the trace should look like this:



Notice on the top, there is a StepMarker process followed by NeuronDevice Execution process. This correlates to the `xm.mark_step()` call which executes the collected graph of our model on Neuron. For the Operator Level Trace (“operator”), we’ll be profiling the model operators that occur on Neuron. In other words, the profiler will zoom into the NeuronDevice Execution process, if the user specifies `profile_type='trace'`.

## Using Named Blocks for the Trace

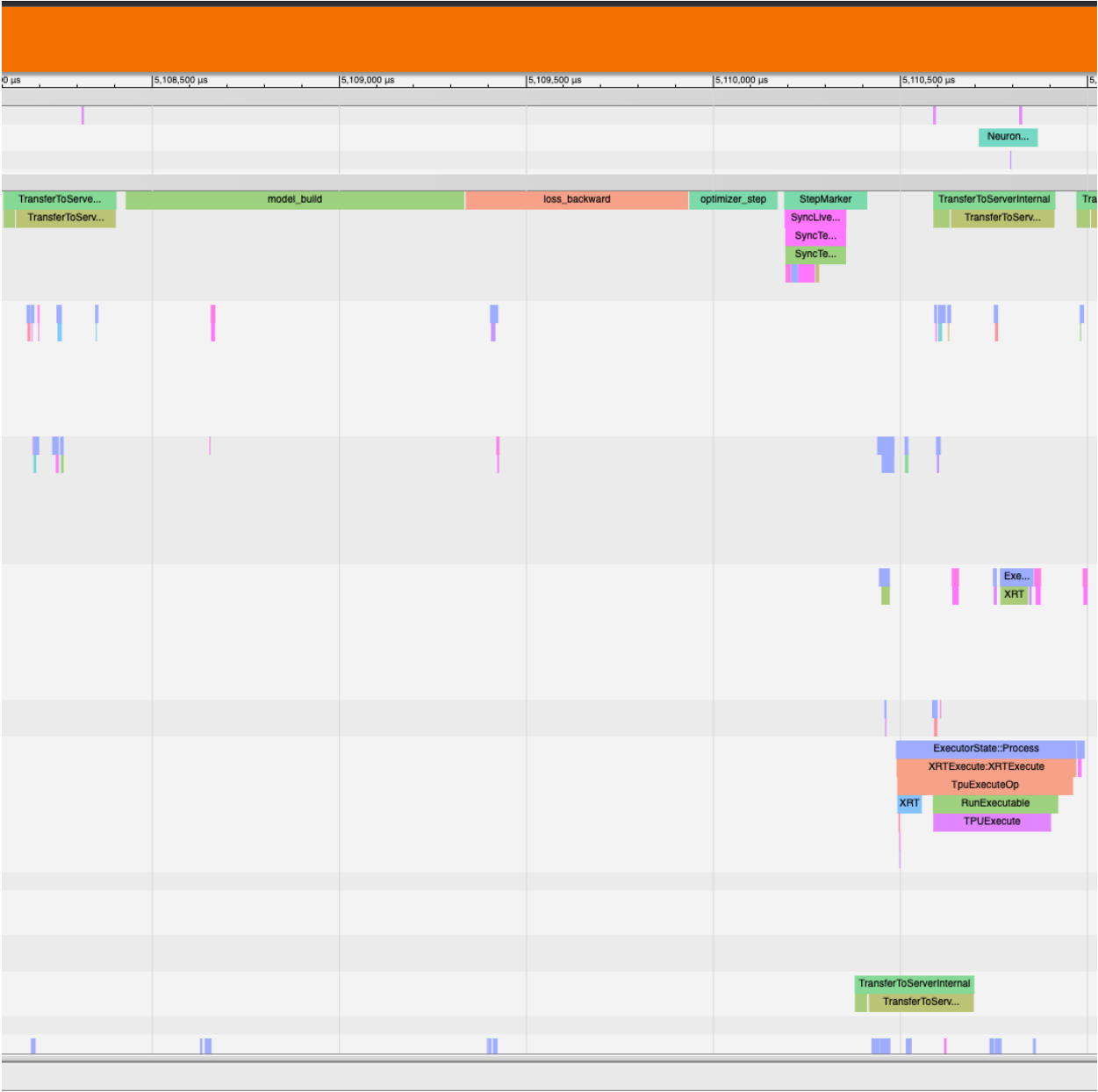
What we've produced so far is the default behavior of the profiler, however it would be more useful to profile specific blocks of our code to narrow down onto performance bottlenecks. To do this, use `xp.Trace` context manager. Replace the respective code in the training loop with the following:

```
...
optimizer.zero_grad()
train_x = torch.randn(1,10).to(device)
train_label = torch.tensor([1]).to(device)

with xp.Trace("model_build"):
    loss = loss_fn(model(train_x), train_label)
with xp.Trace("loss_backward"):
    loss.backward()
with xp.Trace("optimizer_step"):
    optimizer.step()

# XLA: collect ops and run them in XLA runtime
xm.mark_step()
...
```

Run the script, and follow the same TensorBoard steps. Afterwards, the trace should look like this:



As seen, the `model_build`, `loss_backward` and `optimizer_step` sections have been profiled.

**Note:** If you are running your training script in a docker container, to view the tensorboard, you should launch the



docker container using flag: `-network host` eg. `docker run -network host my_image:my_tag`

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer Guide

- [Developer Guide for Training with PyTorch NeuronX](#)
- [How to debug models in PyTorch NeuronX](#)
- [Developer Guide for Profiling with PyTorch NeuronX](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Misc (Training - torch-neuronx)

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch Neuron (torch-neuronx) - Supported Operators

### Table of Contents

- [Operator support](#)

## Operator support

The following list the aten operators supported by torch-neuronx.

|                                |
|--------------------------------|
| aten::_s_where                 |
| aten::_softmax                 |
| aten::_softmax_backward_data   |
| aten::_unsafe_view             |
| aten::add                      |
| aten::addcddiv_                |
| aten::addcmul                  |
| aten::addmm                    |
| aten::bernoulli_               |
| aten::bmm                      |
| aten::constant_pad_nd          |
| aten::div                      |
| aten::embedding                |
| aten::embedding_dense_backward |
| aten::empty                    |
| aten::expand                   |
| aten::fill_                    |
| aten::index_select             |

continues on next page

Table 2.1 – continued from previous page

|                                  |
|----------------------------------|
| aten::_log_softmax               |
| aten::_log_softmax_backward_data |
| aten::lt                         |
| aten::mm                         |
| aten::mul                        |
| aten::native_batch_norm          |
| aten::native_batch_norm_backward |
| aten::neg                        |
| aten::permute                    |
| aten::relu                       |
| aten::rsub                       |
| aten::select                     |
| aten::slice                      |
| aten::sqrt                       |
| aten::sum                        |
| aten::t                          |
| aten::tanh                       |
| aten::tanh_backward              |
| aten::threshold_backward         |
| aten::transpose                  |
| aten::unsqueeze                  |
| aten::view                       |
| aten::zero_                      |

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## How to prepare trn1.32xlarge for multi-node execution

EFA is a low latency transport that is used for inter-node communication. Multi-node jobs, such as distributed training, requires EFA to be enabled on every participating trn1/trn1n 32xlarge instance. Please note that EFA is currently not available on the smaller instances sizes and they cannot be used for running multi-node jobs.

trn1.32xlarge has 8 EFA devices, trn1n.32xlarge has 16 EFA devices. The rest of the document will refer to trn1.32xlarge but everything in the document also applies to trn1n.32xlarge except for the different number of EFA devices.

## Launching an instance

Before launching trn1 you need to create a security group that allows EFA traffic between the instances. Follow Step1 here: <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/efa-start.html#efa-start-security> and note the newly created security group ID. It will be used on the next step.

Determine the region, the AMI, the key and the subnet that will be used to launch trn1.

At the moment launching Trn1 instances with EFA support from the console is not recommended. The instances must be launched using AWS CLI. To launch trn1.32xlarge instance:

```
export AMI=<ami>
export SUBNET=<subnet id>
```

(continues on next page)

(continued from previous page)

```

export SG=<security group created on the previous step>
export REG=<AWS region>
export KEY=<the key>

aws ec2 run-instances --region ${REG} \
--image-id ${AMI} --instance-type trn1.32xlarge \
--key-name ${KEY} \
--tag-specifications "ResourceType=instance,Tags=[{Key=Name,Value=\"friendly name\"}]" \
--network-interfaces \
"NetworkCardIndex=0,DeviceIndex=0,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=1,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=2,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=3,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=4,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=5,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=6,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa" \
"NetworkCardIndex=7,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa"

```

Note that one of the cards is assigned DeviceIndex 0 and the rest are assigned DeviceIndex 1. Cloud-init will configure instance routing to route outgoing traffic prioritized by the device index field. I.e the outbound traffic will always egress from the interface with DeviceIndex 0. That avoids network connectivity problems when multiple interfaces are attached to the same subnet.

To launch trn1n.32xlarge instance:

```

export AMI=<ami>
export SUBNET=<subnet id>
export SG=<security group created on the previous step>
export REG=<AWS region>
export KEY=<the key>

aws ec2 run-instances --region ${REG} \
--image-id ${AMI} --instance-type trn1.32xlarge \
--key-name ${KEY} \
--tag-specifications "ResourceType=instance,Tags=[{Key=Name,Value=\"friendly name\"}]" \
--network-interfaces \
  NetworkCardIndex=0,DeviceIndex=0,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=1,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=2,DeviceIndex=2,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=3,DeviceIndex=3,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=4,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=5,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=6,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=7,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=8,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=9,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=10,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=11,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=12,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=13,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=14,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa \
  NetworkCardIndex=15,DeviceIndex=1,Groups=${SG},SubnetId=${SUBNET},InterfaceType=efa

```

## Assigning public IP address

Multi-interface instances are not assigned public IP automatically. If you require access to the newly launched `trn1` from the Internet you need to assign Elastic IP to the interface with `DeviceIndex = 0`. To find the right interface either parse the output of the instance launch command or use `describe-instances` command:

```
$ aws ec2 describe-instances --instance-ids i-01b17afa1e6021d6c
{
  "Reservations": [
    {
      "Groups": [],
      "Instances": [
        {
          "AmiLaunchIndex": 0,
          "ImageId": "ami-01257e71ecb2f431c",
          "InstanceId": "i-01b17afa1e6021d6c",
          "InstanceType": "trn1.32xlarge",
          .....
          "NetworkInterfaces": [
            {
              "Attachment": {
                "AttachTime": "2023-05-19T17:37:26.000Z",
                "AttachmentId": "eni-attach-03730388baedd4b96",
                "DeleteOnTermination": true,
                "DeviceIndex": 0,
                "Status": "attached",
                "NetworkCardIndex": 4
              },
              "Description": "",
              .....
              "InterfaceType": "efa"
            },
            {
              "Attachment": {
                "AttachTime": "2023-05-19T17:37:26.000Z",
                "AttachmentId": "eni-attach-0e1242371cd2532df",
                "DeleteOnTermination": true,
                "DeviceIndex": 0,
                "Status": "attached",
                "NetworkCardIndex": 3
              },
              "Description": "",
              .....
            }
          ]
        }
      ]
    }
  ]
}
```

The second entry in “NetworkInterfaces” in this example has “DeviceIndex” 0 and should be used to attach EIP.

## Software installation

The software required for EFA operation is distributed via aws-efa-installer package. The package is preinstalled on Neuron DLAMI. If you'd like to install the latest or if you are using your own AMI follow these steps:

```
curl -O https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz
wget https://efa-installer.amazonaws.com/aws-efa-installer.key && gpg --import aws-efa-
↪ installer.key
cat aws-efa-installer.key | gpg --fingerprint
wget https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz.sig && gpg --
↪ verify ./aws-efa-installer-latest.tar.gz.sig
tar -xvf aws-efa-installer-latest.tar.gz
cd aws-efa-installer && sudo bash efa_installer.sh --yes
cd
sudo rm -rf aws-efa-installer-latest.tar.gz aws-efa-installer
```

## Containers

aws-efa-installer package must be installed on the instance. That installs both the efa kernel module and the libraries. The libraries must be accessible to an application running inside a container. This can be accomplished by either installing aws-efa-installer package inside the container or by making on the instance library installation path available inside a container.

If installing aws-efa-installer package inside a container pass the flag that disables the kernel module installation:

```
sudo bash efa_installer.sh --yes --skip-kmod
```

The location of the libraries is distribution specific:

```
/opt/amazon/efa/lib    # Ubuntu
/opt/amazon/efa/lib64  # AL2
```

## Application execution environment

When running an application make sure the following environment variables are set:

```
FI_PROVIDER=efa
FI_EFA_USE_DEVICE_RDMA=1
FI_EFA_FORK_SAFE=1  # only required when running on AL2
```

## Appendix - trn1 instance launch example script

```
#!/bin/bash

set -e

# AWS CLI v2 Installation instructions for Linux:
# curl "https://awscli.amazonaws.com/awscli-exe-linux-x86_64.zip" -o "awscliv2.zip"
# unzip awscliv2.zip
```

(continues on next page)

(continued from previous page)

```

# sudo ./aws/install
# $ aws --version
# aws-cli/2.11.20 Python/3.11.3 Linux/5.15.0-1034-aws exe/x86_64.ubuntu.20 prompt/off
# Someone with AWS console admin privileges can create an access key ID and secret for
→ this:
# Configure credentials: aws configure

# Search the AWS AMIs for the most recent "Deep Learning Base Neuron AMI (Ubuntu 20.04)
→ <Latest_Date>"
# This one is 2023-05-17 - ami-01257e71ecb2f431c
AMI= ... # the ami
KEYNAME= ... # your key
SG= ... # the security group
SUBNET= ... # the subnet
REGION=us-west-2

# Launch instances
echo "Starting instances..."
output=$(aws ec2 --region $REGION run-instances \
--tag-specifications 'ResourceType=instance,Tags=[{Key=Name,Value=_Trainium-Big}]' \
--count 1 \
--image-id $AMI \
--instance-type trn1.32xlarge \
--key-name $KEYNAME \
--network-interfaces "NetworkCardIndex=0,DeviceIndex=0,Groups=$SG,SubnetId=$SUBNET,
→ InterfaceType=efa" \
"NetworkCardIndex=1,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=2,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=3,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=4,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=5,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=6,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa" \
"NetworkCardIndex=7,DeviceIndex=1,Groups=$SG,SubnetId=$SUBNET,InterfaceType=efa")

# Parse the output to get the instance IDs
instance_ids=$(echo $output | jq -r '.Instances[].InstanceId')
echo "Got created instance IDs: $instance_ids"

# Loop through each instance ID
public_ips=""
for instance_id in $instance_ids; do
    echo "Waiting for instance $instance_id to be running..."
    aws ec2 wait instance-running --instance-ids $instance_id --region $REGION

    echo "Creating SSH public IP newtork inteface for instance $instance_id..."
    interface_id=""
    INSTANCE_INFO=$(aws ec2 describe-instances --region $REGION --instance-ids $instance_
→ id)
    OUTPUT=$(echo "$INSTANCE_INFO" | jq -r '.Reservations[0].Instances[0].
→ NetworkInterfaces[] | "\(.Attachment.DeviceIndex),\(.NetworkInterfaceId)"')
    echo $OUTPUT

```

(continues on next page)

(continued from previous page)

```

for pair in $OUTPUT; do
    IFS="," read -r device_idx ni_id <<< $pair
    if [ "$device_idx" == "0" ]; then
        interface_id=$ni_id
        break
    fi
done
if [ "$interface_id" == "" ]; then
    exit -1
fi
echo $interface_id

echo "Checking for unassociated Elastic IPs..."
unassociated_eips=$(aws ec2 describe-addresses --region $REGION | jq -r '.Addresses[] |
↪ | select(.AssociationId == null) | .AllocationId')
if [[ -z "$unassociated_eips" ]]; then
    echo "No unassociated Elastic IPs found. Allocating new Elastic IP..."
    eip_output=$(aws ec2 allocate-address --domain vpc --region $REGION)
    eip_id=$(echo $eip_output | jq -r .AllocationId)
    echo "Allocated Elastic IP ID: $eip_id"
    eip_public_ip=$(echo $eip_output | jq -r .PublicIp)
    echo "Allocated Elastic IP Public IP: $eip_public_ip"
    echo "Note that this newly allocated Elastic IP will persist even after the
↪ instance termination"
    echo "If the Elastic IP is not going to be reused do not forget to delete it"
else
    # use the first unassociated Elastic IP found
    eip_id=$(echo "$unassociated_eips" | head -n 1)
    echo "Found unassociated Elastic IP ID: $eip_id"
    eip_public_ip=$(aws ec2 describe-addresses --allocation-ids $eip_id --region
↪ $REGION | jq -r .Addresses[0].PublicIp)
    echo "Elastic IP Public IP: $eip_public_ip"
fi
public_ips+="{eip_public_ip} "

echo "Associating Elastic IP with network interface $interface_id..."
aws ec2 associate-address --allocation-id $eip_id --network-interface-id $interface_id
↪ --region $REGION
echo "Associated Elastic IP with network interface."
done

echo "The instance has been launched.\nYou can now SSH into $public_ips with key
↪ $KEYNAME.\n"

```

**Note:** if you face connectivity issues after launching trn1\trn1n 32xlarge instance on Ubuntu, please follow the troubleshooting instructions mentioned [here](#).

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Trn1, Trn2

## PyTorch Neuron (torch-neuronx) for Training Troubleshooting Guide

**Table of contents**

- *General Troubleshooting*
- *Possible Error Conditions*
  - *Eager debug mode fails with “urllib3.exceptions.URLSchemeUnknown: Not supported URL scheme http+unix”*
  - *Eager debug mode fails with “TypeError: HTTPConnection.request() got an unexpected keyword argument ‘chunked’”*
  - *Non-Fatal Error OpKernel (‘op: “TPU\*” device\_type: “CPU”’)*
  - *XLA runtime error: “Invalid argument: Cannot assign a device for operation”*
  - *Error: “Could not start gRPC server”*
  - *Failed compilation result in the cache*
  - *Compilation errors when placing NeuronCache home directory on NFS/EFS/FSx mounted drive*
  - *Compilation error: “Expect ap datatype to be of type float32 float16 bfloat16 uint8”*
  - *NeuronCore(s) not available - Requested:1 Available:0*
  - *TDRV error “TDRV:exec\_consume\_infer\_status\_notification”*
  - *TDRV error “TDRV:tdrv\_one\_tmpbuf\_reserve Number of ONE TMPBUF pages requested exceeded the max number of pages allowed (requested: <N>, max allowed: 16).”*
  - *Could not open the ndX, close device failed, TDRV not initialized*
  - *Multiworker execution hangs during NCCL init*
  - *NRT init error “One or more engines are running. Please restart device by reloading driver”*
  - *NRT error “ERROR TDRV:kbl\_model\_add Attempting to load an incompatible model!”*
  - *NRT error “ERROR HAL:aws\_hal\_sprot\_config\_remap\_entry SPROT remap destination address must be aligned size”*
  - *NCCL warning : “NCCL WARN Timeout waiting for RX (waited 120 sec) - retrying”*
  - *RPC error: “RPC failed with status = ‘UNAVAILABLE: Socket closed’”*
  - *Error “Assertion ‘listp->slotinfo[cnt].gen <= GL(dl\_tls\_generation)’ failed” followed by ‘RPC failed with status = “UNAVAILABLE: Connection reset by peer”’*
  - *RPC connection error: “RPC failed with status = UNAVAILABLE: Connection reset by peer” not preceded by any error*
  - *Runtime errors “Missing infer\_status notification” followed by “inference timeout”*
  - *Protobuf Error “TypeError: Descriptors cannot not be created directly.”*
  - *TDRV error “Timestamp program stop timeout”*
  - *Compiler error “module ‘numpy’ has no attribute ‘asscalar’”*
  - *Import errors ‘generic\_type: type “IrValue” is already registered!’ or ‘generic\_type: type “XlaBuilder” is already registered!’*



- Import error “import \_XLAC ImportError: <>/site-packages/\_XLAC.cpython-38-x86\_64-linux-gnu.so: undefined symbol”
- Network Connectivity Issue on trn1/trn1n 32xlarge with Ubuntu
- “Too many open files” when running training job
- “undefined symbol”

This document shows common issues users may encounter while using PyTorch-Neuron and provides guidance how to resolve or work-around them.

## General Troubleshooting

For setting up EFA that is needed for multi-node training, please see *How to prepare trn1.32xlarge for multi-node execution*

For XLA-related troubleshooting notes see *How to debug models in PyTorch Neuron* and *PyTorch-XLA troubleshooting guide*.

If your multi-worker training run is interrupted, you may need to kill all the python processes (WARNING: this kills all python processes and reload the driver):

```
killall -9 python
killall -9 python3
sudo rmmmod neuron; sudo modprobe neuron
```

To turn on RT debug:

```
os.environ["NEURON_RT_LOG_LEVEL"] = "INFO"
```

To turn on Neuron NCCL debug:

```
os.environ["NCCL_DEBUG"] = "WARN"
os.environ["NCCL_DEBUG_SUBSYS"] = "ALL"
```

If some process crashed during training, you can enable core dumps using `ulimit` command:

```
ulimit -S -c unlimited
```

To see the type of signals that would cause core dumps, see <https://www.man7.org/linux/man-pages/man7/signal.7.html>.

Note that core dumps take significant amount of storage, so make sure there is enough free disk space before enabling core dumps.

On Ubuntu, if Appport is not running, core dump file name is by default “core” in the local directory. To change file location and name format, modify `/proc/sys/kernel/core_pattern` (see <https://www.kernel.org/doc/html/latest/admin-guide/sysctl/kernel.html#core-pattern> for pattern info). For example, to dump to /tmp with executable filename and process ID:

```
echo '/tmp/core.%.e.%p' | sudo tee /proc/sys/kernel/core_pattern
```

For containers, install appropriate dependencies during docker build (“`apt-get update && apt-get -y install build-essential gdb`”) and start the container with `--ulimit core=-1` to enable core dump and `-v /tmp:/tmp` to ensure core dumps to /tmp are preserved when container is stopped or deleted. Dependencies can also be installed after container is started.

On Ubuntu, core dumps can also be handled by Apport which is disabled by default. To enable Apport, run `sudo service apport start`. The `/proc/sys/kernel/core_pattern` is updated by Apport service. After a crash, look in `/var/log/apport.log` for the core dump file name, which should be located in `/var/lib/apport/coredump/`.

Once you have the core dump, you can use `gdb` to debug further (for Python applications, `<executable>` is `python` or `python3`):

```
gdb <executable> <core file>
```

If some process (i.e. XRT server) is killed due to out-of-memory on host (i.e. you see `Out of memory: Killed process <PID>` in `/var/log/syslog` or output of `dmesg`), there won't be any core dump generated. However, you can change it to kernel panic mode to trigger core dump by setting `/proc/sys/vm/panic_on_oom` to value of 1 on the host or from inside container.

On the host where you need `sudo` (this change will be reflected inside the container also):

```
echo 1 | sudo tee /proc/sys/vm/panic_on_oom
```

From inside container where `sudo` doesn't work (this change will be reflected on the host also):

```
echo 1 > /proc/sys/vm/panic_on_oom
```

## Possible Error Conditions

### Eager debug mode fails with “`urllib3.exceptions.URLSchemeUnknown: Not supported URL scheme http+unix`”

When running with eager debug mode (`NEURON_USE_EAGER_DEBUG_MODE=1`) using `torch-neuronx` and `neuronx-cc` from releases 2.19.1 and 2.20, you may see the following error:

```
urllib3.exceptions.URLSchemeUnknown: Not supported URL scheme http+unix
```

This error is due to `requests` version `>= 2.32`. While `neuronx-cc` pins `requests` package version to be less than 2.32, installing other packages like `transformers` could bring in a newer version of `requests`. To work-around this, you can pin `requests` to version 2.31.0 with the following command, which also includes `urllib3` pinning due to a related issue noted in the next note:

```
pip install requests==2.31.0 urllib3==1.26.20
```

### Eager debug mode fails with “`TypeError: HTTPConnection.request() got an unexpected keyword argument 'chunked'`”

When running with eager debug mode (`NEURON_USE_EAGER_DEBUG_MODE=1`) using `torch-neuronx` and `neuronx-cc` from releases 2.19.1 and 2.20, you may see the following error:

```
TypeError: HTTPConnection.request() got an unexpected keyword argument 'chunked'
```

This error is due to `urllib3` version `>= 2.*` and can be a dependency of `requests < 2.32`. To work-around this, you can pin `urllib3` to version 1.26.20 with the following command (which also includes `requests` pinning due to a related issue noted in the previous note):

```
pip install requests==2.31.0 urllib3==1.26.20
```

## Non-Fatal Error OpKernel ('op: "TPU\*" device\_type: "CPU"')

During execution using PyTorch Neuron, you may see these non-fatal error messages:

```
E tensorflow/core/framework/op_kernel.cc:1676] OpKernel ('op: "TPURoundRobin" device_
↪type: "CPU"') for unknown op: TPURoundRobin
E tensorflow/core/framework/op_kernel.cc:1676] OpKernel ('op: "TpuHandleToProtoKey"
↪device_type: "CPU"') for unknown op: TpuHandleToProtoKey
```

They don't affect operation of the PyTorch Neuron and can be ignored.

## XLA runtime error: "Invalid argument: Cannot assign a device for operation"

```
RuntimeError: tensorflow/compiler/xla/xla_client/xrt_computation_client.cc:490 : Check_
↪failed: session->session()->Run(session_work->feed_inputs, session_work->outputs_
↪handles, &outputs) == ::tensorflow::Status::OK() (INVALID_ARGUMENT: Cannot assign a
↪device for operation XRTAllocateFromTensor: {{node XRTAllocateFromTensor}} was
↪explicitly assigned to /job:localservice/replica:0/task:0/device:TPU:0 but available
↪devices are [ /job:localservice/replica:0/task:0/device:CPU:0, /job:localservice/
↪replica:0/task:0/device:TPU_SYSTEM:0, /job:localservice/replica:0/task:0/device:XLA_
↪CPU:0 ]. Make sure the device specification refers to a valid device.
[[XRTAllocateFromTensor]] vs. OK)
*** Begin stack trace ***
tensorflow::CurrentStackTrace()

xla::util::MultiWait::Complete(std::function<void ()> const&)

clone
*** End stack trace ***
```

The above error indicates that the framework was not able to initialize the neuron runtime. If you get the above error, check for the following:

1. No other process is taking the neuron cores. If yes, you may have to kill that process.
2. If no process is running, try reloading the driver using `sudo rmmod neuron; sudo modprobe neuron`

## Error: "Could not start gRPC server"

If you get "Could not start gRPC server" error, please check if there are any leftover python processes from a previous interrupted run and terminate them before restarting run.

```
E0207 17:22:12.592127280 30834 server_ckpt2.cc:40] {"created": "@1644254532.
↪592081429", "description": "No address added out of total 1 resolved", "file": "external/
↪com_github_grpc_grpc/src/core/ext/t
ransport/ckpt2/server/ckpt2_server.cc", "file_line": 395, "referenced_errors": [{"created":
↪"@1644254532.592078907", "description": "Failed to add any wildcard listeners", "file":
↪"external/com_github_grpc_grpc/s
rc/core/lib/iomgr/tcp_server_posix.cc", "file_line": 342, "referenced_errors": [{"created":
↪"@1644254532.592072626", "description": "Unable to configure socket", "fd": 10, "file":
↪"external/com_github_grpc_grpc/src/c
ore/lib/iomgr/tcp_server_utils_posix_common.cc", "file_line": 216, "referenced_errors": [{
```

(continues on next page)

(continued from previous page)

```

↪ "created": "@1644254532.592068939", "description": "Address already in use", "errno": 98,
↪ "file": "external/com_github_grpc_grpc/src/core/lib/iomgr/tcp_server_utils_posix_common.
↪ cc", "file_line": 189, "os_error": "Address already in use", "syscall": "bind"}]]}, {"created":
↪ "@1644254532.592078512", "description": "Unable to configure socket"
, "fd": 10, "file": "external/com_github_grpc_grpc/src/core/lib/iomgr/tcp_server_utils_posix_
↪ common.cc", "file_line": 216, "referenced_errors": [{"created": "@1644254532.592077123",
↪ "description": "Address already in
use", "errno": 98, "file": "external/com_github_grpc_grpc/src/core/lib/iomgr/tcp_server_
↪ utils_posix_common.cc", "file_line": 189, "os_error": "Address already in use", "syscall":
↪ "bind"}]]}}]]}
2022-02-07 17:22:12.592170: E tensorflow/core/distributed_runtime/rpc/grpc_server_lib.cc:
↪ 545] Unknown: Could not start gRPC server

```

### Failed compilation result in the cache

All compilation results are by default saved in Neuron Persistent Cache. If the Neuron Compiler fails to compile a graph, we save the failed result in the cache. The reason for doing so is, if the user tries to run the same script, we want the users to error out early rather than wait for the compilation to progress and see an error at the later stage. However, there could be certain cases under which a failed compilation may be do you some environment issues. One possible reason of failure could be, during compilation the process went out of memory. This can happen if you are running multiple processes in parallel such that not enough memory is available for compilation of graph. Failure due to such reasons can be easily mitigated by re-running the compilation. In case, you want to retry a failed compilation, you can do that by passing `--retry_failed_compilation` as follows:

```

os.environ['NEURON_CC_FLAGS'] = os.environ.get('NEURON_CC_FLAGS', '') + ' --retry_failed_
↪ compilation'

```

This would retry the compilation and would replace a failed result in the cache with a successful compilation result.

### Compilation errors when placing NeuronCache home directory on NFS/EFS/FSx mounted drive

Currently, NeuronCache default root directory is `/var/tmp` which is local to the instance you are running on. You can modify the location of the NeuronCache root directory using `NEURON_CC_FLAGS='--cache_dir=<root_dir>'`. However, when the NeuronCache directory is placed in a directory that is part of a NFS mounted drive shared among multiple instances, you may encounter file errors such as file not found, file corruption, or `KeyError` when running multi-instance training:

```

KeyError: 'neff_cache2/neuron-compile-cache/USER_neuroncc-1.0.48875.0+7437fbf18/MODULE_
↪ 7223055628515330524/MODULE_0_SyncTensorsGraph.14_7223055628515330524_compute1-dy-
↪ training-2-1-e859998e-3035-5df63dab5ce63'

```

This is a result of limitations to file locking on NFS. EFS/FSx also exhibit similar limitation. The workaround is to setup separate NeuronCache root directories for each worker instance, such as `NEURON_CC_FLAGS="--cache_dir=$HOME/neuron_cache/bert/`hostname`"`, where the home directory is shared among worker instances as in `ParallelCluster`.

Consider the use case of a `ParallelCluster` with SLURM cluster management. The home directory of the head node is shared via NFS with worker instances. Also, SLURM would terminate the idle worker instances when the cluster is configured as dynamic auto-scaling cluster, and the default cache in the terminated worker instance's `/var/tmp` is deleted. So to persist the cache across runs separated by a cluster idle period, we use the workaround above to create separate NeuronCache root directories for each worker instance. For example, see [BERT ParallelCluster script](#).

## Compilation error: “Expect ap datatype to be of type float32 float16 bfloat16 uint8”

If an XLA example fails to run because of failed compilation and one of the error messages is “Expect ap datatype to be of type float32 float16 bfloat16 uint8”, then please set the environment variable XLA\_USE\_32BIT\_LONG=1 in your script:

```
os.environ['XLA_USE_32BIT_LONG'] = '1'
```

```
11/18/2021 04:51:25 PM WARNING 34567 [StaticProfiler]: matmul-based transposes inserted_
↳by penguin takes up 93.66 percent of all matmul computation
terminate called after throwing an instance of 'std::runtime_error'
  what():  === BIR verification failed ===
Reason: Expect ap datatype to be of type float32 float16 bfloat16 uint8
Instruction: I-545-0
Opcode: Matmult
Input index: 0
Argument AP:
Access Pattern: [[1,8],[1,1],[1,1]]
Offset: 0
Memory Location: {compare.85-t604_i0}@SB<0,0>(8x2)#Internal DebugInfo: <compare.
↳85|uint16|UNDEF|[8, 1, 1]>
```

## NeuronCore(s) not available - Requested:1 Available:0

When you see “NeuronCore(s) not available” please terminate processes that may be holding the NeuronCores and terminate any neuron-top sessions that are running. Also check if someone else is using the system. Then do “sudo rmmod neuron; sudo modprobe neuron” to reload the driver.

```
2021-Nov-15 15:21:28.0231 7245:7245 ERROR NRT:nrt_allocate_neuron_cores NeuronCore(s)_
↳not available - Requested:nc1-nc1 Available:0
2021-11-15 15:21:28.231864: F ./tensorflow/compiler/xla/service/neuron/neuron_runtime.h:
↳1037] Check failed: status == NRT_SUCCESS NEURONPOC : nrt_init failed. Status = 1
```

Often when you run multi-worker training, there can be many python processes leftover after a run is interrupted. To kill all python processes, run the follow (WARNING: this kills all python processes on the system) then reload the driver:

```
killall -9 python
killall -9 python3
sudo rmmod neuron; sudo modprobe neuron
```

## TDRV error “TDRV:exec\_consume\_infer\_status\_notification”

If you see TDRV error “TDRV:exec\_consume\_infer\_status\_notification”, try reloading the driver using sudo modprobe -r neuron; sudo modprobe neuron;.

```
2022-Mar-10 18:51:19.07392022-Mar-10 18:51:19.0739 17821:17931 ERROR TDRV:exec_consume_
↳infer_status_notifications 17822:18046 ERROR TDRV:exec_consume_infer_status_
↳notifications Unexpected number of CC notifications: mod->cc_op_count=1, cc_start_
↳cnt=0, cc_end_cnt=0Unexpected number of CC notifications: mod->cc_op_count=1, cc_
```

(continues on next page)

(continued from previous page)

```

→start_cnt=0, cc_end_cnt=0

2022-Mar-10 18:51:19.07392022-Mar-10 18:51:19.0739 17821:17931 ERROR   TDRV:exec_consume_
→infer_status_notifications 17822:18046 ERROR   TDRV:exec_consume_infer_status_
→notifications (NON-FATAL, Ignoring) inference timeout (180000 ms) on Neuron Device 0.
→NC 0, waiting for cc status notifications.

(NON-FATAL, Ignoring) inference timeout (180000 ms) on Neuron Device 0 NC 1, waiting for.
→cc status notifications.

```

### TDRV error “TDRV:tdrv\_one\_tmpbuf\_reserve Number of ONE TMPBUF pages requested exceeded the max number of pages allowed (requested: <N>, max allowed: 16).”

If you see the TDRV error “TDRV:tdrv\_one\_tmpbuf\_reserve Number of ONE TMPBUF pages requested exceeded the max number of pages allowed (requested: <N>, max allowed: 16)”, it maybe due to model tensors requiring more device memory then available. A solution is to try training with a smaller data batch size.

```

ERROR   TDRV:tdrv_one_tmpbuf_reserve           Number of ONE TMPBUF pages requested.
→exceeded the max number of pages allowed (requested: 28, max allowed: 16).
ERROR   TDRV:copy_and_stage_mr                 Failed to reserve one tmpbuf memory
ERROR   TDRV:kbl_model_add                     copy_and_stage_mr() error
W tensorflow/core/distributed_runtime/rpc/grpc_remote_master.cc:157] RPC failed with
→status = "UNAVAILABLE: Socket closed" and grpc_error_string = "{"created":"@1669183391.
→155135683","description":"Error received from peer ipv4:172.31.58.24:43941","file":
→"external/com_github_grpc_grpc/src/core/lib/surface/call.cc","file_line":1056,"grpc_
→message":"Socket closed","grpc_status":14}", maybe retrying the RPC

```

### Could not open the ndX, close device failed, TDRV not initialized

If you see error messages stating “Could not open the ndX” (where X is an integer from 0..15), please run `neuron-ls` and ensure that you are able to see all 16 Neuron devices in the output. If one or more devices are missing please report the issue to [aws-neuron-support@amazon.com](mailto:aws-neuron-support@amazon.com) with the instance ID and a screen capture of `neuron-ls` output.

```

2021-Nov-11 15:33:20.0161 7912:7912 ERROR   TDRV:tdrv_init_mla_phase1
→ Could not open the nd0
2021-Nov-11 15:33:20.0161 7912:7912 ERROR   TDRV:tdrv_destroy_one_mla
→ close device failed
2021-Nov-11 15:33:20.0161 7912:7912 ERROR   TDRV:tdrv_destroy
→ TDRV not initialized
2021-Nov-11 15:33:20.0161 7912:7912 ERROR   NRT:nrt_init
→ Failed to initialize devices, error:1
2021-11-11 15:33:20.161331: F ./tensorflow/compiler/xla/service/neuron/neuron_runtime.h:
→1033] Check failed: status == NRT_SUCCESS NEURONPOC : nrt_init failed. Status = 1

```

## Multiworker execution hangs during NCCL init

When your multi-worker execution hangs during NCCL init, you can try to reserve the port used by environment variable `NEURON_RT_ROOT_COMM_ID` by (here we use host:port localhost:48620 as an example but you can use any free port and root node's host IP):

```
sudo sysctl -w net.ipv4.ip_local_reserved_ports=48620
```

Then set the environment variable `NEURON_RT_ROOT_COMM_ID` in your script:

```
os.environ["NEURON_RT_ROOT_COMM_ID"] = "localhost:48620"
```

## NRT init error “One or more engines are running. Please restart device by reloading driver”

If you see an error stating “One or more engines are running. Please restart device by reloading driver” please follow the instruction and reload the driver using “`sudo modprobe -r neuron; sudo modprobe neuron;`”.

```
2021-Nov-15 20:23:27.0280 3793:3793 ERROR TDRV:tpb_eng_init_hals_v2 CRITICAL HW ERROR:
↳ One or more engines are running. Please restart device by reloading driver:
sudo modprobe -r neuron; sudo modprobe neuron;
2021-Nov-15 20:23:27.0280 3793:3793 ERROR TDRV:tdrv_init_one_mla_phase2 nd0 nc0 HAL init
↳ failed. error:1
```

## NRT error “ERROR TDRV:kbl\_model\_add Attempting to load an incompatible model!”

If you see an NRT error “ERROR TDRV:kbl\_model\_add Attempting to load an incompatible model!” this means that the compiler `neuronx-cc` used to compile the model is too old. See installation instruction to update to latest compiler.

## NRT error “ERROR HAL:aws\_hal\_sprot\_config\_remap\_entry SPROT remap destination address must be aligned size”

If you see an NRT error “ERROR HAL:aws\_hal\_sprot\_config\_remap\_entry SPROT remap destination address must be aligned size”, please check the kernel version and upgrade it to the distribution's latest kernel.

For example, on Ubuntu 18.04.6 LTS, the kernel version 4.15.0-66-generic is known to cause this error when running MLP tutorial. This is due to a known bug in the kernel in aligned memory allocation. To fix this issue, please upgrade your kernel to latest version (i.e. 4.15.0-171-generic):

```
uname -a
sudo apt-get update
sudo apt-get upgrade
sudo apt-get dist-upgrade
```

Please reboot after the upgrade. Use “`uname -a`” to check kernel version again after reboot.



**NCCL warning : “NCCL WARN Timeout waiting for RX (waited 120 sec) - retrying”**

When running multi-worker training, if a graph has collective communication operator like an `all_reduce`, it requires all the workers involved in the collective communication to load the graph in the runtime at approximately same time. If any of the worker doesn't load the graph within a 120 sec window from the first model load by any of the worker, you would see warnings like `NCCL WARN Timeout waiting for RX (waited 120 sec) - retrying`. When you see such warnings check for the following in the log messages:

1. One of the workers is compiling a graph: In multi-worker training, there is a chance that each worker builds a slightly different graph. This would result in cache miss and can result in compilation. Since the compilations during training run are serialized, the first worker can compile and load the graph with collective communication. It would then wait for 120 secs for other works to join. If they don't show up because they are compiling their own graphs, first worker would start throwing a warning message as above. The warning in this case is `non-fatal` and would go away once all workers have compiled their respective graphs and then loaded them. To identify this scenario, look for `No candidate found under . . . .` logs around the warning. You should also see `. . . .` which indicates compilation is in progress.
2. Server on one of the nodes crashed: In distributed training across multiple nodes, if the server on one node crashed, the workers on other nodes would keep waiting on model load and you would see above `timeout` logs on those nodes. To identify if the server crashed, check if you see the following error on any of the nodes:

```
`RPC failed with status = "UNAVAILABLE: Socket closed" and grpc_error_string = "{"created
↪": "@1664146011.016500243", "description": "Error received from peer ipv4:10.1.24.109:
↪37379", "file": "external/com_github_grpc_grpc/src/core/lib/surface/call.cc", "file_line":
↪1056, "grpc_message": "Socket closed", "grpc_status": 14}", maybe retrying the RPC`
```

If you see the above error, then it means there is a server crash and you need to cancel the training run.

**RPC error: “RPC failed with status = ‘UNAVAILABLE: Socket closed’”**

When you see the above error, it means that the xrt server crashed. When you see such an error, look for the following:

1. Check for any error logs before the `RPC error`. That should indicate the root cause of server crash. Note: The actual error log might be buried because of all the `RPC error` logs that swamp the logs.
2. Sometimes the server can crash because of host OOM. This can happen when we are loading and saving checkpoints. In such cases, you only see `RPC errors` and no other log. You can check if any instance is going out of memory by using tools like `dmesg`

**Error “Assertion ‘listp->slotinfo[cnt].gen <= GL(dl\_tls\_generation)’ failed” followed by ‘RPC failed with status = “UNAVAILABLE: Connection reset by peer””**

The error “Assertion ‘listp->slotinfo[cnt].gen <= GL(dl\_tls\_generation)’ failed” is intermittent and occurs when using glibc 2.26. To find out the glibc version you have, you can run `ldd --version`. The workaround is to use Ubuntu 20 where glibc is 2.27.

```
INFO: Inconsistency detected by ld.so: ../elf/dl-tls.c: 488: _dl_allocate_tls_init:
↪Assertion ‘listp->slotinfo[cnt].gen <= GL(dl_tls_generation)’ failed!
INFO: 2022-10-03 02:16:04.488054: W tensorflow/core/distributed_runtime/rpc/grpc_remote_
↪master.cc:157] RPC failed with status = "UNAVAILABLE: Connection reset by peer" and
↪grpc_error_string = "{"created": "@1664763364.487962663", "description": "Error received
↪from peer ipv4:10.0.9.150:41677", "file": "external/com_github_grpc_grpc/src/core/lib/
↪surface/call.cc", "file_line": 1056, "grpc_message": "Connection reset by peer", "grpc_
↪status": 14}", maybe retrying the RPC
```



## RPC connection error: “RPC failed with status = UNAVAILABLE: Connection reset by peer” not preceded by any error

This error may not be preceded by another error like shown in the previous section. In this case, the RPC connection error usually happens when we do distributed training across multiple nodes. When you see such error, please wait for a few minutes. It might be because some node is taking time to setup and hence the other node is not able to connect to it just yet. Once, all nodes are up, training should resume.

## Runtime errors “Missing infer\_status notification” followed by “inference timeout”

If you get a timeout error like below:

```
ERROR   TDRV:exec_consume_tpb_status_notifications   Missing infer_status notification:␣
→(end:4)
ERROR   TDRV:exec_consume_infer_status_notifications (FATAL-RT-UNDEFINED-STATE) inference␣
→timeout (6000000 ms) on Neuron Device 4 NC 1, waiting for execution completion␣
→notification
```

It maybe due to long graph execution time causing synchronization delays exceeding the default timeout. Please try increasing the timeout to larger value using NEURON\_RT\_EXEC\_TIMEOUT (unit in seconds) and see if the problem is resolved.

## Protobuf Error “TypeError: Descriptors cannot not be created directly.”

If you install torch-neuronx after neuronx-cc, you may get the Protobuf error “TypeError: Descriptors cannot not be created directly.”. To fix this, please reinstall neuronx-cc using “pip install --force-reinstall neuronx-cc”.

```
Traceback (most recent call last):
  File "./run_glue.py", line 570, in <module>
    main()
  File "./run_glue.py", line 478, in main
    data_collator=data_collator,
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→transformers/trainer.py", line 399, in __init__
    callbacks, self.model, self.tokenizer, self.optimizer, self.lr_scheduler
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→transformers/trainer_callback.py", line 292, in __init__
    self.add_callback(cb)
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→transformers/trainer_callback.py", line 309, in add_callback
    cb = callback() if isinstance(callback, type) else callback
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→transformers/integrations.py", line 390, in __init__
    from torch.utils.tensorboard import SummaryWriter # noqa: F401
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→torch/utils/tensorboard/__init__.py", line 10, in <module>
    from .writer import FileWriter, SummaryWriter # noqa: F401
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→torch/utils/tensorboard/writer.py", line 9, in <module>
    from tensorboard.compat.proto.event_pb2 import SessionLog
  File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
```

(continues on next page)

(continued from previous page)

```

→ tensorboard/compat/proto/event_pb2.py", line 17, in <module>
    from tensorboard.compat.proto import summary_pb2 as tensorboard_dot_compat_dot_proto_
→ dot_summary__pb2
    File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→ tensorboard/compat/proto/summary_pb2.py", line 17, in <module>
    from tensorboard.compat.proto import tensor_pb2 as tensorboard_dot_compat_dot_proto_
→ dot_tensor__pb2
    File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→ tensorboard/compat/proto/tensor_pb2.py", line 16, in <module>
    from tensorboard.compat.proto import resource_handle_pb2 as tensorboard_dot_compat_
→ dot_proto_dot_resource__handle__pb2
    File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→ tensorboard/compat/proto/resource_handle_pb2.py", line 16, in <module>
    from tensorboard.compat.proto import tensor_shape_pb2 as tensorboard_dot_compat_dot_
→ proto_dot_tensor__shape__pb2
    File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→ tensorboard/compat/proto/tensor_shape_pb2.py", line 42, in <module>
        serialized_options=None, file=DESCRIPTOR),
    File "/home/ec2-user/aws_neuron_venv_pytorch_p37_exp/lib64/python3.7/site-packages/
→ google/protobuf/descriptor.py", line 560, in __new__
        _message.Message._CheckCalledFromGeneratedFile()
TypeError: Descriptors cannot not be created directly.
If this call came from a _pb2.py file, your generated code is out of date and must be
→ regenerated with protoc >= 3.19.0.
If you cannot immediately regenerate your protos, some other possible workarounds are:
    1. Downgrade the protobuf package to 3.20.x or lower.
    2. Set PROTOCOL_BUFFERS_PYTHON_IMPLEMENTATION=python (but this will use pure-Python
→ parsing and will be much slower).

```

### TDRV error “Timestamp program stop timeout”

If you see TDRV error “Timestamp program stop timeout”, i.e. when rerunning a training script after it was interrupted, try first reloading the driver using `sudo modprobe -r neuron; sudo modprobe neuron`; (make sure neuron-top and/or neuron-monitor are not running).

```

2022-Aug-31 04:59:21.0546 117717:117717 ERROR   TDRV:tsync_wait_eng_stop
→ nd0 nc0 Timestamp program stop timeout (1000 ms)
2022-Aug-31 04:59:21.0546 117717:117717 ERROR   TDRV:tsync_wait_nc_stop
→ nd0 nc0 Error while waiting for timestamp program to end on TPB eng 0
2022-Aug-31 04:59:21.0546 117717:117717 ERROR   TDRV:tsync_timestamps_finish
→ nd0 nc0 Failed to stop neuron core
2022-Aug-31 04:59:21.0546 117717:117717 ERROR   TDRV:tdrv_tsync_timestamps
→ nd0 nc0 Failed to end timestamp sync programs
2022-Aug-31 04:59:22.0768 117717:117717 ERROR   TDRV:tdrv_destroy
→ TDRV not initialized
2022-Aug-31 04:59:22.0768 117717:117717 ERROR   NRT:nrt_init
→ Failed to initialize devices, error:5

```

## Compiler error “module ‘numpy’ has no attribute ‘asscalar’”

When you have a newer version of numpy in the Python environment, compilations may fail with the “error module ‘numpy’ has no attribute ‘asscalar’”. Please note the neuronx-cc has the following dependency on numpy “numpy<=1.20.0,>=1.13.3”. To workaround this error, please do “pip install –force-reinstall neuronx-cc” to reinstall neuronx-cc with the proper dependencies.

```
ERROR 227874 [neuronx-cc]:
↳ *****
ERROR 227874 [neuronx-cc]: An Internal Compiler Error has occurred
ERROR 227874 [neuronx-cc]:
↳ *****
ERROR 227874 [neuronx-cc]:
ERROR 227874 [neuronx-cc]: Error message: module 'numpy' has no attribute 'asscalar'
ERROR 227874 [neuronx-cc]:
ERROR 227874 [neuronx-cc]: Error class:      AttributeError
ERROR 227874 [neuronx-cc]: Error location: Unknown
ERROR 227874 [neuronx-cc]: Version information:
ERROR 227874 [neuronx-cc]:   NeuronX Compiler version 2.1.0.76+2909d26a2
ERROR 227874 [neuronx-cc]:
ERROR 227874 [neuronx-cc]:   HWM version 2.1.0.7-64eaede08
ERROR 227874 [neuronx-cc]:   NEFF version Dynamic
ERROR 227874 [neuronx-cc]:   TVM not available
ERROR 227874 [neuronx-cc]:   NumPy version 1.23.3
ERROR 227874 [neuronx-cc]:   MXNet not available
ERROR 227874 [neuronx-cc]:
```

## Import errors ‘generic\_type: type “IrValue” is already registered!’ or ‘generic\_type: type “XlaBuilder” is already registered!’

When you encounter a PyTorch import error ‘import \_XLAC ... generic\_type: type “IrValue” is already registered!’ or ‘import \_XLAC ... generic\_type: type “XlaBuilder” is already registered!’, please check that TensorFlow and/or JAX are not installed in the Python environment. If they are installed, please uninstall them.

## Import error “import \_XLAC ImportError: <>/site-packages/\_XLAC.cpython-38-x86\_64-linux-gnu.so: undefined symbol”

When you encounter a PyTorch import error “import \_XLAC ImportError: <>/site-packages/\_XLAC.cpython-38-x86\_64-linux-gnu.so: undefined symbol” during execution, please check:

1. TensorFlow and/or JAX are not installed in the Python environment. If they are installed, please uninstall them.
2. The installed PyTorch (torch) package major/minor versions match the installed torch-neuronx package’s major/minor versions (ie. 1.11). If they don’t match, please install the version of PyTorch that matches torch-neuronx.

```
Traceback (most recent call last):
  File "/opt/ml/mlp_train.py", line 11, in <module>
    import torch_xla.core.xla_model as xm
  File "/usr/local/lib/python3.8/site-packages/torch_xla/__init__.py", line 117, in
↳ <module>
```

(continues on next page)

(continued from previous page)

```
import _XLAC
ImportError: /usr/local/lib/python3.8/site-packages/_XLAC.cpython-38-x86_64-linux-gnu.so:
↳ undefined symbol: _ZNK3c1010TensorImpl7stridesEv
```

### NaNs seen with transformers version >= 4.21.0 when running HF BERT fine-tuning or pretraining with XLA\_USE\_BF16=1 or XLA\_DOWNCAST\_BF16=1

When running HuggingFace BERT (any size) fine-tuning tutorial or pretraining tutorial with transformers version >= 4.21.0 and using XLA\_USE\_BF16=1 or XLA\_DOWNCAST\_BF16=1, you will see NaNs in the loss immediately at the first step. More details on the issue can be found at [pytorch/xla#4152](https://pytorch/xla#4152). The workaround is to use 4.20.0 or earlier (the tutorials currently recommend version 4.15.0) or add `transformers.modeling_utils.get_parameter_dtype = lambda x: torch.bfloat16` to the Python script.

### Network Connectivity Issue on trn1/trn1n 32xlarge with Ubuntu

#### Description

Ubuntu distributions have network connectivity issues when multiple interfaces are connected to the same subnet. trn1/trn1n 32xlarge comes with 8/16 network interfaces. (To launch trn1/trn1n with 8/16 interfaces please follow [here](#))

AWS publishes a package that installs a helper service to address the issue. This service runs at the startup, creates the appropriate netplan files, updates the netplan and the the instance networking and terminates.

Note that the following fix is only required on instances launched using generic Ubuntu AMIs. Neuron AMIs and instances launched via ParallelCluster do not require the fix.

#### Patch to fix networking on a multi-interface instance

```
wget -O /tmp/aws-ubuntu-eni-helper.deb 'https://github.com/aws-samples/aws-efa-nccl-
↳ baseami-pipeline/blob/master/nvidia-efa-ami_base/networking/aws-ubuntu-eni-helper_0.3-
↳ 1_all.deb?raw=true'
sudo apt install /tmp/aws-ubuntu-eni-helper.deb -y
sudo systemctl enable aws-ubuntu-eni-helper.service
sudo systemctl start aws-ubuntu-eni-helper.service
```

#### How to apply the patch?

The following steps could be followed to resolve this issue:

- Launch trn1.32xl from AWS console (starts with `single interface`, does not suffer from the multi-interface issue)
- Apply the patch on this newly launched single-interface instance
- Create a new AMI from this instance
- Launch an 8 or 16 interface instance using that AMI.

**Note:** The patch installs and enables the service but does not run it. This is intentional. The service will run at the startup when the AMI is used to launch a multi-interface instance.

#### FAQs

**Note:** Neuron DLAMI has the patch installed, users are always encouraged to launch the instances using the DLAMI which does not require any fix. Please refer to the [Set Up Guide](#) to know how to launch an instance using DLAMI.

### “Too many open files” when running training job

When running a large model training with several workers, it can result in errors like the following.

```
2023-Jun-14 19:05:29.0312 4112959:4113326 [23] bootstrap.cc:106 CCOM WARN Call to accept_
↳ failed : Too many open files
2023-Jun-14 19:05:29.0312 4112959:4113263 [14] include/socket.h:438 CCOM WARN Net :_
↳ Socket creation failed : Too many open files
2023-Jun-14 19:05:29.0312 4112959:4113326 ERROR ENC:ncclBootstrapRecv _
↳ failed neuronBootstrapRecv request to NCCL
2023-Jun-14 19:05:29.0312 4112959:4113249 [12] bootstrap.cc:106 CCOM WARN Call to accept_
↳ failed : Too many open files
2023-Jun-14 19:05:29.0312 4112959:4113263 ERROR ENC:ncclBootstrapSend _
↳ failed neuronBootstrapSend request to NCCL
2023-Jun-14 19:05:29.0312 4112959:4113270 [15] bootstrap.cc:106 CCOM WARN Call to accept failed : Too_
↳ many open files
```

This can result when the default OS limits is low. The hard and soft limits can be set on OS using the following commands or by manually opening and setting the limits.

```
sudo sed -i 'H;1h;$!d;x;/hard *nofile/!s/$/\n* hard nofile 65536/' /etc/security/limits.
↳ conf
sudo sed -i 'H;1h;$!d;x;/soft *nofile/!s/$/\n* soft nofile 65536/' /etc/security/limits.
↳ conf
sudo sed -i 's/^#*\(\*\|s*\*\)\s*soft\s*nofile\s*[0-9]\+$/\1 soft nofile 65536/' /etc/
↳ security/limits.conf
sudo sed -i 's/^#*\(\*\|s*\*\)\s*hard\s*nofile\s*[0-9]\+$/\1 hard nofile 65536/' /etc/
↳ security/limits.conf
sudo sed -i 's/^#*\(\*\|s*\*\)\s*soft\s*nofile\s*[0-9]\+$/\1 soft nofile 65536/' /etc/
↳ security/limits.d/01_efa.conf || true
sudo sed -i 's/^#*\(\*\|s*\*\)\s*hard\s*nofile\s*[0-9]\+$/\1 hard nofile 65536/' /etc/
↳ security/limits.d/01_efa.conf || true
```

The `01_efa.conf` file is created as part of the EFA installation and needs to be updated. If EFA driver is not installed the file `01_efa.conf` doesn't exist and the sed commands will fail with *No such file or directory*. If there are other files under `limits.d` with file limits they need to be updated as well.

### “undefined symbol”

To maintain compatibility with the packages vended publicly in Pypi, AWS Neuron python packages contain binary extensions that are compiled with the pre-2011 libstdc++ application binary interface (ABI). If a custom version of a package - such as *torch* - is compiled using a modern compiler, it can result in “undefined symbol” errors due to mismatches between the package and AWS Neuron package.

To support this situation, we provide alternative versions of AWS Neuron packages that are compiled according to the newer 2011 ABI. For information on how to use these packages, see `pytorch-install-cxx11`.

*This document is relevant for: Trn1, Trn2*

- [PyTorch Neuron \(torch-neuronx\) - Supported Operators](#)
- [How to prepare trn1.32xlarge for multi-node execution](#)
- [PyTorch Neuron \(torch-neuronx\) for Training Troubleshooting Guide](#)
- [PyTorch Neuron \(torch-neuronx\) release notes](#)

This document is relevant for: **Inf2, Trn1, Trn2**

Setup (**torch-neuronx**)

### Tutorials

- [Hugging Face BERT Pretraining Tutorial \(Data-Parallel\)](#)
- [Multi-Layer Perceptron Training Tutorial](#)
- [PyTorch Neuron for Trainium Hugging Face BERT MRPC task finetuning using Hugging Face Trainer API](#)
- [Fine-tune T5 model on Trn1](#)
- [ZeRO-1 Tutorial](#)
- [Analyze for Training Tutorial](#)
- [Neuron Custom C++ Operators in MLP Training](#)
- [Neuron Custom C++ Operators Performance Optimization](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

### Additional Examples

- [AWS Neuron Reference for Nemo Megatron GitHub Repository](#)
- [AWS Neuron Samples for EKS](#)
- [AWS Neuron Samples for AWS ParallelCluster](#)
- [AWS Neuron Samples GitHub Repository](#)

### API Reference Guide

- [PyTorch NeuronX neuron\\_parallel\\_compile CLI](#)
- [Neuron Persistent Cache](#)
- [PyTorch NeuronX Environment Variables](#)
- [PyTorch NeuronX Profiling API](#)

## Developer Guide

- [Developer Guide for Training with PyTorch NeuronX](#)
- [How to debug models in PyTorch NeuronX](#)
- [Developer Guide for Profiling with PyTorch NeuronX](#)

## Misc

- [PyTorch Neuron \(torch-neuronx\) - Supported Operators](#)
- [How to prepare trn1.32xlarge for multi-node execution](#)
- [PyTorch Neuron \(torch-neuronx\) for Training Troubleshooting Guide](#)
- [PyTorch Neuron \(torch-neuronx\) release notes](#)

*This document is relevant for:* Trn1, Trn2

## 2.1.5 PyTorch NeuronX

PyTorch NeuronX for training on Trn1 and Trn2

Pytorch NeuronX for inference on Inf2, Trn1, and Trn2

## 2.1.6 PyTorch Neuron

PyTorch Neuron for inference on Inf1

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## 2.2 JAX Neuron (beta)

JAX Neuron is a software package containing Neuron-specific JAX features, such as the Neuron NKI JAX interface. It also serves as a meta-package for providing a tested combination of the `jax-neuronx`, `jax`, `jaxlib`, `libneuronxla`, and `neuronx-cc` packages.

*This document is relevant for:* Inf2, Trn1, Trn2

### 2.2.1 JAX Neuron plugin Setup

The JAX Neuron plugin is a set of modularized JAX plugin packages integrating AWS Trainium and Inferentia machine learning accelerators into JAX as pluggable devices. It includes the following Python packages, all hosted on the AWS Neuron pip repository.

- `libneuronxla`: A package containing Neuron's integration into JAX's runtime `PJRT`, built using the `PJRT C-API plugin` mechanism. Installing this package enables using Trainium and Inferentia natively as JAX devices.
- `jax-neuronx`: A package containing Neuron-specific JAX features, such as the `Neuron NKI JAX` interface. It also serves as a meta-package for providing a tested combination of the `jax-neuronx`, `jax`, `jaxlib`, `libneuronxla`, and `neuronx-cc` packages. Making proper use of the features provided in `jax-neuronx` will unleash the full potential of Trainium and Inferentia.

**Note:** If you are facing a connectivity issue during the model loading process on a Trn1 instance with Ubuntu, that could probably be because of Ubuntu limitations with multiple interfaces. To solve this problem, please follow the steps mentioned [here](#).

Users are highly encouraged to use DLAMI to launch the instances, since DLAMIs come with the required fix.

---

### Launch the Instance

- To launch an instance, follow the instructions at [launch an Amazon EC2 Instance](#). Make sure to select the correct instance type on the EC2 console.
- For more information about instance sizes and pricing, see [Amazon EC2 Trn1 Instances](#) and [Amazon EC2 Inf2 Instances](#)
- Select Ubuntu Server 22 AMI.
- When launching a Trn1, adjust your primary EBS volume size to a minimum of 512GB.
- After launching the instance, follow the instructions in [Connect to your instance](#) to connect to the instance.

### Install Drivers and Tools

Ubuntu

```
# Configure Linux for Neuron repository updates
. /etc/os-release
sudo tee /etc/apt/sources.list.d/neuron.list > /dev/null <<EOF
deb https://apt.repos.neuron.amazonaws.com ${VERSION_CODENAME} main
EOF
wget -qO - https://apt.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-NEURON.PUB | \
  sudo apt-key add -

# Update OS packages
sudo apt-get update -y

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install git
sudo apt-get install git -y

# install Neuron Driver
sudo apt-get install aws-neuronx-dkms=2.* -y

# Install Neuron Runtime
sudo apt-get install aws-neuronx-collectives=2.* -y
sudo apt-get install aws-neuronx-runtime-lib=2.* -y

# Install Neuron Tools
sudo apt-get install aws-neuronx-tools=2.* -y

# Add PATH
export PATH=/opt/aws/neuron/bin:$PATH
```



Amazon Linux 2023

```
# Configure Linux for Neuron repository updates
sudo tee /etc/yum.repos.d/neuron.repo > /dev/null <<EOF
[neuron]
name=Neuron YUM Repository
baseurl=https://yum.repos.neuron.amazonaws.com
enabled=1
metadata_expire=0
EOF
sudo rpm --import https://yum.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-NEURON.
↪PUB

# Update OS packages
sudo yum update -y

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install git
sudo yum install git -y

# install Neuron Driver
sudo yum install aws-neuronx-dkms-2.* -y

# Install Neuron Runtime
sudo yum install aws-neuronx-collectives-2.* -y
sudo yum install aws-neuronx-runtime-lib-2.* -y

# Install Neuron Tools
sudo yum install aws-neuronx-tools-2.* -y

# Add PATH
export PATH=/opt/aws/neuron/bin:$PATH
```

## Install the JAX Neuron Plugin

We provide two methods for installing the JAX Neuron plugin. The first is to install the `jax-neuronx` meta-package from the AWS Neuron pip repository. This method provides a production-ready JAX environment where `jax-neuronx`'s major dependencies, namely `jax`, `jaxlib`, `libneuronxla`, and `neuronx-cc`, have undergone thorough testing by the AWS Neuron team and will have their versions pinned during installation.

```
python3 -m pip install jax-neuronx[stable] --extra-index-url=https://pip.repos.neuron.
↪amazonaws.com
```

The second is to install packages `jax`, `jaxlib`, `libneuronxla`, and `neuronx-cc` separately, with `jax-neuronx` being an optional addition. Because `libneuronxla` supports a broad range of `jaxlib` versions through the PJRT C-API mechanism, this method provides flexibility when choosing `jax` and `jaxlib` versions, enabling JAX users to bring the JAX Neuron plugin into their own JAX environments.

```
python3 -m pip install jax==0.4.31 jaxlib==0.4.31 jax-neuronx libneuronxla neuronx-cc==2.
↪* --extra-index-url=https://pip.repos.neuron.amazonaws.com
```

We can now run some simple JAX programs on the Trainium or Inferentia accelerators.

```
~$ python3 -c 'import jax; print(jax.numpy.multiply(1, 1))'
Platform 'neuron' is experimental and not all JAX functionality may be correctly
↳ supported!
.
Compiler status PASS
1
```

Compatibility between packages `jaxlib` and `libneuronxla` can be determined from [PJRT C-API version](#). For more information, see [PJRT integration guide](#).

To determine compatible JAX versions, you can use the `libneuronxla.supported_clients` API for querying known supported client packages and their versions.

Help on function `supported_clients` in module `libneuronxla.version`:

```
supported_clients()
    Return a description of supported client (jaxlib, torch-xla, etc.) versions,
    as a list of strings formatted as ``<package> <version> (PJRT C-API <c-api version>)
↳ ``.
    For example,
    >>> import libneuronxla
    >>> libneuronxla.supported_clients()
    ['jaxlib 0.4.31 (PJRT C-API 0.54)', 'torch_xla 2.2.0 (PJRT C-API 0.35)', 'torch_xla_
↳ 2.3.0 (PJRT C-API 0.46)']
```

Note that the list of supported client packages and versions covers known versions only and may be incomplete. More versions could be supported, including Google's future `jaxlib` releases, assuming the PJRT C-API stays compatible with the current release of `libneuronxla`. As a result, we avoid specifying any dependency relationship between `libneuronxla` and `jaxlib`. This provides more freedom when coordinating `jax` and `libneuronxla` installations.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## 2.2.2 JAX NeuronX Known Issues

- Threefry RNG algorithm is not completely supported. Use `rbg` algorithm instead. This can be configured by setting the following config option `jax.config.update("jax_default_prng_impl", "rbg")`
- For JAX versions older than 0.4.34, caching does not work out of the box. Use the following to enable caching support,

```
import jax
import jax_neuronx
from jax._src import compilation_cache

compilation_cache.set_cache_dir('./cache_directory')
```

- For JAX versions older than 0.4.34, Buffer donation does not work out of the box. Add the following snippet to your script to enable it - `jax._src.interpreters.mlir._platforms_with_donation.append('neuron')`
- `jax.random.randint` does not produce expected distribution of `randint` values. Run it on CPU instead.
- Dynamic loops are not supported for `jax.lax.while_loop`. Only static while loops are supported.

- `jax.lax.cond` is not supported.
- Host callbacks are not supported. As a result APIs based on callbacks from `jax.debug` and `jax.experimental.checkify` are not supported.
- Mesh configurations which use non-connected Neuron cores might crash during execution. You might observe compilation or Neuron runtime errors for such configurations. Device connectivity can be determined by using `neuron-ls --topology`.
- Not all dtypes supported by JAX work on Neuron. Check [Data Types](#) for supported data types.
- `jax.dlpack` is not supported.
- `jax.experimental.sparse` is not supported.
- `jax.lax.sort` only supports comparators with LE, GE, LT and GT operations.
- `jax.lax.reduce_precision` is not supported.
- Certain operations (for example, rng weight initialization) might result in slow compilations. Try to run such operations on the CPU backend or by setting the following environment variable `NEURON_RUN_TRIVIAL_COMPUTATION_ON_CPU=1`.
- Neuron only supports `float8_e4m3` and `float8_e5m2` for FP8 dtypes.
- Complex dtypes (`jnp.complex64` and `jnp.complex128`) are not supported.
- Variadic reductions are not supported.
- Out of bound access for scatter/gather operations can result in runtime errors.
- Dot operations on int dtypes are not supported.
- `lax.DotAlgorithmPreset` is not always respected. Dot operations occur in operand dtypes. This is a configurable parameter for `jax.lax.dot` and `jax.lax.dot_general`.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 2.2.3 API Reference Guide for JAX Neuronx

*This document is relevant for: Inf2, Trn1, Trn2*

### JAX NeuronX Environment Variables

Environment variables allow modifications to JAX NeuronX behavior without requiring code change to user script. It is recommended to set them in code or just before invoking the python process, such as `NEURON_RT_VISIBLE_CORES=8 python3 <script>` to avoid inadvertently changing behavior for other scripts. Environment variables specific to JAX Neuronx are:

#### NEURON\_CC\_FLAGS

- Compiler options. Full compiler options are described in the [mixed-precision-casting-options](#).

#### XLA\_FLAGS

- When set to `"--xla_dump_hlo_snapshots --xla_dump_to=<dir>"`, this environmental variable enables dumping snapshots in `<dir>` directory. See [Snapshotting With Torch-Neuronx 2.1](#) section for more information. The snapshotting interface for JAX and Pytorch are identical.

- When set to "--xla\_dump\_hlo\_as\_text --xla\_dump\_hlo\_as\_proto --xla\_dump\_to=<dir> --xla\_dump\_hlo\_pass\_re='.\*'", this environmental variable enables dumping HLOs in proto and text formats after each XLA pass. The dumped \*.hlo.pb files are in HloProto format.

#### NEURON\_FORCE\_PJRT\_PLUGIN\_REGISTRATION

- When NEURON\_FORCE\_PJRT\_PLUGIN\_REGISTRATION=1, the Neuron PJRT plugin will be registered in JAX regardless of the instance type.

#### NEURON\_RUN\_TRIVIAL\_COMPUTATION\_ON\_CPU

- When NEURON\_RUN\_TRIVIAL\_COMPUTATION\_ON\_CPU=1, the Neuron PJRT plugin will compile and execute “trivial” computations on CPU instead of Neuron cores. A “trivial” computation is defined as an HLO program that does not contain any collective-compute instructions. The HLO program will be compiled by the XLA CPU compiler and outputs of the computation will be allocated on Neuron cores. The following HLO instructions are considered as collective-compute instructions.

- all-gather
- all-gather-done
- all-gather-start
- all-reduce-done
- all-reduce-start
- all-to-all
- collective-permute
- partition-id
- replica-id
- recv
- recv-done
- reduce-scatter
- send
- send-done

#### NEURON\_PJRT\_PROCESSES\_NUM\_DEVICES

- Should be set to a comma-separated list stating the number of NeuronCores used by each worker process. It is used to construct a global device array with its size equal to the sum of the list. This gets reported to the XLA PJRT runtime when requested. Must be set for multi-process executions. It can be used in conjunction with NEURON\_RT\_VISIBLE\_CORES to expose a limited number of NeuronCores to each worker process. If NEURON\_RT\_VISIBLE\_CORES is not set, it should be set to available number of NeuronCores on the host. NEURON\_PJRT\_PROCESSES\_NUM\_DEVICES must be less than or equal to NEURON\_RT\_VISIBLE\_CORES.

#### NEURON\_PJRT\_PROCESS\_INDEX

- An integer stating the index (or rank) of the current worker process. This is required for multi-process environments where all workers need to know information on all participating processes. Must be set for multi-process executions. The value should be between 0 and NEURON\_PJRT\_PROCESS\_INDEX - 1.

#### NEURON\_RT\_STOCHASTIC\_ROUNDING\_EN [Neuron Runtime]

- When NEURON\_RT\_STOCHASTIC\_ROUNDING\_EN=1, JAX Neuron will use stochastic rounding instead of round-nearest-even for all internal rounding operations when casting from FP32 to a reduced precision data type (FP16, BF16, FP8, TF32). This feature has been shown to improve training convergence for reduced precision training jobs. To switch to round-nearest-even mode, set NEURON\_RT\_STOCHASTIC\_ROUNDING\_EN=0.

**NEURON\_RT\_STOCHASTIC\_ROUNDING\_SEED [Neuron Runtime]**

- Sets the seed for the random number generator used in stochastic rounding (see previous section). If this environment variable is not set, the seed is set to 0 by default. Please set NEURON\_RT\_STOCHASTIC\_ROUNDING\_SEED to a fixed value to ensure reproducibility between runs.

**NEURON\_RT\_VISIBLE\_CORES [Neuron Runtime]**

- Integer range of specific NeuronCores needed by the process (for example, 0-3 specifies NeuronCores 0, 1, 2, and 3). Use this environment variable when launching processes to limit the launched process to specific consecutive NeuronCores.

Additional Neuron runtime environment variables are described in [NeuronX Runtime Configuration](#).

*This document is relevant for:* Inf2, Trn1, Trn2

- [JAX NeuronX Environment Variables](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 2.2.4 JAX NeuronX (jax-neuronx) release notes

### Table of Contents

- [Release \[0.6.0.1.0.\\*\]](#)
- [Release \[0.5.3.1.0.\\*\]](#)
- [Release \[0.1.3\]](#)
- [Release \[0.1.2\]](#)
- [Release \[0.1.1\]](#)

JAX NeuronX is a software package containing Neuron-specific JAX features, such as the Neuron NKI JAX interface. It also serves as a meta-package for providing a tested combination of the `jax-neuronx`, `jax`, `jaxlib`, `libneuronxla`, and `neuronx-cc` packages.

### Release [0.6.0.1.0.\*]

Date: 06/20/2025

### Summary

- This release supports JAX versions up to `0.6.0`.
- Known issues are listed within [JAX NeuronX Known Issues](#).

### Release [0.5.3.1.0.\*]

Date: 05/20/2025

#### Summary

- This release supports JAX versions up to 0.5.3.
- `jax_neuronx.nki_call` is no longer supported. Use `neuronxcc.nki.jit` instead.
- Known issues are listed within *JAX NeuronX Known Issues*.

### Release [0.1.3]

Date: 04/03/2025

#### Summary

- This release supports JAX versions up to 0.5.0.
- Known issues are listed within *JAX NeuronX Known Issues*.

### Release [0.1.2]

Date: 12/20/2024

#### Summary

This release supports JAX versions up to 0.4.35.

#### What's new in this release

- Support for JAX versions up to 0.4.35.
- Support for JAX caching API for versions 0.4.30+.

### Release [0.1.1]

Date: 09/16/2024

## Summary

This is the initial beta release of JAX NeuronX that contains Neuron-specific JAX features, such as the Neuron NKI JAX interface

## What's new in this release

Announcing the first JAX NeuronX release

- JAX interface for Neuron NKI

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

- [JAX Neuron plugin Setup](#)
- [JAX NeuronX Known Issues](#)
- [API Reference Guide for JAX Neuronx](#)
- [JAX NeuronX \(jax-neuronx\) release notes](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1

## 2.3 TensorFlow Neuron

TensorFlow Neuron unlocks high-performance and cost-effective deep learning acceleration on AWS Trainium-based and Inferentia-based Amazon EC2 instances.

TensorFlow Neuron enables native TensorFlow models to be accelerated on Neuron devices, so you can use your existing framework application and get started easily with minimal code changes.

*This document is relevant for:* Inf1, Inf2, Trn1

### 2.3.1 Tensorflow Neuron Setup

Tensorflow Neuron (tensorflow-neuronx) Setup for Inf2, Trn1/Trn1n Instances      Tensorflow Neuron (tensorflow-neuron) Setup for Inf1 Instances      *This document is relevant for:* Inf1, Inf2, Trn1

*This document is relevant for:* Inf2, Trn1

### 2.3.2 Inference on Inf2 & Trn1/Trn1n (tensorflow-neuronx)

*This document is relevant for:* Inf2, Trn1

## Tutorials (tensorflow-neuronx)

### Running Huggingface Roberta-Base with TensorFlow-NeuronX

This tutorial demonstrates how to compile the Huggingface roberta-base model and infer on a trn1.2xlarge instance with tensorflow-neuronx. To compile larger models like roberta-large, please consider using an inf2 instance.

#### Setup

To run this tutorial please follow the instructions for [TensorFlow-NeuronX Setup](#) and the [Jupyter Notebook Quickstart](#) and set your kernel to “Python (tensorflow-neuronx)”.

Next, install some additional dependencies.

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
      ↪detect
      !pip install transformers
```

### Download From Huggingface and Compile for AWS-Neuron

```
[ ]: import tensorflow as tf
import tensorflow_neuronx as tfnx
from transformers import RobertaTokenizer, TFRobertaModel
from transformers import BertTokenizer, TFBertModel

# Create a wrapper for the roberta model that will accept inputs as a list
# instead of a dictionary. This will allow the compiled model to be saved
# to disk with the model.save() fucntion.
class RobertaWrapper(tf.keras.Model):
    def __init__(self, model):
        super().__init__()
        self.model = model
    def __call__(self, example_inputs):
        return self.model({'input_ids' : example_inputs[0], 'attention_mask' : example_
        ↪inputs[1]})

tokenizer = RobertaTokenizer.from_pretrained('roberta-base')
model = RobertaWrapper(TFRobertaModel.from_pretrained('roberta-base'))

batch_size = 16

# create example inputs with a batch size of 16
text = ["Paris is the <mask> of France."] * batch_size
encoded_input = tokenizer(text, return_tensors='tf', padding='max_length', max_length=64)

# turn inputs into a list
example_input = [encoded_input['input_ids'], encoded_input['attention_mask']]

#compile
model_neuron = tfnx.trace(model, example_input)
```

(continues on next page)



(continued from previous page)

```
print("Running on neuron:", model_neuron(example_input))

# save the model to disk to save recompilation time for next usage
model_neuron.save('./roberta-neuron-b16')
```

### Run Basic Inference Benchmarking

```
[ ]: import numpy as np
import concurrent.futures
import time

reloaded_neuron_model = tf.keras.models.load_model('./roberta-neuron-b16')
print("Reloaded model running on neuron:", reloaded_neuron_model(example_input))

num_threads = 4
num_inferences = 1000

latency_list = []
def inference_with_latency_calculation(example_input):
    global latency_list
    start = time.time()
    result = reloaded_neuron_model(example_input)
    end = time.time()
    latency_list.append((end-start) * 1000)
    return result

start = time.time()
with concurrent.futures.ThreadPoolExecutor(max_workers=num_threads) as executor:
    futures = []
    for i in range(num_inferences):
        futures.append(executor.submit(inference_with_latency_calculation, example_
→input))
    for future in concurrent.futures.as_completed(futures):
        get_result = future.result()
end = time.time()

total_time = end - start

print(f"Throughput was {(num_inferences * batch_size)/total_time} samples per second.")
print(f"Latency p50 was {np.percentile(latency_list, 50)} ms")
print(f"Latency p90 was {np.percentile(latency_list, 90)} ms")
print(f"Latency p99 was {np.percentile(latency_list, 99)} ms")
```

*This document is relevant for: Inf2, Trn1*

## Using NEURON\_RT\_VISIBLE\_CORES with TensorFlow Serving

TensorFlow serving allows customers to scale-up inference workloads across a network. TensorFlow Neuron Serving uses the same API as normal TensorFlow Serving with two differences: (a) the saved model must be compiled for neuron and (b) the entry point is a different binary named `tensorflow_model_server_neuronx`. Follow the steps below to install the package using apt-get or yum. This will be pre-installed in a future release.

### Install TensorFlow Model Server and Serving API

Follow the steps in the `install-neuronx-tensorflow`.

Then ensure you install using either apt-get or yum.

```
sudo apt-get install tensorflow-model-server-neuronx
```

or

```
sudo yum install tensorflow-model-server-neuronx
```

Also, you would need TensorFlow Serving API (use `--no-deps` to prevent installation of regular tensorflow).

```
pip install --no-deps tensorflow_serving_api
```

For the example image preprocessing using Keras preprocessing, the Python Imaging Library Pillow is required:

```
pip install pillow
```

To workaround h5py issue <https://github.com/aws/aws-neuron-sdk/issues/220>:

```
pip install "h5py<3.0.0"
```

### Export and Compile Saved Model

The following example shows graph construction followed by the addition of Neuron compilation step before exporting to saved model.

```
import tensorflow as tf
import tensorflow_neuronx as tfnx
import numpy as np

tf.keras.backend.set_learning_phase(0)
tf.keras.backend.set_image_data_format('channels_last')
image_sizes = [224, 224]
model = tf.keras.applications.ResNet50(weights='imagenet')
example_inputs = tf.random.uniform([1, *image_sizes, 3], dtype=tf.float32)

model_neuron = tfnx.trace(model, example_inputs)
# run the model once to define the forward pass and allow for saving
model_neuron(example_inputs)
tf.keras.models.save_model(model_neuron, './resnet50_neuron/1')
```

## Serving Saved Model

User can now serve the saved model with the `tensorflow_model_server_neuron` binary. To utilize multiple Neuron-Cores, it is recommended to launch multiple tensorflow model servers that listen to the same gRPC port:

```
export NEURON_RT_VISIBLE_CORES=0 # important to set this environment variable before_
↪launching model servers
tensorflow_model_server_neuron --model_name=resnet50_neuron \
    --model_base_path=$(pwd)/resnet50_neuron/ --port=8500

# then to run another server on a different neuron core open another
# window and run this, except this time set NEURON_RT_VISIBLE_CORES=1
# you can keep doing this up to the number of Neuron Cores on your machine

export NEURON_RT_VISIBLE_CORES=1
tensorflow_model_server_neuron --model_name=resnet50_neuron \
    --model_base_path=$(pwd)/resnet50_neuron/ --port=8500
```

The compiled model is staged in neuron DRAM by the server to prepare for inference.

## Generate inference requests to the model server

Now run inferences via GRPC as shown in the following sample client code:

```
import numpy as np
import grpc
import tensorflow as tf
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications.resnet50 import preprocess_input
from tensorflow_serving.apis import predict_pb2
from tensorflow_serving.apis import prediction_service_pb2_grpc
from tensorflow.keras.applications.resnet50 import decode_predictions

tf.keras.backend.set_image_data_format('channels_last')

if __name__ == '__main__':
    channel = grpc.insecure_channel('localhost:8500')
    stub = prediction_service_pb2_grpc.PredictionServiceStub(channel)
    img_file = tf.keras.utils.get_file(
        "./kitten_small.jpg",
        "https://raw.githubusercontent.com/awsmlabs/mxnet-model-server/master/docs/images/
↪kitten_small.jpg")
    img = image.load_img(img_file, target_size=(224, 224))
    img_array = preprocess_input(image.img_to_array(img)[None, ...])
    request = predict_pb2.PredictRequest()
    request.model_spec.name = 'resnet50_neuron'
    request.inputs['input_1'].CopyFrom(
        tf.make_tensor_proto(img_array, shape=img_array.shape))
    result = stub.Predict(request)
    prediction = tf.make_ndarray(result.outputs['output_1'])
    print(decode_predictions(prediction))
```

*This document is relevant for: Inf2, Trn1*

- [HuggingFace Roberta-Base \[html\]](#) [\[notebook\]](#)
- [Using NEURON\\_RT\\_VISIBLE\\_CORES with TensorFlow Serving](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

*This document is relevant for: Inf2, Trn1*

*This document is relevant for: Inf2, Trn1*

### API Reference Guide (tensorflow-neuronx)

*This document is relevant for: Inf2, Trn1*

### TensorFlow 2.x (tensorflow-neuronx) Tracing API

The Neuron tracing API enables tracing TensorFlow 2.x models for deployment on trn1 and inf2 AWS machine learning accelerators.

#### Method

`tensorflow_neuronx.trace`

#### Description

Trace a `keras.Model` or a Python callable that can be decorated by `tf.function`, and return an AWS-Neuron-optimized `keras.Model` that can execute on trn1 and inf2 AWS machine learning accelerators. Tracing is ideal for `keras.Model` that accepts a list of `tf.Tensor` objects and returns a list of `tf.Tensor` objects. It is expected that users will provide example inputs, and the `trace` function will execute `func` symbolically and convert it to a `keras.Model`.

The returned `keras.Model` will support inference only. Attributes or variables held by the original function or `keras.Model` will be dropped.

The returned `keras.Model` can be exported as `SavedModel` and served using TensorFlow Serving. Please see [tensorflow-serving](#) for more information about exporting to saved model and serving using TensorFlow Serving.

The returned `keras.Model` has an `.on_neuron_ratio` attribute which shows the percentage of ops mapped to neuron hardware. This calculation ignores `PlaceholderOp`, `IdentityOp`, `ReadVariableOp` and `NoOp`.

Options can be passed to Neuron compiler via the environment variable `NEURON_CC_FLAGS`. For example, the syntax `env NEURON_CC_FLAGS="--workdir ./artifacts"` directs the Neuron compiler to dump artifacts in the `artifacts` directory for debugging. See [Neuron Compiler CLI Reference Guide \(neuronx-cc\)](#) for more information about compiler options.

## Arguments

- **func:** The `keras.Model` or function to be traced.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.

- **subgraph\_builder\_function:** (Optional) A callable with signature

```
subgraph_builder_function(node : NodeDef) -> bool (NodeDef is defined in tensorflow/core/framework/node_def.proto)
```

that is used as a call-back function to determine which part of the tensorflow `GraphDef` given by tracing `func` will be placed on Machine Learning Accelerators.

If `subgraph_builder_function` is not provided, then `trace` will automatically place operations on Machine Learning Accelerators or on CPU to maximize the execution efficiency.

If it is provided, and `subgraph_builder_function(node)` returns `True`, and placing `node` on Machine Learning Accelerators will not cause deadlocks during execution, then `trace` will place `node` on Machine Learning Accelerators. If `subgraph_builder_function(node)` returns `False`, then `trace` will place `node` on CPU.

## Special Flags

These are flags that get passed directly to the Neuron tracing API (rather than the Neuron Compiler). The flags are still passed via the environment variable `NEURON_CC_FLAGS`.

- **workdir:** example usage - `NEURON_CC_FLAGS='--workdir ./artifacts'` will create a folder named `artifacts` in the current directory and save artifacts that can be used for debug.
- **dynamic-batch-size:** example usage - `NEURON_CC_FLAGS='--dynamic-batch-size'` A flag to allow Neuron graphs to consume variable sized batches of data. Dynamic sizing is restricted to the 0th dimension of a tensor.
- **extract-weights (Beta):** example usage - `NEURON_CC_FLAGS='--extract-weights trn1.2xlarge'` will reduce the compiled model's protobuf size by taking the weights out of the protobuf. Useful for compiling large models that would exceed the 2GB protobuf size limit. This feature is in beta. Model performance is not guaranteed and the flag does not work in combination with `--neuroncore-pipeline-cores`, `--dynamic-batch-size`, models with multiple NEFFs, and models that are 16GB or greater. Compiles models for different neuron instances depending on the instance type passed. Supports all `trn1` and `inf2` instance types except for `trn1n`.

## Returns

- An AWS-Neuron-optimized `keras.Model`.

## Example Usage

```
import tensorflow as tf
import tensorflow_neuronx as tfnx

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
model = tf.keras.Model(inputs=[input0], outputs=[dense0])
example_inputs = tf.random.uniform([1, 3])
model_neuron = tfnx.trace(model, example_inputs) # trace
# check to see how much of the model was compiled successfully
print(model_neuron.on_neuron_ratio)

model_dir = './model_neuron'
model_neuron.save(model_dir)
model_neuron_reloaded = tf.keras.models.load_model(model_dir)
```

## Example Usage with Manual Device Placement Using `subgraph_builder_function`

```
import tensorflow as tf
import tensorflow_neuronx as tfnx

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
reshape0 = tf.keras.layers.Reshape([1, 3])(dense0)
output0 = tf.keras.layers.Dense(2)(reshape0)
model = tf.keras.Model(inputs=[input0], outputs=[output0])
example_inputs = tf.random.uniform([1, 3])

def subgraph_builder_function(node):
    return node.op == 'MatMul'

model_neuron = tfnx.trace(
    model, example_inputs,
    subgraph_builder_function=subgraph_builder_function,
)
```

*This document is relevant for: Inf2, Trn1*

*This document is relevant for: Inf2, Trn1*

## TensorFlow 2.x (tensorflow-neuronx) Auto Multicore Replication (Beta)

The Neuron auto multicore replication Python API enables modifying TensorFlow 2.x models trace by `tensorflow_neuronx.trace` so that they can be automatically replicated across multiple cores.

### Table of contents

- [TensorFlow 2.x \(tensorflow-neuron TF2.x\) Auto Multicore Replication Python API \(Beta\)](#)

- *TensorFlow Neuron TF2.x (tensorflow-neuronx TF2.x) Auto Multicore Replication CLI (Beta)*

## TensorFlow 2.x (tensorflow-neuron TF2.x) Auto Multicore Replication Python API (Beta)

### Method

`tensorflow.neuron.auto_multicore` on models traced by `tensorflow_neuronx.trace`

### Description

Converts an existing AWS-Neuron-optimized `keras.Model` and returns an auto-replication tagged AWS-Multicore-Neuron-optimized `keras.Model` that can execute on AWS Machine Learning Accelerators. Like the traced model, the returned `keras.Model` will support inference only. Attributes or variables held by the original function or `keras.Model` will be dropped.

The auto model replication feature in TensorFlow-Neuron enables you to create a model once and the model parallel replication would happen automatically. The desired number of cores can be less than the total available NeuronCores on an trn1 or inf2 instance but not less than 1. This reduces framework memory usage as you are not loading the same model multiple times manually. Calls to the returned model will execute the call on each core in a round-robin fashion.

The returned `keras.Model` can be exported as SavedModel and served using TensorFlow Serving. Please see [tensorflow-serving](#) for more information about exporting to saved model and serving using TensorFlow Serving.

Note that the automatic replication will only work on models compiled with pipeline size 1: via `--neuroncore-pipeline-cores=1`. If auto replication is not enabled, the model will default to replicate on up to 4 cores.

See *Neuron Compiler CLI Reference Guide (neuronx-cc)* for more information about compiler options.

### Arguments

- **func:** The `keras.Model` or function to be traced.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.
- **num\_cores:** The desired number of cores where the model will be automatically replicated across

### Returns

- An AWS-Multicore-Neuron-optimized `keras.Model`.

**Example Python API Usage for TF2.x traced models:**

```
import tensorflow as tf
import tensorflow.neuron as tfn
import tensorflow_neuronx as tfnx

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
inputs = [input0]
outputs = [dense0]
model = tf.keras.Model(inputs=inputs, outputs=outputs)
input0_tensor = tf.random.uniform([1, 3])
model_neuron = tfnx.trace(model, input0_tensor)

# a trn1.2xlarge has 2 neuron cores
num_cores = 2
multicore_model = tfn.auto_multicore(model_neuron, input0_tensor, num_cores=num_cores)
multicore_model(input0_tensor)
```

**Example Python API Usage for TF2.x saved models:**

```
from tensorflow.python import saved_model

input0_tensor = tf.random.uniform([1, 3])
num_cores = 4
reload_model = saved_model.load(model_dir)
multicore_model = tfn.auto_multicore(reload_model, input0_tensor, num_cores=num_cores)
```

**TensorFlow Neuron TF2.x (tensorflow-neuronx TF2.x) Auto Multicore Replication CLI (Beta)**

The Neuron auto multicore replication CLI enables modifying Tensorflow 2.x traced saved models so that they can be automatically replicated across multiple cores. By performing this call on Tensorflow Saved Models, we can support Tensorflow-Serving without significant modifications to the code.

**Method**

```
tf-neuron-auto-multicore MODEL_DIR --num_cores NUM_CORES --new_model_dir NEW_MODEL_DIR
```

**Arguments**

- **MODEL\_DIR:** The directory of a saved AWS-Neuron-optimized `keras.Model`.
- **NUM\_CORES:** The desired number of cores where the model will be automatically replicated across
- **NEW\_MODEL\_DIR:** The directory of where the AWS-Multicore-Neuron-optimized `keras.Model` will be saved



**Example CLI Usage for Tensorflow-Serving saved models:**

```
tf-neuron-auto-multicore ./resnet --num_cores 8 --new_model_dir ./modified_resnet
```

*This document is relevant for: Inf2, Trn1*

*This document is relevant for: Inf2, Trn1*

**TensorFlow 2.x (tensorflow-neuronx) analyze\_model API****Method**

```
tensorflow_neuronx.analyze_model
```

**Description**

Analyzes a `keras.Model` or a Python callable that can be decorated by `tf.function` for its compatibility with Neuron. It displays supported vs. unsupported operators in the model as well as percentages and counts of each operator and returns a dictionary with operator statistics.

**Arguments**

- **func:** The `keras.Model` or function to be analyzed.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.

**Returns**

- A results dict with these keys: `percent_supported`, `supported_count`, `total_count`, `supported_operators`, `unsupported_operators`, `operators`, `operator_count`.

**Example Usage**

```
import tensorflow as tf
import tensorflow_neuronx as tfnx

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
model = tf.keras.Model(inputs=[input0], outputs=[dense0])
example_inputs = tf.random.uniform([1, 3])
results = tfnx.analyze_model(model, example_inputs)
print(results)
```

(continues on next page)

(continued from previous page)

```
# expected output
"""
BiasAdd
  MatMul
  100.00% of all operations (2 of 2) are supported
  {'percent_supported': 100.0, 'supported_count': 2, 'total_count': 2,
   'supported_operators': {'BiasAdd', 'MatMul'}, 'unsupported_operators': [],
   'operators': ['BiasAdd', 'MatMul'], 'operator_count': {'MatMul': 1, 'BiasAdd': 1}}
"""
```

This document is relevant for: Inf2, Trn1

- *TensorFlow 2.x (tensorflow-neuronx) Tracing API*
- *TensorFlow 2.x (tensorflow-neuronx) Auto Multicore Replication (Beta)*
- *TensorFlow 2.x (tensorflow-neuronx) analyze\_model API*

This document is relevant for: Inf2, Trn1

This document is relevant for: Inf2, Trn1

## Misc (tensorflow-neuronx)

This document is relevant for: Inf2, Trn1, Trn2

## TensorFlow 2.x (tensorflow-neuronx) Release Notes

### Table of contents

- *tensorflow-neuronx 2.x release [2.1.0]*
- *tensorflow-neuronx 2.10 release [2.0.0]*
- *tensorflow-neuronx 2.10 release [1.0.0]*

This document lists the release notes for the tensorflow-neuronx 2.x packages.

### tensorflow-neuronx 2.x release [2.1.0]

Date: 09/15/2023

- Minor updates

Date: 05/1/2023

- Added support for tracing models larger than 2 GB through the environment variable `NEURON_CC_FLAGS='--extract-weights INSTANCE_TYPE'` for all trn1 and inf2 instance types.
- tensorflow-neuronx now supports tensorflow 2.7, 2.8, and 2.9 (In addition to the already supported 2.10).
- Neuron release 2.10 release will be the last release that will include support for tensorflow-neuronx version 2.7. Future Neuron releases will not include tensorflow-neuronx version 2.7.

**tensorflow-neuronx 2.10 release [2.0.0]**

Date: 03/28/2023

The second release of tensorflow-neuronx 2.10 includes the following features:

- Dynamic batching

The following features are not included in this release:

- Support for tracing models larger than 2 GB

**tensorflow-neuronx 2.10 release [1.0.0]**

Date: 2/24/2023

The initial release of tensorflow-neuronx 2.10 includes the following features:

- Initial support for TensorFlow 2.10 inference on Inf2 and Trn1
- Trace API (tensorflow\_neuronx.trace)
- Automatic partitioning of model into CPU vs NeuronCore parts
- Automatic data parallel on multiple NeuronCores (beta)
- Python 3.7, 3.8 and 3.9 support
- HuggingFace Roberta tutorial

The following features are not included in this release:

- Dynamic batching
- Support for tracing models larger than 2 GB

*This document is relevant for: Inf2, Trn1, Trn2*

- *[TensorFlow 2.x \(tensorflow-neuronx\) Release Notes](#)*

*This document is relevant for: Inf2, Trn1*

Setup (tensorflow-neuronx)

**Tutorials (tensorflow-neuronx)**

- HuggingFace Roberta-Base [\[html\]](#) [\[notebook\]](#)
- *[Using NEURON\\_RT\\_VISIBLE\\_CORES with TensorFlow Serving](#)*

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
-

**API Reference Guide (tensorflow-neuronx)**

- *TensorFlow 2.x (tensorflow-neuronx) Tracing API*
- *TensorFlow 2.x (tensorflow-neuronx) Auto Multicore Replication (Beta)*
- *TensorFlow 2.x (tensorflow-neuronx) analyze\_model API*

**Misc (tensorflow-neuronx)**

- *TensorFlow 2.x (tensorflow-neuronx) Release Notes*

*This document is relevant for:* Inf2, Trn1

*This document is relevant for:* Inf1

**2.3.3 Inference on Inf1 (tensorflow-neuron)**

*This document is relevant for:* Inf1

**Tutorials (tensorflow-neuron)**

*This document is relevant for:* Inf1

**Natural Language Processing (NLP) Tutorials (tensorflow-neuron)**

- Tensorflow 2.x - HuggingFace DistilBERT with Tensorflow2 Neuron [\[html\]](#) [\[notebook\]](#)

*This document is relevant for:* Inf1

**[Broken] Running TensorFlow BERT-Large with AWS Neuron**

This example shows a Neuron compatible BERT-Large implementation that is functionally equivalent to open source BERT-Large model. This demo uses TensorFlow-Neuron, BERT-Large weights fine tuned for MRPC and also shows the performance achieved by the Inf1 instance. For users who want to use public BERT SavedModels please also follow the steps described [Using public BERT SavedModels](#).

**Launch EC2 instances**

For this demo, launch two EC2 instances :

- a c5.4xlarge instance for compiling the BERT-Large Model and
- an inf1.xlarge instance for running inference

For both of these instances choose the latest Ubuntu 18 Deep Learning AMI (DLAMI).

## Compiling Neuron compatible BERT-Large

First connect to a c5.4xlarge instance and update tensorflow-neuron and neuron-cc

### Update compilation EC2 instance

Update to the latest neuron software by executing the instructions at [install-neuron-tensorflow](#).

Note: if your tensorflow-neuron version on the inference instance is lower than 1.15.0.1.0.1333.0, you will need to run this demo on inf1.2xlarge instead of inf1.xlarge.

### Compile open source BERT-Large saved model using Neuron compatible BERT-Large implementation

Neuron software works with TensorFlow saved models. Users should bring their own BERT-Large saved model for this section. This demo will run inference for the MRPC task and the saved model should be fine tuned for MRPC. Users who need additional help to fine-tune the model for MRPC or to create a saved model can refer to [Appendix 1](#).

In the same environment and directory bert\_demo scripts, run the following :

```
git clone https://github.com/aws/aws-neuron-sdk
cd ~/aws-neuron-sdk/src/examples/tensorflow/bert_demo/
export BERT_LARGE_SAVED_MODEL="/path/to/user/bert-large/savedmodel"
pip install tensorflow_neuron==1.15.5.2.8.9.0 --extra-index-url=https://pip.repos.neuron.
↳amazonaws.com/
pip install neuron_cc==1.13.5.0 --extra-index-url=https://pip.repos.neuron.amazonaws.com
python bert_model.py --input_saved_model $BERT_LARGE_SAVED_MODEL --output_saved_model ./
↳bert-saved-model-neuron --batch_size=6 --aggressive_optimizations
```

This compiles BERT-Large pointed to by \$BERT\_LARGE\_SAVED\_MODEL for an input size of 128 and batch size of 6. The compilation output is stored in bert-saved-model-neuron. Copy this to your Inf1 instance for inferencing.

The bert\_model.py script encapsulates all the steps necessary for this process. For details on what is done by bert\_model.py please refer to [Appendix 2](#).

### Running the inference demo

Connect to your inf1.xlarge instance and update tensorflow-neuron, aws-neuron-runtime and aws-neuron-tools.

### Update inference EC2 instance

Update to the latest neuron software by executing the instructions at [install-neuron-tensorflow](#).

## Launching the BERT-Large demo server

Copy the compiled model (bert-saved-model-neuron) from your c5.4xlarge to your inf1.xlarge instance. Place the model in the same directory as the bert\_demo scripts. Then from the same conda environment launch the BERT-Large demo server :

```
cd ~/aws-neuron-sdk/src/examples/tensorflow/bert_demo/  
pip install tensorflow_neuron==1.15.5.2.8.9.0 --extra-index-url=https://pip.repos.neuron.  
amazonaws.com/  
python bert_server.py --dir bert-saved-model-neuron --batch 6 --parallel 4
```

This loads 4 BERT-Large models, one into each of the 4 NeuronCores found in an inf1.xlarge instance. For each of the 4 models, the BERT-Large demo server opportunistically stitches together asynchronous requests into batch 6 requests. When there are insufficient pending requests, the server creates dummy requests for batching.

Wait for the bert\_server to finish loading the BERT-Large models to Inferentia memory. When it is ready to accept requests it will print the inferences per second once every second. This reflects the number of real inferences only. Dummy requests created for batching are not credited to inferentia performance. Once the inferences are done you can send a keyboard interrupt to print out the average throughput of your run.

## Sending requests to server from multiple clients

Wait until the bert demo server is ready to accept requests. Then on the same inf1.xlarge instance, launch a separate linux terminal. From the bert\_demo directory execute the following commands :

```
source activate aws_neuron_tensorflow_p36  
cd ~/aws-neuron-sdk/src/examples/tensorflow/bert_demo/  
for i in {1..96}; do python bert_client.py --cycle 128 & done
```

This spins up 96 clients, each of which sends 128 inference requests.

## Printing latency metrics

After all your requests have been sent to your server you can run the following command:

```
python latency_printer.py
```

## Using public BERT SavedModels

We are now providing a compilation script that has better compatibility with various flavors of BERT SavedModels generated from <https://github.com/google-research/bert>. Here are the current limitations:

1. You did not change `modeling.py`
2. BERT SavedModel is generated using `estimator.export_saved_model`
3. BERT SavedModel uses fixed sequence length 128 (you may check by `saved_model_cli show --dir /path/to/user/bert/savedmodel --all`)
4. neuron-cc version is at least 1.0.12000.0
5. aws-neuron-runtime version is at least 1.0.7000.0
6. The `--batch_size` argument specified in this script is at most 4

Example usage is shown below:

```
export BERT_LARGE_SAVED_MODEL="/path/to/user/bert-large/savedmodel"
cd ~/aws-neuron-sdk/src/examples/tensorflow/bert_demo/
python bert_no_model.py --input_saved_model $BERT_LARGE_SAVED_MODEL --output_saved_model_
↪ ./bert-saved-model-neuron --batch_size=1
```

## Appendix 1

Users who need help finetuning BERT-Large for MRPC and creating a saved model may follow the instructions here.

Connect to the c5.4xlarge compilation EC2 instance you started above and download these three items :

1. clone [this](#) github repo.
2. download GLUE data as described [here](#). Do not run the finetuning command.
3. download a desired pre-trained BERT-Large checkpoint from [here](#). This is the model we will fine tune.

Next edit run\_classifier.py in the cloned bert repo to apply the patch described in the following git diff.

```
diff --git a/run_classifier.py b/run_classifier.py
index 817b147..c9426bc 100644
--- a/run_classifier.py
+++ b/run_classifier.py
@@ -955,6 +955,18 @@ def main(_):
     drop_remainder=predict_drop_remainder)

     result = estimator.predict(input_fn=predict_input_fn)
+   features = {
+       "input_ids": tf.placeholder(shape=[None, FLAGS.max_seq_length], dtype=tf.int32,
+↪ name='input_ids'),
+       "input_mask": tf.placeholder(shape=[None, FLAGS.max_seq_length], dtype=tf.int32,
+↪ name='input_mask'),
+       "segment_ids": tf.placeholder(shape=[None, FLAGS.max_seq_length], dtype=tf.
+↪ int32, name='segment_ids'),
+       "label_ids": tf.placeholder(shape=[None], dtype=tf.int32, name='label_ids'),
+       "is_real_example": tf.placeholder(shape=[None], dtype=tf.int32, name='is_real_
+↪ example'),
+   }
+   serving_input_fn = tf.estimator.export.build_raw_serving_input_receiver_fn(features)
+   estimator._export_to_tpu = False ## !!important to add this
+   estimator.export_saved_model(
+       export_dir_base='./bert_classifier_saved_model',
+       serving_input_receiver_fn=serving_input_fn)

    output_predict_file = os.path.join(FLAGS.output_dir, "test_results.tsv")
    with tf.gfile.GFile(output_predict_file, "w") as writer:
```

NOTE : Users who are interested may refer to this [link](#) for additional background information on the patch but it is not necessary for running this demo.

Then from the bert\_demo directory run the following :

```
source activate aws_neuron_tensorflow_p36
cd ~/aws-neuron-sdk/src/examples/tensorflow/bert_demo/
export BERT_REPO_DIR="/path/to/cloned/bert/repo/directory"
export GLUE_DIR="/path/to/glue/data/directory"
export BERT_BASE_DIR="/path/to/pre-trained/bert-large/checkpoint/directory"
./tune_save.sh
```

The a saved model will be created in `$BERT_REPO_DIR/bert-saved-model/random_number/`. Where, *random\_number* is a random number generated for every run. Use this saved model to continue with the rest of the demo.

## Appendix 2

For all BERT variants, we currently need to augment the standard Neuron compilation process for performance tuning. In the future, we intend to automate this tuning process. This would allow users to use the standard Neuron compilation process, which requires only a one line change in user source code. The standard compilation process is described `/src/examples/mxnet/resnet50/resnet50.ipynb`.

The augmented Neuron compilation process is encapsulated by the `bert_model.py` script, which performs the following things :

1. Define a Neuron compatible implementation of BERT-Large. For inference, this is functionally equivalent to the open source BERT-Large. The changes needed to create a Neuron compatible BERT-Large implementation is described in [Appendix 3](#).
2. Extract BERT-Large weights from the open source saved model pointed to by `-input_saved_model` and associates it with the Neuron compatible model
3. Invoke TensorFlow-Neuron to compile the Neuron compatible model for Inferentia using the newly associated weights
4. Finally, the compiled model is saved into the location given by `-output_saved_model`

## Appendix 3

The Neuron compatible implementation of BERT-Large is functionally equivalent to the open source version when used for inference. However, the detailed implementation does differ and here are the list of changes :

1. Data Type Casting : If the original BERT-Large an FP32 model, `bert_model.py` contains manually defined cast operators to enable mixed-precision. FP16 is used for multi-head attention and fully-connected layers, and fp32 everywhere else. This will be automated in a future release.
2. Remove Unused Operators: A model typically contains training operators that are not used in inference, including a subset of the reshape operators. Those operators do not affect inference functionality and have been removed.
3. Reimplementation of Selected Operators : A number of operators (mainly mask operators), has been reimplemented to bypass a known compiler issue. This will be fixed in a planned future release.
4. Manually Partition Embedding Ops to CPU : The embedding portion of BERT-Large has been partitioned manually to a subgraph that is executed on the host CPU, without noticable performance impact. In near future, we plan to implement this through compiler auto-partitioning without the need for user intervention.

*This document is relevant for: Inf1*



## Running Huggingface DistilBERT with TensorFlow-Neuron

In this tutorial you will compile and deploy DistilBERT version of HuggingFace Transformers BERT for Inference using TensorFlow-Neuron. The full list of HuggingFace's pretrained BERT models can be found in the BERT section on this page [https://huggingface.co/transformers/pretrained\\_models.html](https://huggingface.co/transformers/pretrained_models.html). you can also read about HuggingFace's pipeline feature here: [https://huggingface.co/transformers/main\\_classes/pipelines.html](https://huggingface.co/transformers/main_classes/pipelines.html)

This Jupyter notebook should be run on an instance which is inf1.6xlarge or larger, but in real life scenario the compilation should be done on a compute instance and the deployment on inf1 instance to save costs.

### Setup

To run this tutorial please follow the instructions for [TensorFlow-Neuron Setup](#) and the [Jupyter Notebook Quickstart](#) and set your kernel to "Python (tensorflow-neuron)".

Next, install some additional dependencies.

```
[ ]: %env TOKENIZERS_PARALLELISM=True #Supresses tokenizer warnings making errors easier to
    ↪ detect
    !pip install transformers==4.30.2
    !pip install ipywidgets
```

### Download From Huggingface and Compile for AWS-Neuron

```
[ ]: import tensorflow as tf
import tensorflow_neuron as tfn
from transformers import DistilBertTokenizer, TFDistilBertModel

# Create a wrapper for the roberta model that will accept inputs as a list
# instead of a dictionary. This will allow the compiled model to be saved
# to disk with the model.save() function.
class DistilBertWrapper(tf.keras.Model):
    def __init__(self, model):
        super().__init__()
        self.model = model
    def __call__(self, example_inputs):
        return self.model({'input_ids' : example_inputs[0], 'attention_mask' : example_
    ↪ inputs[1]})

tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-uncased-finetuned-sst-2-
    ↪ english')
model = DistilBertWrapper(TFDistilBertModel.from_pretrained('distilbert-base-uncased-
    ↪ finetuned-sst-2-english'))

batch_size = 16

# create example inputs with a batch size of 16
text = ["Paris is the <mask> of France."] * batch_size
encoded_input = tokenizer(text, return_tensors='tf', padding='max_length', max_length=64)
```

(continues on next page)

(continued from previous page)

```
# turn inputs into a list
example_input = [encoded_input['input_ids'], encoded_input['attention_mask']]

#compile
model_neuron = tfn.trace(model, example_input)

print("Running on neuron:", model_neuron(example_input))

# save the model to disk to save recompilation time for next usage
model_neuron.save('./distilbert-neuron-b16')
```

### Run Basic Inference Benchmarking

```
[ ]: import numpy as np
import concurrent.futures
import time

reloaded_neuron_model = tf.keras.models.load_model('./distilbert-neuron-b16')
print("Reloaded model running on neuron:", reloaded_neuron_model(example_input))

num_threads = 4
num_inferences = 1000

latency_list = []
def inference_with_latency_calculation(example_input):
    global latency_list
    start = time.time()
    result = reloaded_neuron_model(example_input)
    end = time.time()
    latency_list.append((end-start) * 1000)
    return result

start = time.time()
with concurrent.futures.ThreadPoolExecutor(max_workers=num_threads) as executor:
    futures = []
    for i in range(num_inferences):
        futures.append(executor.submit(inference_with_latency_calculation, example_
↪input))
    for future in concurrent.futures.as_completed(futures):
        get_result = future.result()
end = time.time()

total_time = end - start
throughput = (num_inferences * batch_size)/total_time

print(f"Throughput was {throughput} samples per second.")
print(f"Latency p50 was {np.percentile(latency_list, 50)} ms")
print(f"Latency p90 was {np.percentile(latency_list, 90)} ms")
print(f"Latency p95 was {np.percentile(latency_list, 95)} ms")
print(f"Latency p99 was {np.percentile(latency_list, 99)} ms")
```

(continues on next page)

(continued from previous page)

```
assert(throughput >= 1930.0)
```

[ ]:

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

### Utilizing Neuron Capabilities Tutorials (tensorflow-neuron)

- Using NEURON\_RT\_VISIBLE\_CORES with TensorFlow Serving [html]

*This document is relevant for: Inf1*

### Natural Language Processing (NLP) Tutorials

- Tensorflow 2.x - HuggingFace Pipelines distilBERT with Tensorflow2 Neuron [html] [notebook]

### Utilizing Neuron Capabilities Tutorials

- Tensorflow 2.x - Using NEURON\_RT\_VISIBLE\_CORES with TensorFlow Serving [html]

---

**Note:** To use Jupyter Notebook see:

- setup-jupyter-notebook-steps-troubleshooting
  - running-jupyter-notebook-as-script
- 

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

### Additional Examples (tensorflow-neuron)

- [AWS Neuron Samples GitHub Repository](#)

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

### API Reference Guide (tensorflow-neuron)

*This document is relevant for: Inf1*

## TensorFlow 2.x (tensorflow-neuron) Tracing API

The Neuron tracing API enables tracing TensorFlow 2.x models for deployment on AWS Machine Learning Accelerators.

### Method

`tensorflow.neuron.trace`

### Description

Trace a `keras.Model` or a Python callable that can be decorated by `tf.function`, and return an AWS-Neuron-optimized `keras.Model` that can execute on AWS Machine Learning Accelerators. Tracing is ideal for `keras.Model` that accepts a list of `tf.Tensor` objects and returns a list of `tf.Tensor` objects. It is expected that users will provide example inputs, and the `trace` function will execute `func` symbolically and convert it to a `keras.Model`.

The returned `keras.Model` will support inference only. Attributes or variables held by the original function or `keras.Model` will be dropped.

The returned `keras.Model` can be exported as `SavedModel` and served using TensorFlow Serving. Please see `tensorflow-serving` for more information about exporting to saved model and serving using TensorFlow Serving.

The returned `keras.Model` has an `.on_neuron_ratio` attribute which shows the percentage of ops mapped to neuron hardware. This calculation ignores PlaceholderOp, IdentityOp, ReadVariableOp and NoOp.

Options can be passed to Neuron compiler via the environment variable `NEURON_CC_FLAGS`. For example, the syntax `env NEURON_CC_FLAGS="--neuroncore-pipeline-cores=4"` directs Neuron compiler to compile each subgraph to fit in the specified number of NeuronCores. This number can be less than the total available NeuronCores on an Inf1 instance. See *Neuron compiler CLI Reference Guide (neuron-cc)* for more information about compiler options.

### Arguments

- **func:** The `keras.Model` or function to be traced.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.
- **subgraph\_builder\_function:** (Optional) A callable with signature

```
subgraph_builder_function(node : NodeDef) -> bool (NodeDef is defined in tensorflow/core/framework/node_def.proto)
```

that is used as a call-back function to determine which part of the tensorflow `GraphDef` given by tracing `func` will be placed on Machine Learning Accelerators.

If `subgraph_builder_function` is not provided, then `trace` will automatically place operations on Machine Learning Accelerators or on CPU to maximize the execution efficiency.

If it is provided, and `subgraph_builder_function(node)` returns `True`, and placing `node` on Machine Learning Accelerators will not cause deadlocks during execution, then `trace` will place `node` on Machine Learning Accelerators. If `subgraph_builder_function(node)` returns `False`, then `trace` will place `node` on CPU.

## Special Flags

These are flags that get passed directly to the Neuron tracing API (rather than the Neuron Compiler). The flags are still passed via the environment variable `NEURON_CC_FLAGS`.

- **workdir:** example usage - `NEURON_CC_FLAGS='--workdir ./artifacts'` will create a folder named `artifacts` in the current directory and save artifacts that can be used for debug.
- **dynamic-batch-size:** example usage - `NEURON_CC_FLAGS='--dynamic-batch-size'` A flag to allow Neuron graphs to consume variable sized batches of data. Dynamic sizing is restricted to the 0th dimension of a tensor.
- **extract-weights (Beta):** example usage - `NEURON_CC_FLAGS='--extract-weights inf1.2xlarge'` will reduce the compiled model's protobuf size by taking the weights out of the protobuf. Useful for compiling large models that would exceed the 2GB protobuf size limit. This feature is in beta. Model performance is not guaranteed and the flag does not work in combination with `--neuroncore-pipeline-cores`, `--dynamic-batch-size`, models with multiple NEFFs, and models that are 4GB or greater. Compiles models for different neuron instances depending on the instance type passed. Supports all `inf1` instance types.

## Returns

- An AWS-Neuron-optimized `keras.Model`.

## Example Usage

```
import tensorflow as tf
import tensorflow.neuron as tfn

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
model = tf.keras.Model(inputs=[input0], outputs=[dense0])
example_inputs = tf.random.uniform([1, 3])
model_neuron = tfn.trace(model, example_inputs) # trace
# check to see how much of the model was compiled successfully
print(model_neuron.on_neuron_ratio)

model_dir = './model_neuron'
model_neuron.save(model_dir)
model_neuron_reloaded = tf.keras.models.load_model(model_dir)
```

## Example Usage with Manual Device Placement Using `subgraph_builder_function`

```
import tensorflow as tf
import tensorflow.neuron as tfn

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
reshape0 = tf.keras.layers.Reshape([1, 3])(dense0)
output0 = tf.keras.layers.Dense(2)(reshape0)
model = tf.keras.Model(inputs=[input0], outputs=[output0])
example_inputs = tf.random.uniform([1, 3])
```

(continues on next page)

(continued from previous page)

```
def subgraph_builder_function(node):  
    return node.op == 'MatMul'  
  
model_neuron = tfn.trace(  
    model, example_inputs,  
    subgraph_builder_function=subgraph_builder_function,  
)
```

---

**Important:** Although the old API `tensorflow.neuron.saved_model.compile` is still available under `tensorflow-neuron 2.x`, it supports only the limited capabilities of `tensorflow.neuron.trace` and will be deprecated in future releases.

---

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## TensorFlow 2.x (tensorflow-neuron) analyze\_model API

### Method

`tensorflow.neuron.analyze_model`

### Description

Analyzes a `keras.Model` or a Python callable that can be decorated by `tf.function` for its compatibility with Neuron. It displays supported vs. unsupported operators in the model as well as percentages and counts of each operator and returns a dictionary with operator statistics.

### Arguments

- **func:** The `keras.Model` or function to be analyzed.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.

## Returns

- A results dict with these keys: ``percent\_supported`, `supported\_count`, `total\_count`, `supported\_operators`, `unsupported\_operators`, `operators`, `operator\_count``.

## Example Usage

```
import tensorflow as tf
import tensorflow.neuron as tfn

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
model = tf.keras.Model(inputs=[input0], outputs=[dense0])
example_inputs = tf.random.uniform([1, 3])
results = tfn.analyze_model(model, example_inputs)
print(results)

# expected output
'''
BiasAdd
  MatMul
  100.00% of all operations (2 of 2) are supported
  {'percent_supported': 100.0, 'supported_count': 2, 'total_count': 2,
  'supported_operators': {'BiasAdd', 'MatMul'}, 'unsupported_operators': [],
  'operators': ['BiasAdd', 'MatMul'], 'operator_count': {'MatMul': 1, 'BiasAdd': 1}}
'''
```

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## TensorFlow 2.x (tensorflow-neuron) Auto Multicore Replication (Beta)

The Neuron auto multicore replication Python API enables modifying TensorFlow 2.x traced models so that they can be automatically replicated across multiple cores. For Tensorflow-Serving models and TensorFlow 1.x models, see *TensorFlow Neuron TF2.x (tensorflow-neuronx TF2.x) Auto Multicore Replication CLI (Beta)*

### Table of contents

- *TensorFlow 2.x (tensorflow-neuron TF2.x) Auto Multicore Replication Python API (Beta)*
- *TensorFlow Neuron 2.x (tensorflow-neuron) Auto Multicore Replication CLI (Beta)*

## TensorFlow 2.x (tensorflow-neuron TF2.x) Auto Multicore Replication Python API (Beta)

### Method

`tensorflow.neuron.auto_multicore`

### Description

Converts an existing AWS-Neuron-optimized `keras.Model` and returns an auto-replication tagged AWS-Multicore-Neuron-optimized `keras.Model` that can execute on AWS Machine Learning Accelerators. Like the traced model, the returned `keras.Model` will support inference only. Attributes or variables held by the original function or `keras.Model` will be dropped.

The auto model replication feature in TensorFlow-Neuron enables you to create a model once and the model parallel replication would happen automatically. The desired number of cores can be less than the total available NeuronCores on an Inf1 instance but not less than 1. This reduces framework memory usage as you are not loading the same model multiple times manually. Calls to the returned model will execute the call on each core in a round-robin fashion.

The returned `keras.Model` can be exported as `SavedModel` and served using TensorFlow Serving. Please see `tensorflow-serving` for more information about exporting to saved model and serving using TensorFlow Serving.

Note that the automatic replication will only work on models compiled with pipeline size 1: via `--neuroncore-pipeline-cores=1`. If auto replication is not enabled, the model will default to replicate on up to 4 cores.

See *Neuron compiler CLI Reference Guide (neuron-cc)* for more information about compiler options.

### Arguments

- **func:** The `keras.Model` or function to be traced.
- **example\_inputs:** A `tf.Tensor` or a tuple/list/dict of `tf.Tensor` objects for tracing the function. When `example_inputs` is a `tf.Tensor` or a list of `tf.Tensor` objects, we expect `func` to have calling signature `func(example_inputs)`. Otherwise, the expectation is that inference on `func` is done by calling `func(*example_inputs)` when `example_inputs` is a tuple, or `func(**example_inputs)` when `example_inputs` is a dict. The case where `func` accepts mixed positional and keyword arguments is currently unsupported.
- **num\_cores:** The desired number of cores where the model will be automatically replicated across

### Returns

- An AWS-Multicore-Neuron-optimized `keras.Model`.



**Example Python API Usage for TF2.x traced models:**

```

input0 = tf.keras.layers.Input(3)
dense0 = tf.keras.layers.Dense(3)(input0)
inputs = [input0]
outputs = [dense0]
model = tf.keras.Model(inputs=inputs, outputs=outputs)
input0_tensor = tf.random.uniform([1, 3])
model_neuron = tfn.trace(model, input0_tensor)

num_cores = 4
multicore_model = tfn.auto_multicore(model_neuron, input0_tensor, num_cores=num_cores)
multicore_model(input0_tensor)

```

**Example Python API Usage for TF2.x saved models:**

```

from tensorflow.python import saved_model

input0_tensor = tf.random.uniform([1, 3])
num_cores = 4
reload_model = saved_model.load(model_dir)
multicore_model = tfn.auto_multicore(reload_model, input0_tensor, num_cores=num_cores)

```

**TensorFlow Neuron 2.x (tensorflow-neuron) Auto Multicore Replication CLI (Beta)**

The Neuron auto multicore replication CLI enables modifying TensorFlow 1.x and Tensorflow 2.x traced saved models so that they can be automatically replicated across multiple cores. By performing this call on Tensorflow Saved Models, we can support both Tensorflow-Serving and Tensorflow 1.x without significant modifications to the code. Note that the python API does not support Tensorflow 1.x.

**Method**

```
tf-neuron-auto-multicore MODEL_DIR --num_cores NUM_CORES --new_model_dir NEW_MODEL_DIR
```

**Arguments**

- **MODEL\_DIR:** The directory of a saved AWS-Neuron-optimized `keras.Model`.
- **NUM\_CORES:** The desired number of cores where the model will be automatically replicated across
- **NEW\_MODEL\_DIR:** The directory of where the AWS-Multicore-Neuron-optimized `keras.Model` will be saved

*This document is relevant for: Inf1*

- *TensorFlow 2.x (tensorflow-neuron) Tracing API*
- *TensorFlow 2.x (tensorflow-neuron) analyze\_model API*
- *TensorFlow 2.x (tensorflow-neuron) Auto Multicore Replication (Beta)*

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## **Misc (tensorflow-neuron)**

*This document is relevant for: Inf1*

## **TensorFlow 2.x (tensorflow-neuron) Release Notes**

### **Table of contents**

- *Known Issues and Limitations - updated 08/12/2021*
- *tensorflow-neuron 2.x release [2.12.2.0]*
- *tensorflow-neuron 2.x release [2.11.4.0]*
- *tensorflow-neuron 2.x release [2.10.19.0]*
- *tensorflow-neuron 2.x release [2.10.8.0]*
- *tensorflow-neuron 2.x release [2.10.2.0]*
- *tensorflow-neuron 2.x release [2.10.1.0]*
- *tensorflow-neuron 2.x release [2.9.3.0]*
- *tensorflow-neuron 2.x release [2.8.9.0]*
- *tensorflow-neuron 2.x release [2.8.1.0]*
- *tensorflow-neuron 2.x release [2.7.4.0]*
- *tensorflow-neuron 2.x release [2.7.3.0]*
- *tensorflow-neuron 2.x release [2.6.5.0]*
- *tensorflow-neuron 2.x release [2.6.0.0]*
- *tensorflow-neuron 2.x release [2.4.0.0]*
- *tensorflow-neuron 2.x release [2.3.0.0]*
- *tensorflow-neuron 2.x release [2.2.0.0]*
- *tensorflow-neuron 2.x release [2.1.14.0]*
- *tensorflow-neuron 2.x release [2.1.13.0]*
- *tensorflow-neuron 2.x release [2.1.6.0]*
- *tensorflow-neuron 2.x release [2.0.4.0]*
- *tensorflow-neuron 2.x release [2.0.3.0]*
- *tensorflow-neuron 2.x release [1.6.8.0]*

This document lists the release notes for the tensorflow-neuron 2.x packages.

## Known Issues and Limitations - updated 08/12/2021

- Support on serialized TensorFlow 2.x custom operators is currently limited. Serializing some operators registered from tensorflow-text through [TensorFlow Hub](#) is going to cause failure in tensorflow.neuron.trace.
- Memory leak exists on latest releases of TensorFlow Neuron for versions 2.1, 2.2, 2.3, and 2.4.
- Issue: When compiling large models, user might run out of memory and encounter this fatal error.

```
terminate called after throwing an instance of 'std::bad_alloc'
```

Solution: run compilation on a c5.4xlarge instance type or larger.

- Issue: When upgrading tensorflow-neuron with `pip install tensorflow-neuron --upgrade`, the following error message may appear, which is caused by pip version being too low.

```
Could not find a version that satisfies the requirement tensorflow<1.16.0,>=1.15.0 (from ↵
↵ tensorflow-neuron)
```

Solution: run a `pip install pip --upgrade` before upgrading tensorflow-neuron.

- Issue: Some Keras routines throws the following error:

```
AttributeError: 'str' object has no attribute 'decode'.
```

Solution: Please downgrade *h5py* by `pip install 'h5py<3'`. This is caused by <https://github.com/TensorFlow/TensorFlow/issues/44467>.

## tensorflow-neuron 2.x release [2.12.2.0]

Date: 09/16/2024

- Minor updates.

## tensorflow-neuron 2.x release [2.11.4.0]

Date: 07/03/2024

- Minor updates.

## tensorflow-neuron 2.x release [2.10.19.0]

Date: 04/01/2024

- Minor updates.

### **tensorflow-neuron 2.x release [2.10.8.0]**

Date: 12/21/2023

- Minor updates.

### **tensorflow-neuron 2.x release [2.10.2.0]**

Date: 10/15/2023

- Minor updates.

### **tensorflow-neuron 2.x release [2.10.1.0]**

Date: 09/15/2023

- Minor updates.

### **tensorflow-neuron 2.x release [2.9.3.0]**

Date: 7/19/2023

- Minor updates.

### **tensorflow-neuron 2.x release [2.8.9.0]**

Date: 06/14/2023

- Added Python 3.10 support.

### **tensorflow-neuron 2.x release [2.8.1.0]**

Date: 05/01/2023

- Added support for tracing models larger than 2 GB through the environment variable `NEURON_CC_FLAGS='--extract-weights INSTANCE_TYPE'` for all `infl` instance types.
- Neuron release 2.10 release will be the last release that will include support for tensorflow-neuron version 2.7. Future Neuron releases will not include tensorflow-neuron version 2.7.

### **tensorflow-neuron 2.x release [2.7.4.0]**

Date: 04/19/2023

- Minor updates.

**tensorflow-neuron 2.x release [2.7.3.0]**

Date: 03/28/2023

- Introduce the `tfn.analyze_model` function that displays information about the supported and unsupported operators of a traceable model.
- Introduce the `on_neuron_ratio` attribute of AWS Optimized Neuron Models returned by `tfn.trace`, which is the percentage of ops on neuron after compilation.

**tensorflow-neuron 2.x release [2.6.5.0]**

Date: 02/24/2023

- Minor updates.

**tensorflow-neuron 2.x release [2.6.0.0]**

Date: 2/24/2023

- Minor bug fixes.

**tensorflow-neuron 2.x release [2.4.0.0]**

Date: 11/22/2022

- Beta support for tracing models larger than 2 GB through environment variable `NEURON_CC_FLAGS='--extract-weights'`.
- Introduce `tfn.auto_multicore` Python API to enable automatic data parallel on multiple NeuronCores.
- Introduce `tf-neuron-auto-multicore` tool to enable automatic data parallel on multiple NeuronCores.
- Deprecated the `NEURONCORE_GROUP_SIZES` environment variable.
- Minor bug fixes.

**tensorflow-neuron 2.x release [2.3.0.0]**

Date: 04/29/2022

- Added support for Tensorflow 2.8.0.
- Added support for Slice operator
- The graph partitioner now prefers to place less compute intensive operators on CPU if the model already contains a large amount of compute intensive operators.
- Fixed [Github issue #408](#), the fix solves data type handling bug in `tfn.trace` when the model contains Conv2D operators.

**tensorflow-neuron 2.x release [2.2.0.0]**

Date: 03/25/2022

- Updated TensorFlow 2.5 to version 2.5.3.
- Added support for TensorFlow 2.6 and 2.7.
- Added a warning message when calling `tfn.saved_model.compile` API. In tensorflow-neuron 2.x you should call *tensorflow.neuron.trace*. `tfn.saved_model.compile` API supports only partial functionality of *tensorflow.neuron.trace* and will be deprecated in the future.

**tensorflow-neuron 2.x release [2.1.14.0]**

Date: 02/17/2022

- Fixed a bug in TensorFlow Neuron versions 2.1, 2.2, 2.3 and 2.4. The fixed bug was causing a memory leak of 128 bytes for each inference.
- Improved warning message when calling deprecated compilation API under tensorflow-neuron 2.x.

**tensorflow-neuron 2.x release [2.1.13.0]**

Date: 02/16/2022

- Fixed a bug that caused a memory leak. The memory leak was approximately 128b for each inference and exists in all versions of TensorFlow Neuron versions part of Neuron 1.16.0 to Neuron 1.17.0 releases. see pre-release-content for exact versions included in each release. This release only addresses the leak in TensorFlow Neuron 2.5. Future release of TensorFlow Neuron will fix the leak in other versions as well (2.1, 2.2, 2.3, 2.4).

**tensorflow-neuron 2.x release [2.1.6.0]**

Date: 01/20/2022

- Updated TensorFlow 2.5 to version 2.5.2.
- Enhanced auto data parallel (e.g. when using `NEURONCORE_GROUP_SIZES=X,Y,Z,W`) to support edge cases.
- Fixed a bug that may cause tensorflow-neuron to generate in some cases scalar gather instruction with incorrect arguments.

**tensorflow-neuron 2.x release [2.0.4.0]**

Date: 11/05/2021

- Updated Neuron Runtime (which is integrated within this package) to `libnrt 2.2.18.0` to fix a container issue that was preventing the use of containers when `/dev/neuron0` was not present. See details here [neuron-runtime-release-notes](#).

## tensorflow-neuron 2.x release [2.0.3.0]

Date: 10/27/2021

### New in this release

- TensorFlow 2.x (tensorflow-neuron) now support Neuron Runtime 2.x (libnrt.so shared library) only.

---

**Important:**

- You must update to the latest Neuron Driver (aws-neuron-dkms version 2.1 or newer) for proper functionality of the new runtime library.
  - Read *Introducing Neuron Runtime 2.x (libnrt.so)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.
  - Read *Migrate your application to Neuron Runtime 2.x (libnrt.so)* for detailed information of how to migrate your application.
- 

- Updated TensorFlow 2.3.x from TensorFlow 2.3.3 to TensorFlow 2.3.4.
- Updated TensorFlow 2.4.x from TensorFlow 2.4.2 to TensorFlow 2.4.3.
- Updated TensorFlow 2.5.x from TensorFlow 2.5.0 to TensorFlow 2.5.1.

### Resolved Issues

- Fix bug that can cause illegal compiler optimizations
- Fix bug that can cause dynamic-shape operators be placed on Neuron

## tensorflow-neuron 2.x release [1.6.8.0]

Date: 08/12/2021

### New in this release

- First release of TensorFlow 2.x integration, Neuron support now TensorFlow versions 2.1.4, 2.2.3, 2.3.3, 2.4.2, and 2.5.0.
- New public API tensorflow.neuron.trace: trace a TensorFlow 2.x keras.Model or a Python callable that can be decorated by tf.function, and return an AWS-Neuron-optimized keras.Model that can execute on AWS Machine Learning Accelerators.

**Please note** that TensorFlow 1.x SavedModel compilation API tensorflow.neuron.saved\_model.compile is not supported in tensorflow-neuron 2.x . It continues to function in tensorflow-neuron 1.15.x .

- Included versions:
  - tensorflow-neuron-2.5.0.1.6.8.0
  - tensorflow-neuron-2.4.2.1.6.8.0
  - tensorflow-neuron-2.3.3.1.6.8.0

- tensorflow-neuron-2.2.3.1.6.8.0
- tensorflow-neuron-2.1.4.1.6.8.0

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## TensorFlow 2.x (tensorflow-neuron) Accelerated (torch-neuron) Python APIs and Graph Ops

This page lists TensorFlow 2.x Python APIs and graph operators that are accelerated by AWS Neuron. The lists are not exhaustive. TensorFlow 2.x Python APIs or graph operators that are not listed here may still be accelerated if they are composed of accelerated primitives, or they will be executed on CPU without significant acceleration. The TensorFlow Neuron integration contains an automatic operator-device-placement mechanism that strives to maximize the execution efficiency of your deep learning models on AWS Machine Learning ASIC instances.

### Accelerated Python APIs

| Module | Accelerated Python API    | Comments   |
|--------|---------------------------|--|
| tf     | tf.abs                    |  |
|        | tf.add                    |  |
|        | tf.add_n                  |  |
|        | tf.broadcast_static_shape |  |
|        | tf.cast                   |  |
|        | tf.constant               |  |
|        | tf.convert_to_tensor      |  |
|        | tf.cumsum                 | axis must be a compile-time constant.                            |
|        | tf.einsum                 |  |
|        | tf.erf                    |  |
|        | tf.exp                    |  |
|        | tf.identity               |  |
|        | tf.matmul                 | Uses float16/bfloat16 matmul with float32 accumulation.          |
|        | tf.maximum                |  |
|        | tf.minimum                |  |
|        | tf.multiply               |  |
|        | tf.negative               |  |
|        | tf.range                  | start, limit and delta arguments must be compile-time constants. |
|        | tf.realdiv                |  |
|        | tf.reciprocal             |  |
|        | tf.reduce_all             | axis must be a compile-time constant.                            |
|        | tf.reduce_any             | axis must be a compile-time constant.                            |
|        | tf.reduce_max             | axis must be a compile-time constant.                            |
|        | tf.reduce_min             | axis must be a compile-time constant.                            |

continues on next page



Table 2.2 – continued from previous page

| Module                 | Accelerated Python API                     | Comments  |
|------------------------|--|---|
|                        | <code>tf.reduce_prod</code>                | axis must be a compile-time constant.   |
|                        | <code>tf.reduce_sum</code>                 | axis must be a compile-time constant.   |
|                        | <code>tf.reshape</code>                    | shape argument must be a compile-time constant.   |
|                        | <code>tf.rsqrt</code>                      |   |
|                        | <code>tf.scalar_mul</code>                 |   |
|                        | <code>tf.shape</code>                      |   |
|                        | <code>tf.shape_n</code>                    |   |
|                        | <code>tf.sigmoid</code>                    |   |
|                        | <code>tf.size</code>                       |   |
|                        | <code>tf.slice</code>                      | size must be a compile-time constant. In addition, either begin must be a compile-time constant or size must be non-negative. |
|                        | <code>tf.sqrt</code>                       |   |
|                        | <code>tf.square</code>                     |   |
|                        | <code>tf.squared_difference</code>         |   |
|                        | <code>tf.squeeze</code>                    |   |
|                        | <code>tf.stack</code>                      |   |
|                        | <code>tf.stop_gradient</code>              |   |
|                        | <code>tf.strided_slice</code>              |   |
|                        | <code>tf.tanh</code>                       |   |
|                        | <code>tf.tensordot</code>                  |   |
|                        | <code>tf.to_bfloat16</code>                |   |
|                        | <code>tf.to_float</code>                   |   |
|                        | <code>tf.truediv</code>                    |   |
| <code>tf.layers</code> | <code>tf.layers.batch_normalization</code> |   |
|                        | <code>tf.layers.dense</code>               |   |
|                        | <code>tf.layers.flatten</code>             |   |
| <code>tf.nn</code>     | <code>tf.nn.batch_normalization</code>     |   |
|                        | <code>tf.nn.bias_add</code>                |   |
|                        | <code>tf.nn.dropout</code>                 | Always treated as <code>tf.identity</code> during inference.  |
|                        | <code>tf.nn.fused_batch_norm</code>        |   |
|                        | <code>tf.nn.leaky_relu</code>              |   |
|                        | <code>tf.nn.relu</code>                    |   |
|                        | <code>tf.nn.relu6</code>                   |   |
|                        | <code>tf.nn.relu_layer</code>              |   |
|                        | <code>tf.nn.softmax</code>                 |   |

## Accelerated graph operators

```
Add
AddN
AddV2
BatchMatMul
BatchMatMulV2
BiasAdd
Cast
Const
Cumsum
Einsum
Erf
Exp
ExpandDims
FusedBatchNorm
FusedBatchNormV2
FusedBatchNormV3
Greater
Identity
LeakyRelu
MatMul
Max
Maximum
Minimum
Mean
Mul
Neg
Pack
RealDiv
Relu
Relu6
Reshape
Rsqrt
Sigmoid
Softmax
Split
SplitV
Sqrt
Square
SquaredDifference
Squeeze
StridedSlice
Sub
Sum
Tanh
Transpose
Unpack
```

The lists share many commonalities with [Available TensorFlow Ops](#). Portions of this page are modifications based on work created and [shared by Google](#) and used according to terms described in the [Creative Commons 4.0 Attribution License](#).

*This document is relevant for:* Inf1

- *TensorFlow 2.x (tensorflow-neuron) Release Notes*
- *TensorFlow 2.x (tensorflow-neuron) Accelerated (torch-neuron) Python APIs and Graph Ops*

This document is relevant for: **Inf1**

Setup (tensorflow-neuron)

## Tutorials (tensorflow-neuron)

### Natural Language Processing (NLP) Tutorials

- TensorFlow 2.x - HuggingFace Pipelines distilBERT with Tensorflow2 Neuron [\[html\]](#) [\[notebook\]](#)

### Utilizing Neuron Capabilities Tutorials

- TensorFlow 2.x - Using NEURON\_RT\_VISIBLE\_CORES with TensorFlow Serving [\[html\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

## Additional Examples (tensorflow-neuron)

- [AWS Neuron Samples GitHub Repository](#)

## API Reference Guide (tensorflow-neuron)

- *TensorFlow 2.x (tensorflow-neuron) Tracing API*
- *TensorFlow 2.x (tensorflow-neuron) analyze\_model API*
- *TensorFlow 2.x (tensorflow-neuron) Auto Multicore Replication (Beta)*

## Misc (tensorflow-neuron)

- *TensorFlow 2.x (tensorflow-neuron) Release Notes*
- *TensorFlow 2.x (tensorflow-neuron) Accelerated (torch-neuron) Python APIs and Graph Ops*

This document is relevant for: **Inf1**

Tensorflow NeuronX for Inference on Inf2 & Trn1 / Trn1n

Tensorflow Neuron for Inference on Inf1

This document is relevant for: **Inf1, Inf2, Trn1**



## NEURONX DISTRIBUTED (NXD)

*This document is relevant for: Trn1, Trn2*

### 3.1 NxD Training

*This document is relevant for: Trn1, Trn2*

#### 3.1.1 Overview

##### Table of contents

- *NxD Training*
- *Using NxD Training*
  - *Configuration File*
  - *PyTorch Lightning APIs*
  - *NxD Core Primitives*

#### NxD Training

The NeuronX Distributed Training (NxD Training) library is a collection of open-source tools and libraries designed to empower customers to train PyTorch models on AWS Trainium instances. It combines both ease-of-use and access to features built on top of *NxD Core* library. Except for a few Trainium specific features, NxD Training is compatible with training platforms like NVIDIA's NeMo.

Specifically, *NxD Training* offers the following features and productivity flows:

- **Training Workflows:** Developers benefit from turnkey support for multiple workflows such as model Pre-training, Supervised Finetuning (SFT), and Parameter Efficient Finetuning (PEFT) using Low Rank Adapters (LoRA)<sup>1</sup>. For these workflows, precision types supported include (a) FP32 for both baseline and for master weights when using ZeRO-1, and (b) BF16 combined with *stochastic rounding*.
- **Distributed Strategies:** Splitting training workload over multiple nodes shortens the job duration. This is made possible through distributed strategies that are the techniques used to shard large scale models across multiple Neuron Cores. NxD Training Distributed Strategies are implemented in the *NxD Core* library and include: Data

---

<sup>1</sup> Supported through NxD Core.

Parallelism, [Tensor-parallelism](#), [Sequence-Parallelism](#), [Pipeline-parallelism](#) (including 1F1B pipeline schedule and interleaved pipeline schedule), and [ZeRO-1](#).

- **Data Science Modules:** The integration of datasets, dataloaders, tokenizers and other data wrangling tools makes it easy to prepare and use large-scale training data.
- **Data Engineering Modules:** Integrated *Experiment Manager* allows for saving training outputs through checkpointing and evaluating results through enhanced logging. It comes with multiple options for optimally loading/saving checkpoints such as sharded checkpoints, last-K checkpoints, asynchronous checkpoints, auto-resume from checkpoints and storage in S3 buckets.
- **PyTorch Lightning:** NxD Training is integrated with training frameworks like like PyTorch Lightning that help with organizing training code.
- **Models:** Users can start on NxD Training with ready-to-use samples based on HuggingFace and Megatron-LM model formats. It has support for advanced LLM architecture blocks such as [Grouped Query Attention layer](#).
- **SW Releases:** NxD Training code is available on [GitHub](#), both as pip wheel and source code.

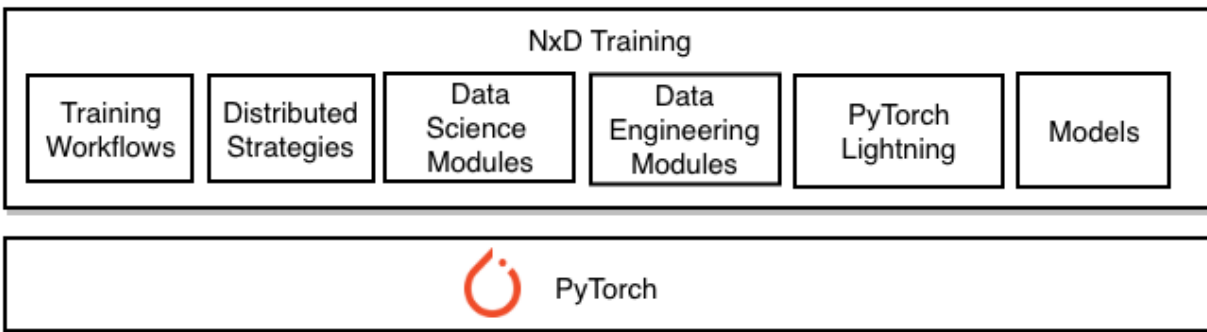


Fig. 3.1: *NxD Training*

## Using NxD Training

ML developers often need access to training code at different levels of abstraction. As shown in [figure](#), using NxD Training is possible using three interfaces:

- High-level [YAML](#) configuration file used in conjunction with models in NxD Training's model hub
- [PyTorch Lightning \(PTL\)](#) APIs and Trainer in conjunction with NxD Core primitives
- [NxD Core](#) foundational API, also referred to as NxD Core primitives

All three usage mechanisms employ the underlying NxD Core library either directly through programming interfaces or configuration files and developers can choose the method that meets their needs.

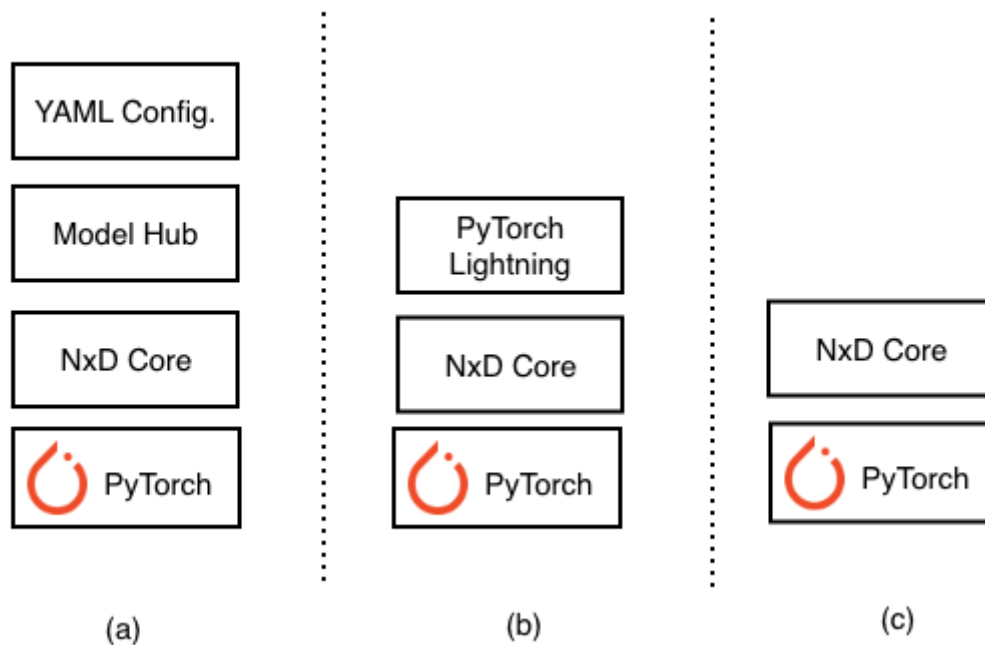


Fig. 3.2: Using Nx/D Training through (a) Configuration Files (b) PyTorch Lightning APIs, and (c) Nx/D Core primitives

### Configuration File

Nx/D Training supports a top-level access for distributed training using YAML based configuration files. This option is available for models that are available in the model hub or custom ones enabled after following the steps listed in [model integration guide](#) inside Nx/D Training. With this usage model, only the configuration parameters inside the YAML file need to be set and no further code changes are necessary. This facilitates easy experimentation with various configuration settings and automating the workflow. Figure below shows the major settings available inside YAML configuration file and more details on how to exercise these options are in [YAML Configuration Settings](#). Existing users of NeuronX NeMo Megatron (NNM) or NVIDIA NeMo can review [NNM](#) and [NeMo](#) migration guides, respectively, to map the configuration parameters to Nx/D Training.

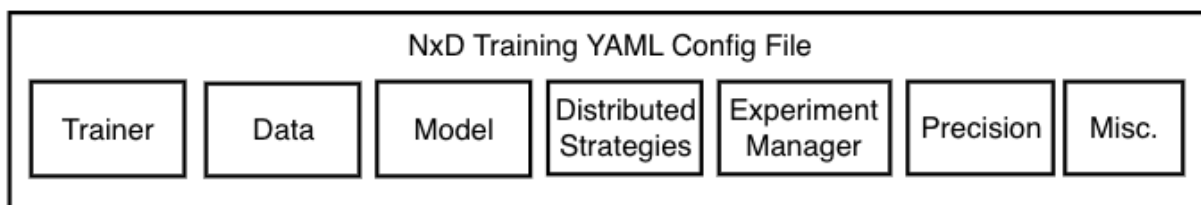


Fig. 3.3: Top level settings for Nx/D Training through configuration file

## PyTorch Lightning APIs

**PyTorch Lightning** is a library that abstracts out model training workflows and eliminates the boilerplate code to setup training loops. Through its inheritable classes for training loops, data and customizable callbacks for checkpointing and distributed strategies, developers can set training workflows in a standardized and compact manner.

As shown in *user interfaces to NxD Training, Figure (b)*, overall training scripts can be built using PyTorch Lightning and making use of NxD Core library. This requires overriding the base classes of PyTorch Lightning such as `LightningModule`, `DataModule`; configuring optimizer and LR scheduler; setting appropriate callbacks; and launching the Trainer. For more details, refer to NxD Core's PyTorch Lightning *developer guide* and *sample tutorial*.

## NxD Core Primitives

NxD Core primitives are basic APIs that can be stitched together to build complete training workflows for AWS Trainium instances. Additionally, these primitives are required for integrating a new custom model into NxD Training or using the model directly via NxD Core library.

NxD Core library has support for all the essential training features - model sharding, handling collective communications, memory reduction, checkpointing, optimizer setting and profiling. For example, tensor parallelism through NxD Core is achieved by converting the linear layers, common in attention modules of transformer-architecture based models, to parallel layers. For pipeline parallelism, NxD Core offers ability for both manual and automatic selection of pipeline cut points in the model graph. Additional options for sequence parallelism and activation recomputation help with memory reduction. For all these parallelism options, NxD Core library automatically ensures efficient management of all the required collective communications across Neuron Cores.

Exact details on how these capabilities can be exercised are described in *NxD Core developer guide*. For background information and description of NxD Core primitives, users are referred to NxD Core's *app notes*, and *API guide*, respectively. Following these steps, once a new model is onboarded using NxD Core APIs, its training workflow can be streamlined using NxD Training's experiment manager and data science/engineering modules.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### 3.1.2 Setup

Neuronx Distributed Training framework is built on top of **NeuronxDistributed (NxD)**, **NeMo** libraries and **PyTorch-Lightning**. The guide below will provide a step-by-step instructions on how to setup the environment to run training using NeuronX Distributed Training framework. Alternatively, you can use the Neuronx Distributed Training virtual environment found in the Neuron DLAMI without running any of these setup steps. See *Neuron DLAMI User Guide*.

#### Table of contents

- *Setup a python Virtual Environment*
- *Installing Neuron Dependencies*
- *Building Apex*
- *Installing the requirements*
- *Installing Neuronx Distributed Training framework*
- *Common failures during installation*



## Setup a python Virtual Environment

Let's first setup a virtual env for our development. This can be done using the command below:

```
python3 -m venv env
source env/bin/activate
```

## Installing Neuron Dependencies

Install the neuron packages using the command:

```
pip install -U pip
pip install --upgrade neuronx-cc==2.* torch-neuronx torchvision neuronx_distributed --
↪extra-index-url https://pip.repos.neuron.amazonaws.com
```

## Building Apex

NxD Training uses the NeMo toolkit, which requires you to install additional dependencies. One of these dependencies is the [Apex](#) library. The NeMo toolkit uses this library for several fused module implementations.

**Note:** NeMo used to use Apex for all distributed training APIs. Since we are using NxD for the same purpose, the use of Apex for this framework is very minimal. It's been added as a dependency since some of the minor imports inside NeMo will break without it. Hence, when building Apex, we build a slim CPU version using the instructions below:

1. Clone Apex repo

```
git clone https://github.com/NVIDIA/apex.git
cd apex
git checkout 23.05
```

2. Replace the contents of the `setup.py` with the following contents:

```
import sys
import warnings
import os
from packaging.version import parse, Version

from setuptools import setup, find_packages
import subprocess

import torch
from torch.utils.cpp_extension import BuildExtension, CppExtension, CUDAExtension, CUDA_
↪HOME, load

setup(
    name="apex",
    version="0.1",
    packages=find_packages(
        exclude=("build", "csrc", "include", "tests", "dist", "docs", "tests", "examples
↪", "apex.egg-info",)
```

(continues on next page)

(continued from previous page)

```

    ),
    install_requires=["packaging>20.6",],
    description="PyTorch Extensions written by NVIDIA",
)

```

3. Install python dependencies:

```
pip install packaging wheel
```

4. Build the wheel using the command:

```
python setup.py bdist_wheel
```

5. After this, you should see the wheel at `dist/`. You can use this for installation in the next section.

6. Come out of the apex directory using `cd ...`

## Installing the requirements

Download the `requirements.txt` using the command:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed-training/master/
↪requirements.txt
```

We can now install the dependencies of the library using the following command:

```
pip install -r requirements.txt ~/apex/dist/apex-0.1-py3-none-any.whl
```

After installing the requirements, we need to patch some of the installations so run

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed-training/master/
↪install_setup.sh
chmod +x install_setup.sh
./install_setup.sh
```

You may see some warnings related to the installations, but those can be ignored.

## Installing Neuronx Distributed Training framework

To install the library, one can run the following command:

```
pip install neuronx_distributed_training --extra-index-url https://pip.repos.neuron.
↪amazonaws.com
```

## Common failures during installation

This section goes over the common failures one can see during setup and how to resolve them.

### 1. ``ModuleNotFoundError: No module named 'Cython'``

You may have to install Cython explicitly using `pip install Cython`

### 2. Error while building ``youtokentome``

If you get an error that says `Python.h file not found`, you may have to install `python-dev` and recreate the virtual env. To install `python-dev`, you can use the command: `sudo apt-get install python-dev`

### 3. Mismatched torch and torch-xla version

When you see an error that looks like:

```
ImportError: env/lib/python3.10/site-packages/_XLAC.cpython-310-x86_64-linux-gnu.so:
↳ undefined symbol: _ZN3c109TupleTypeC1EST6vectorINS_4Type24SingletonOrSharedTypePtrIS2_
↳ EESaIS4_EENS_8optionalINS_13QualifiedNameEEEESt10shared_ptrINS_14FunctionSchemaEE
```

It indicates that the major versions of torch **and** torch-xla don't match.

**Note:** If you install torch again, make sure to install the corresponding torchvision version else that would have a conflict.

### 4. Torch vision version error

The below error indicates incorrect torchvision version. If installing `torch=2.1`, install `torchvision=0.16` (This [link](#) shows which version of torchvision is compatible with which version of torch).

```
ValueError: Could not find the operator torchvision::nms. Please make sure you have
↳ already registered the operator
and (if registered from C++) loaded it via torch.ops.load_library.`
```

### 5. Matplotlib lock error

If you see the below error:

```
TimeoutError: Lock error: Matplotlib failed to acquire the following lock file

This error means there is some contention in compute/worker nodes to access the
↳ matplotlib cache, and hence the timeout
error. To resolve this error, add or run ``python -c 'import matplotlib.pyplot as plt'``
↳ command as part of your setup.
This will create a matplotlib cache and avoid the race condition.
```

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### 3.1.3 App Notes

*This document is relevant for: Trn1, Trn2*

#### Introducing NxD Training

##### Table of contents

- *What are we introducing?*
- *I currently use NeuronX Distributed (NxD Core). How does NxD Training release affect me?*
- *Should the current Neuron NeMo Megatron (NNM) users continue to use NNM?*
- *I am new to Neuron and have training workloads, what toolkits or libraries should I use?*
- *Additional Resources*

#### What are we introducing?

Starting with the Neuron 2.20 release, we are introducing NxD Training. In doing so, we are expanding NeuronX Distributed library (previously called NxD that will now be called NxD Core) to NxD Training with data science/engineering modules, and end to end examples. NxD Training is a PyTorch based distributed training library that enables customers to train large-scale models. Some key distributed strategies supported by NxD Training include 3D-parallelism (data parallelism, tensor parallelism and pipeline parallelism) and ZeRO-1 (where optimizer states are partitioned across workers).

NxD Training supports model training workflows like pretraining, supervised finetuning (SFT) and parameter efficient finetuning (PEFT) using Low-Rank Adapter (LoRA) techniques<sup>1</sup>. For developers, NxD Training offers both API level access through NxD Core and PyTorch Lightning and an intuitive interface via YAML based configuration files. NxD Training offers a flexible approach that enables customers to leverage only the functionalities that align with their unique workflows and seamlessly integrate their machine learning training software at the appropriate level within NxD Training, ensuring a user experience tailored to their specific requirements. This is a beta preview version of NxD Training and feedback from the developer community is strongly encouraged for upcoming releases.

#### I currently use NeuronX Distributed (NxD Core). How does NxD Training release affect me?

Existing NxD Core customers can continue to use NxD Core APIs available under NxD Training. If workflows based on NxD Core meet your needs, you do not need to do anything different with NxD Training's introduction. NxD Core APIs and functionalities for NxD Core continue to be available to you as before. You can choose to [install NxD Core only](#) and skip all subsequent installation steps for NxD Training. However, NxD Training has additional support for YAML based configuration, a model hub and integration with PyTorch Lightning. If these capabilities are of interest to you, you may choose to evaluate and start using NxD Training.

<sup>1</sup> Supported through NxD Core.

## Should the current Neuron NeMo Megatron (NNM) users continue to use NNM?

NxD Training offers same capabilities as Neuron NeMo Megatron (NNM). Additionally, NNM will go into maintenance mode in the next release. If you are currently using NNM, the introduction of NxD Training toolkit means that you should start evaluating NxD Training for your training needs. With its YAML interface, NxD Training is very close in terms of usability to NNM and NeMo. Migrating from NNM to NxD Training should involve a relatively minor effort and instructions for doing so are provided [here](#).

## I am new to Neuron and have training workloads, what toolkits or libraries should I use?

If you are starting with Neuron and looking for solutions to your model pretraining or finetuning needs, then NxD Training is the recommended toolkit for you. Please start from [NxD Training page](#) for overview, installation and usage instructions.

## Additional Resources

Multiple NxD Training resources on getting started, using it and getting required support are listed below. If you encounter issues or have product related questions, please refer to FAQs and troubleshooting guides. Additionally, please feel free to reach out to us using resources in Support section.

[How to get started](#)

[Release notes](#)

[Main section](#)

[Troubleshooting](#)

[Support](#)

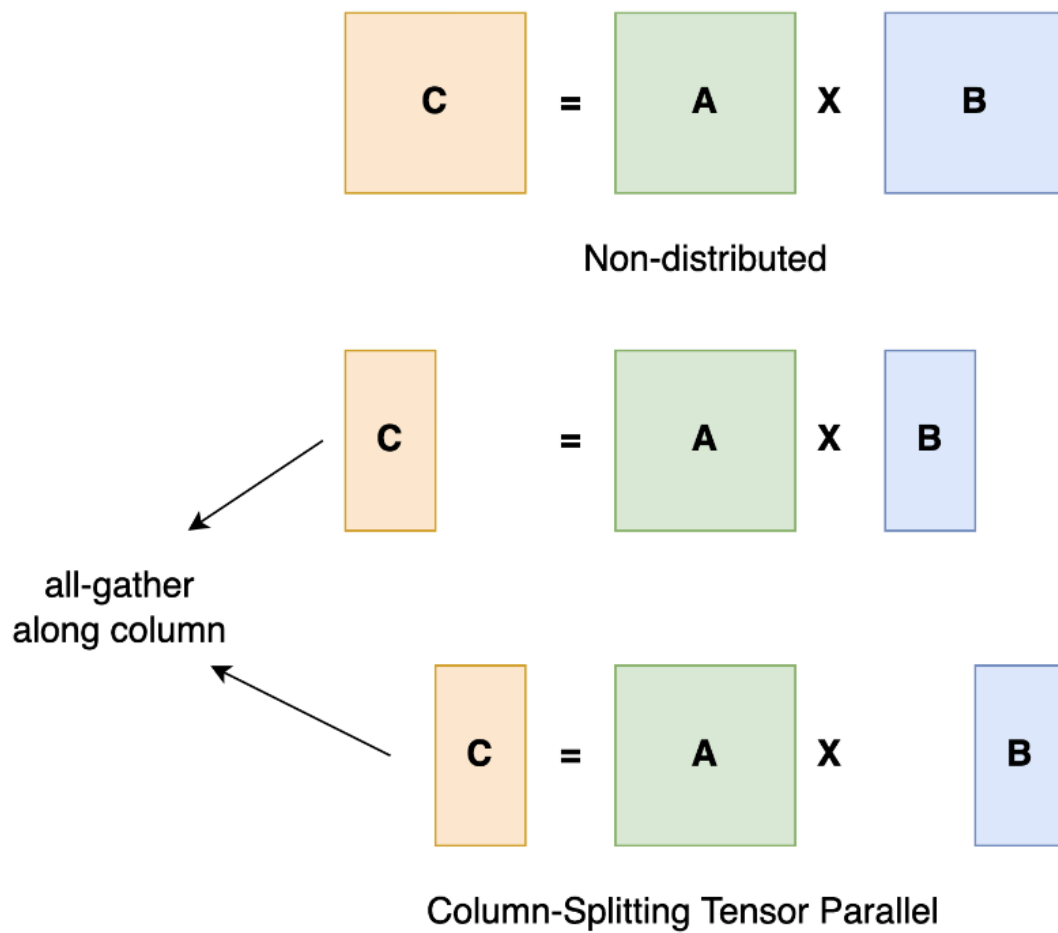
*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Tensor Parallelism Overview

Tensor Parallelism is a technique in which a tensor is split into  $N$  chunks along a particular dimension such that each device only holds  $1/N$  chunk of the tensor. Computation is performed using this partial chunk so as to get partial output. These partial outputs are collected from all devices ensuring the correctness of the computation is maintained.

Taking a general matrix multiplication as an example, let's say we have  $C = AB$ . We can split  $B$  along the column dimension into  $[B_0 B_1 B_2 \dots B_n]$  and each device holds a column. We then multiply  $A$  with each column in  $B$  on each device, we will get  $[AB_0 AB_1 AB_2 \dots AB_n]$ . At this moment, each device still holds partial results, e.g. device rank 0 holds  $AB_0$ . To make sure the result is correct, we need to all-gather the partial result and concatenate the tensor along the column dimension. In this way, we are able to distribute the tensor over devices while making sure the computation flow remains correct.

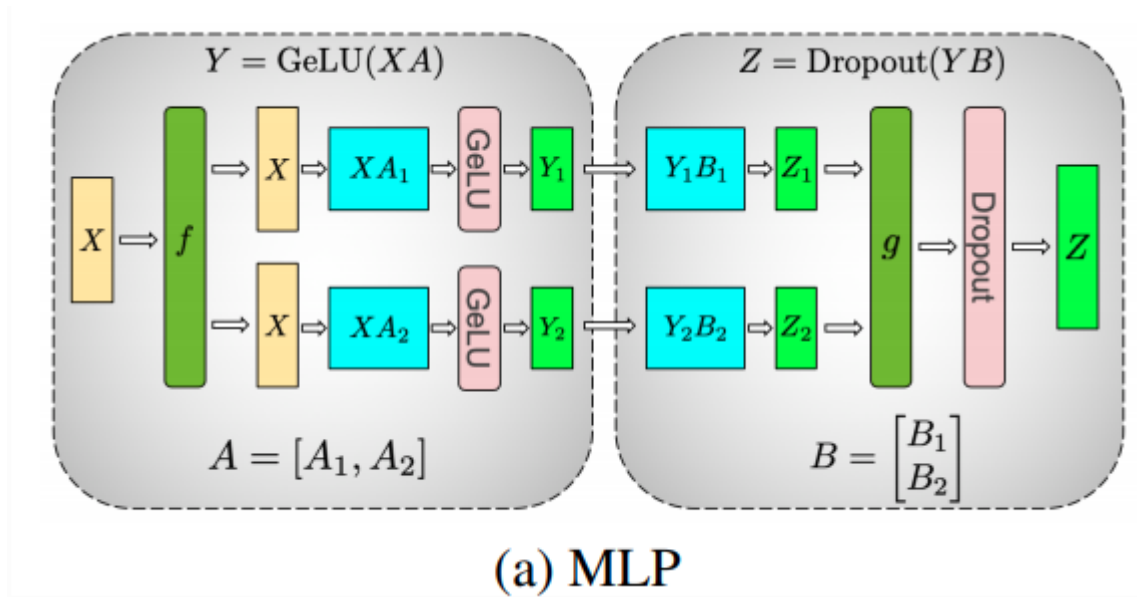


Tensor parallel illustration

Fig and TP explanation is borrowed from [https://colossalai.org/docs/concepts/paradigms\\_of\\_parallelism/#tensor-parallel](https://colossalai.org/docs/concepts/paradigms_of_parallelism/#tensor-parallel)

Similarly we can perform the partition along the row dimensions and create a RowParallel Linear layer. In RowParallelLinear layer, we partition the weight matrix along the row dimension. Let's say we have  $C = AB$ . We can split  $B$  along the row dimension into  $[B_0 B_1 B_2 \dots B_n]$  and each device holds a row. We then multiply each column of  $A$  on each device, we will get  $[A_0B_0 A_1B_1 A_2B_2 \dots A_nB_n]$ . At this moment, each device still holds partial results, e.g. device rank 0 holds  $A_0B_0$ . To make sure the result is correct, we need to all-reduce sum the partial result from all devices to produce the final output.

Using this principle of sharded linear layers, we can construct MLPs of arbitrary depth until the need to operate on the whole output tensor, in which case we would have to construct the output but gathering it from all devices.



Here is an illustration from the Megatron-LM paper. In the above case, as you can see two linear layers are implemented using Column Parallel and Row Parallel linear layers, wherein the ColumnParallel Linear shards along the columns and then it is followed by RowParallel Linear layer which takes in parallel inputs (sharded outputs from ColumnParallelLinear). Consider the example shown in the above diagram,  $Z = (XA)B$ . In this case we split the first matrix multiplication over column dimension such that each device after first matrix multiplication holds partial result of  $Y_0=XA_0, Y_1=XA_1$  and so on. For the second matrix multiplication, we partition the weight matrix over row dimension and since the inputs are already columns sharded and we can multiply them to produce partial outputs. These outputs finally requires an all-reduce sum, since we want to sum up the single column\*row result.

Tensor Parallelism for Transformers:

A transformer block

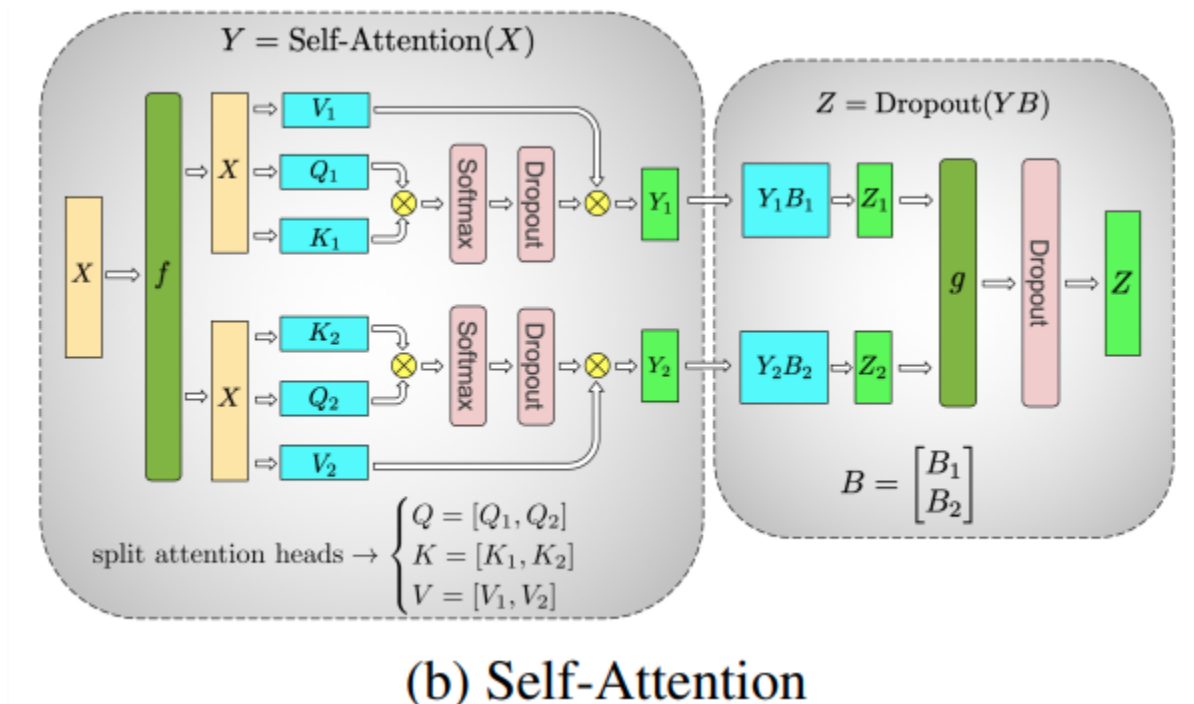


Fig: Taken from Megatron-LM paper.

As seen from the figure above, a simple self attention block has the QKV linear layer followed by MLP. Using the same Column and Row Parallel linear layers, we can partition the self-attention block across devices thereby reducing the memory footprint on each device, since each device now only holds partial parameters. This weight distribution strategy allows us to scale large model training across devices.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Pipeline Parallelism Overview

Pipeline parallelism is a technique used in deep learning model training to improve efficiency and reduce the training time of large neural networks. Currently NeuronxDistributed's pipeline parallelism is built specially for transformer based models, where each Neuron core will be assigned with a subset of transformer layers. Pipelining is a technique to achieve true parallelization in pipeline parallelism, by having the Neuron cores compute simultaneously on different data samples, and to overcome the performance loss due to sequential computation. When you use pipeline parallelism, training job is executed in a pipelined fashion over microbatches to maximize device usage.

## Model partitioning

In NeuronxDistributed, we use [Pytorch's FX](#) to trace the model and do partition on the FX IR. User simply needs to specify where to cut the pipeline stages, and our algorithm will cut the pipeline stages and assign the corresponding modules to each Neuron core automatically. Currently we require user to provide model partition decision but auto-partition will be supported in the future. Here is an example of simple partition with 5 linear layers

```
# original NN module
class MyModule(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.linears = torch.nn.ModuleList([torch.nn.Linear(4, 4) for _ in range(5)])

    def forward(self, x):
        for lin in self.linears:
            x = lin(x)
        return x

m = MyModule()
gm = torch.fx.symbolic_trace(m)
print(gm)
"""
GraphModule(
  (linears): Module(
    (0): Linear(in_features=4, out_features=4, bias=True)
    (1): Linear(in_features=4, out_features=4, bias=True)
    (2): Linear(in_features=4, out_features=4, bias=True)
    (3): Linear(in_features=4, out_features=4, bias=True)
    (4): Linear(in_features=4, out_features=4, bias=True)
  )
)

def forward(self, x):
```

(continues on next page)



(continued from previous page)

```

linear_0 = getattr(self.linears, "0")(x); x = None
linear_1 = getattr(self.linears, "1")(linear_0); linear_0 = None
linear_2 = getattr(self.linears, "2")(linear_1); linear_1 = None
linear_3 = getattr(self.linears, "3")(linear_2); linear_2 = None
linear_4 = getattr(self.linears, "4")(linear_3); linear_3 = None
return linear_4
"""

```

If user decide to cut the pipeline stage at the 3rd linear call, after partition there will be 2 submodules, where *submod\_0* contains first 3 linear layers and *submod\_1* contains last 2 linear layers.

```

After Split module
GraphModule(
  (submod_0): GraphModule(
    (linear_0): Linear(in_features=4, out_features=4, bias=True)
    (linear_1): Linear(in_features=4, out_features=4, bias=True)
    (linear_2): Linear(in_features=4, out_features=4, bias=True)
  )
  (submod_1): GraphModule(
    (linear_3): Linear(in_features=4, out_features=4, bias=True)
    (linear_4): Linear(in_features=4, out_features=4, bias=True)
  )
)

def forward(self, x):
    submod_0 = self.submod_0(x); x = None
    submod_1 = self.submod_1(submod_0); submod_0 = None
    return submod_1

```

## Pipeline Execution Schedule

Pipelining is based on splitting a mini-batch into microbatches, which are fed into the training pipeline one-by-one and follow an execution schedule defined by the library runtime. A microbatch is a smaller subset of a given training mini-batch. The pipeline schedule determines which microbatch is executed by which device for every time slot.

For example, depending on the pipeline schedule and the model partition, Neuron core *i* might perform (forward or backward) computation on microbatch *b* while Neuron core *i+1* performs computation on microbatch *b+1*, thereby keeping both Neuron cores active at the same time. An example taken from Megatron paper is showed as below



*This document is relevant for:* Trn1, Trn2

*This document is relevant for:* Trn1, Trn2

## Activation Memory Reduction

There are three major contributors to high device memory utilization: *Parameters*, *Optimizer states* and *Activation Memory*. To reduce the size of parameter/optimizer states memory, one can use parallelism techniques like *Tensor-parallelism*, *Pipeline-parallelism* or *Zero1*. However, as the hidden size and sequence length increases, the size of the activation memory keeps growing linearly with hidden size and quadratically with sequence length.

The total activation memory without any parallelism comes to about:

$$\text{Activations memory per layer} = sbh \left( 34 + \frac{5as}{h} \right)$$

where,

- $a$ : Number of attention heads
- $b$ : microbatch size
- $h$ : hidden dimension size
- $s$ : sequence length

When we use tensor-parallelism, it not only helps to reduce the parameter and optimizer states size on each device, but it also helps to reduce the activation memory. For a transformer model, where we apply the tensor-parallel sharding on the attention block (more info [here](#)), the activation memory within the attention block also drops by a factor of tensor-parallel degree ( $t$ ). However, since the layernorms and dropouts (which are outside these attention blocks) are not parallelised and they are replicated on each device. These layernorms and dropouts are computationally inexpensive, however, they increase the overall activation memory on each device. Moreover, since we only parallelize within the attention block or within the MLP block ( $h \rightarrow 4h$  projection), the inputs to the QKV multiplies and the MLP are still unsharded. This overall adds to about  $10sbh$  of total activation memory. To reduce this activation memory, one can use 2 methods:

- [Sequence-Parallelism](#)
- [Activation Recomputation](#)

## Sequence Parallelism

Sequence-Parallelism was proposed by [Shenggui and et.al](#). The authors propose to parallelize the compute along all the sequence dimension in an attempt to reduce increasing the memory pressure due to high sequence-lengths. Sequence-parallelism can be combined with tensor-parallelism to reduce the activation memory pressure due to increasing sequence-lengths.

Tensor-parallelism parallelizes the parts of the transformer which are computationally heavy, however, it leaves the layer-norms, dropouts and some MLP block intact. The activation memory for this block adds up to a factor of  $10sbh$ . [Vijay Korthikanti et.al](#) noticed that the compute in the non-tensor parallel region is independent in the sequence dimension. This property can be leveraged to shard the compute along the sequence dimension. The main advantage of sharding these non-tensor parallel block is reducing the activation memory. We can use the same tensor-parallel degree to partition, thereby reducing the overall activation memory by a factor of  $t$ . However, this partitioning comes at a cost. Since we are partitioning the non-tensor parallel region along sequence dimension, we have to collect the activations before we feed to the tensor-parallel block. This requires an introduction of [all-gather](#) collective operation which would gather the activations along the sequence dimension. Similarly, after the tensor-parallel block, since we would have to split the activations along the sequence dimension and distribute among the devices. Since, the tensor-parallel block in the transformer module already uses an all-reduce (Row-parallel linear layer used for MLP), we can replace the all-reduce operation with a [reduce-scatter](#) operation.

libraries/nxd-training/app\_notes/images/sequence\_parallel.png

Ref: [Reducing Activation Recomputation in Large Transformer Models](#)

In the figure,  $g$  is all-gather operation and  $g^{-1}$  is the reduce-scatter operation.  $g$  and  $g^{-1}$  are conjugates and in the backward pass,  $g^{-1}$  becomes an all-gather operation and  $g$  becomes the reduce-scatter operation. At first glance, it appears that sequence-parallelism when combined with tensor-parallelism introduces an extra communication operation, however, in a ring all-reduce, the op is broken down into all-gather and reduce-scatter. Hence, the bandwidth required for sequence-parallelism is same as tensor-parallelism only. Hence, we are not losing out on compute but end up saving the activation memory per device. The final activation memory when sequence-parallelism is combined with tensor-parallelism:

$$\text{Activations memory per layer} = sbh \left( \frac{10}{t} + \frac{24}{t} + \frac{5as}{ht} \right) = \frac{sbh}{t} \left( 34 + \frac{5as}{h} \right)$$

## Activation Recomputation

The total required memory in the above equation can still be high as we increase the sequence length and hidden size. We would have to keep increasing the tensor-parallel degree to accommodate this requirement. Increasing the tensor-parallel degree might soon start producing diminishing returns as the model would start becoming bandwidth bottlenecked because of the extra collective communication operations. [Activation recomputation](#) can help to alleviate this problem. In this method, we recompute a part of the forward pass during the backward pass, thereby avoiding the need to save the activations during the forward pass. Activation recomputation is a trade-off between duplicate computation vs memory. It allows you to save on memory at the cost of extra recompute. This trade-off becomes valuable when we can fit larger models at the expense of recomputing forward pass activations.

Ideally one can recompute the entire forward pass, thereby resulting in an activation memory of  $2sbh$  per transformer layer. This method is called *Full-activation checkpointing*. This memory can further go down by a factor of  $t$  if we use tensor-parallelism. In the activation memory equation, we have a quadratic term of  $5abs^2$ . As the sequence length, this term will grow at a much faster rate. This quadratic term comes from the softmax computation. [Vijay Korthikanti et.al](#) propose *Selective activation checkpointing* where they only recompute the softmax and attention computation and thereby avoid saving the activations coming from softmax and attention computation. This completely gets rid of the quadratic term and brings down the activation memory per layer to  $34sbh/t$ . The LLama-7B example in [this tutorial](#) used selective activation checkpointing.

This document is relevant for: Trn1, Trn2

- [Introducing NxD Training](#)
- [Tensor Parallelism Overview](#)
- [Pipeline Parallelism Overview](#)
- [Activation Memory Reduction](#)

This document is relevant for: Trn1, Trn2

This document is relevant for: Trn1, Trn2

### 3.1.4 API Reference Guide

*This document is relevant for: Trn1, Trn2*

#### YAML Configuration Settings

The library allows configuring a bunch of parameters in the YAML file to run large scale training. The important categories and parameters are highlighted below. At the top level, we have the following keys:

```
name:
    # Name of the experiment
model_source:
    # Model source code, could be megatron or hf
seed:
    # Random seed to be used for the entire experiment
trainer:
    # Settings to configure the PyTorch-Lightning trainer
exp_manager:
    # Settings to configure logging/checkpointing
distributed_strategy:
    # Settings to configure how the model is to be distributed across devices
data:
    # Settings to configure the dataset/dataloader
model:
    # Settings to configure the model architecture and the optimizer
precision:
    # Settings to configure the model precision
compiler_flags:
    # Neuron compiler flags to be used
compiler_cache_url:
    # Cache to be used to save the compiled artifacts
async_exec_max_inflight_requests:
    # Used to configure the runtime queue
bucket_size_collectives:
    # Collectives are batched into tensors of this size (in MBs)
neuron_rt_exec_timeout:
    # Runtime timeout
neuron_experimental_compress_rg:
    # To use compress replica group
```

#### Trainer

Neuronx Distributed Trainer framework is built on top of `PyTorch-Lightning` and this key allows users to configure the trainer.

```
devices: 32
num_nodes: 1
max_epochs: -1
max_steps: 20000
log_every_n_steps: 1
val_check_interval: 20000
```

(continues on next page)

(continued from previous page)

```

check_val_every_n_epoch: null
num_sanity_val_steps: 0
limit_val_batches: 1
limit_test_batches: 1
gradient_clip_val: 1.0
lnc: 2
sequential_move_factor: 11

```

**Note:** All the above trainer parameters follow the exact same definition of the PyTorch-Lightning Trainer. More information about each of them can be found [here](#).

### devices

Number of devices to be used for training. If using torchrun, this is equal to `nproc_per_node * num_nodes`.

- **Type:** integer
- **Required:** True

### lnc

Neuron-specific setting that specifies the logical-to-physical Neuron Core mapping ratio. This parameter determines the number of physical Neuron cores used for each logical Neuron Core.

Values:

- lnc: 1 - Each node exposes 128 logical devices, with a 1:1 mapping between logical and physical Neuron Cores.
- lnc: 2 - Implements a 2:1 mapping between logical and physical Neuron Cores.
  - **Type:** integer
  - **Required:** False
  - **Default:** None (must be explicitly set)

### num\_nodes

Number of nodes to be used for training

- **Type:** integer
- **Required:** True

### max\_epochs

Maximum number of epochs to run. A value of -1 means that the number of training steps would be inferred from `max_steps`

- **Type:** integer
- **Required:** True

### log\_every\_n\_steps

How often to log loss values

- **Default value:** 1
- **Type:** integer
- **Required:** True

**val\_check\_interval**

How often to run validation step. Using this parameter one can run validation step after X training steps.

- **Type:** integer
- **Required:** True

**check\_val\_every\_n\_epoch**

Another parameter that controls the frequency of validation step. Using this parameter, one can run validation step after X epochs.

- **Type:** integer
- **Required:** True

**num\_sanity\_val\_steps**

How many sanity validation steps to run. Keeping it to 0 would not run validation step at the start of training.

- **Type:** integer
- **Required:** True

**limit\_val\_batches**

Number of batches to run validation step on.

- **Type:** integer
- **Required:** True

**gradient\_clip\_val**

Float value to clip gradients at.

- **Type:** float
- **Required:** True

**sequential\_move\_factor**

Number of ranks/devices participating in initializing the model weights in parallel. Useful to reduce init time when using TP-PP config. The value can be increased upto the number of `trainer.devices` being used.

- **Default value:** 11
- **Type:** integer
- **Required:** False

**Experiment Manager**

This setting is mainly for configuring different aspects of experiment management like checkpointing, experiment logging directory, which parameters to log and how often to log, etc.

```
log_local_rank_0_only: True
create_tensorboard_logger: True
explicit_log_dir: null
exp_dir: null
name: megatron_llama
resume_if_exists: True
resume_ignore_no_checkpoint: True
```

(continues on next page)

(continued from previous page)

```

create_checkpoint_callback: True
checkpoint_callback_params:
    monitor: step
    save_top_k: 1
    mode: max
    save_last: False
    filename: 'megatron_llama--{step}-{consumed_samples}'
    every_n_train_steps: 200
    use_master_weights_in_ckpt: False
log_parameter_norm: True
log_gradient_norm: True
enable_recovery_time_instrumentation: False
save_xser: True
load_xser: True
async_checkpointing: False
resume_from_checkpoint: null

```

**log\_local\_rank\_0\_only**

Log only on rank 0. The recommended setting should be True

- **Type:** bool
- **Default:** False
- **Required:** False

**create\_tensorboard\_logger**

Setting this True would log the loss and other parameters to tensorboard.

- **Type:** bool
- **Default:** False
- **Required:** False

**exp\_log\_dir**

Explicitly specify the logging directory. Otherwise, the framework would save to current directory as default.

- **Type:** str
- **Default:** null
- **Required:** False

**resume\_if\_exists**

Set this to True to resume from an existing checkpoint. This config will be useful when we want to auto-resume from a failed training job.

- **Type:** bool
- **Default:** False
- **Required:** False

**resume\_ignore\_no\_checkpoint**

Experiment manager errors out if `resume_if_exists` is True and no checkpoint could be found. This behaviour can be disabled, in which case `exp_manager` will print a message and continue without restoring, by setting `resume_ignore_no_checkpoint` to True.

- **Type:** bool
- **Default:** False
- **Required:** False

**checkpoint\_callback\_params.save\_top\_k**

How many checkpoints to keep around. Example: If set to 1, only 1 checkpoint at any given time would be kept around. The framework would automatically keep deleting checkpoints.

- **Type:** int
- **Required:** True

**checkpoint\_callback\_params.every\_n\_train\_steps**

How often we want to checkpoint.

- **Type:** int
- **Required:** True

**checkpoint\_callback\_params.use\_master\_weights\_in\_ckpt**

Whether or not to save master weights when checkpointing.

- **Type:** bool
- **Default:** False
- **Required:** False

**log\_parameter\_norm**

Set this to log parameter norm across model parallel ranks.

- **Type:** bool
- **Default:** False
- **Required:** False

**log\_gradient\_norm**

Set this to log gradient norm across model parallel ranks.

- **Type:** bool
- **Default:** False
- **Required:** False

**enable\_recovery\_time\_instrumentation**

Set this if you don't want to default to not printing the detailing timing for recovery.

- **Type:** bool
- **Default:** False
- **Required:** False

**save\_xser**

Set this to save with torch xla serialization to reduce time saving, it's recommended to enable `xser` for significantly faster save/load. Note that if the checkpoint is saved with `xser`, it can only be loaded with `xser`, vice versa.

- **Type:** bool
- **Default:** False



- **Required:** False

#### **load\_xser**

Set this to load with torch xla serialization to reduce time saving, it's recommended to enable `xser` for significantly faster save/load. Note that if the checkpoint is saved with `xser`, it can only be loaded with `xser`, vice versa.

- **Type:** bool
- **Default:** False
- **Required:** False

#### **async\_checkpointing**

Set this if you want to use async checkpointing. Under the hood the library uses the async checkpointing feature provided by NeuronxDistributed's [save API](#).

- **Type:** bool
- **Default:** False
- **Required:** False

#### **resume\_from\_checkpoint**

Set this as the checkpoint file to load from. Check the SFT/DPO/ORPO example config under `conf` on how to use it.

- **Type:** str
- **Default:** null
- **Required:** False

#### **ckpt\_ptl\_version**

Set this only if your checkpoint does not contain the pytorch-lightning version in it. This version is the pytorch-lightning version the checkpoint was saved with.

- **Type:** str
- **Default:** "2.5.0"
- **Required:** False

### **Distributed Strategy**

```
tensor_model_parallel_size: 8
pipeline_model_parallel_size: 1
virtual_pipeline_model_parallel_size: 1
zero1: True
sequence_parallel: True
kv_replicator: 4
```

This setting allows users to configure the sharding strategy to be used for distributing the model across workers.

#### **tensor\_model\_parallel\_size**

Tensor `parallel degree` to be used for sharding models.

- **Type:** int
- **Required:** True

**pipeline\_model\_parallel\_size**

Pipeline parallel degree to be used for sharding models.

- **Type:** int
- **Required:** True

**virtual\_pipeline\_model\_parallel\_size**

Interleaved pipeline parallel degree. Use a value of 1 if no pipeline parallelism is used.

- **Type:** int
- **Required:** True

**context\_parallel\_size**

Context parallel degree to be used for sharding sequence. When context\_parallel\_size is greater than 1, fusions.ring\_attention must be set to True.

- **Type:** int
- **Required:** False
- **Default:** 1

**zero1**

Wraps the optimizer with zero1.

- **Type:** bool
- **Required:** True

**sequence\_parallel**

To shard along the sequence dimension. Sequence Parallel is always used in conjunction with tensor parallel. The sequence dimension will be sharded with the same degree as the tensor\_model\_parallel\_size.

- **Type:** bool
- **Required:** True

**kv\_replicator**

This parameter is used together with qkv\_linear parameter. It is used to configure the [GQAQKVLinear module](#)

- **Type:** bool
- **Required:** True

**Data**

This is where we configure the dataset/dataloader. This config is dependent on the dataloader/dataset been used. Users can add custom keys in this config and read inside the CustomDataModule using `cfg.data`. Currently the library adds support for 3 kinds of data modules: MegatronDataModule, ModelAlignmentDataModule and HFDataModule. To learn about the config parameters of MegatronDataModule please check the `megatron_llama_7B_config.yaml`, for ModelAlignmentDataModule check the `megatron_llama2_7B_SFT_config.yaml` and for HFDataModule, refer to `hf_llama3_8B_config.yaml`.

The parameters that are common across all the configs are documented below.

```
micro_batch_size: 1
global_batch_size: 1024
```

**micro\_batch\_size**

The batch is distributed across multiple data parallel ranks and within each rank, we accumulate gradients. Micro batch size is the size that is used for each of those gradient calculation steps.

- **Type:** int
- **Required:** True

**global\_batch\_size**

This config along with micro-batchsize decides the gradient accumulation number automatically.

- **Type:** int
- **Required:** True

**Model**

This is where we can configure the model architecture. When building custom models, this config can be used to parameterize the custom model. The below parameters are taken from an example of the Megatron model config. Depending on the model and required parameters, this config can change.

**HF Model**

Let's start with the config for the HF model:

```
# model architecture
model_config: /home/ubuntu/config.json
encoder_seq_length: 4096
max_position_embeddings: ${.encoder_seq_length}
num_layers: 4
hidden_size: 4096
qkv_linear: False

# Miscellaneous
use_cpu_initialization: True

## Activation Checkpointing
activations_checkpoint_granularity: selective
activations_checkpoint_recompute: [CoreAttention]

fusions:
  softmax: True
  flash_attention: False

do_layer_norm_weight_decay: False

optim:
  name: adamw_fp32OptState
  lr: 3e-4
  weight_decay: 0.01
  capturable: False
  betas:
    - 0.9
```

(continues on next page)

(continued from previous page)

```
- 0.999
sched:
  name: LinearAnnealingWithWarmUp
  warmup_steps: 100
  max_steps: ${trainer.max_steps}
```

**model\_config**

Points to the `config.json` path required by the transformers model implementation. One such example of `config.json` is [here](#)

- **Type:** str
- **Required:** True

**encoder\_seq\_length**

Setting the sequence length for the training job. This parameter is common for all models supported in the library.

- **Type:** int
- **Required:** True

**num\_layers**

This config will override the number of layers inside the `config.json` in the `model_config`. This is exposed so that one can quickly increase/decrease the size of the model. This parameter is common for all models supported in the library.

- **Type:** int
- **Required:** True

**hidden\_size**

This config will override the `hidden_size` inside the `config.json` in the `model_config`. This parameter is common for all models supported in the library.

- **Type:** int
- **Required:** True

**qkv\_linear**

This needs to be set if users want to use the [GQAQKVLinear module](#)

- **Type:** bool
- **Required:** True

**fuse\_qkv**

This is set if users want to use fused q, k and v tensors in [GQAQKVLinear module](#). Using `fuse_qkv` can improve throughput. This parameter is True by default.

- **Type:** bool
- **Required:** False

**transpose\_nki\_inputs**

This is set if users want to transpose the inputs to NKI FlashAttention function. To be used only when `fusions.flash_attention` is True. Using `transpose_nki_inputs` with `fusions.flash_attention` can improve throughput. This parameter is True by default for all models, unless used otherwise.

- **Type:** bool

- **Required:** False

### **pipeline\_cuts**

This is set as a list of layer names if users want to specify manual cut points for pipeline parallelism. One example is ['model.layers.10', 'model.layers.20'] in the case of PP=3.

- **Type:** List[str]
- **Required:** False

---

**Note:** When using this param, the number of pipeline cuts should always be `pipeline_model_parallel_size-1`.

---

### **use\_cpu\_initialization**

Setting this flag to True will initialize the weights on CPU and then move to device. It is recommended to set this flag to True. This parameter is common for all models supported in the library.

- **Type:** bool
- **Required:** True

### **activations\_checkpoint\_granularity**

This flag controls which module needs to be recomputed during the backward pass.

Values:

- **selective** - Enables selective recomputation of specified modules in *activations\_checkpoint\_recompute* during the backward pass.
- **full** - Saves activations at layer boundaries and recomputes the entire layer during the backward pass.
- **null** - Disables activation checkpointing.

More information on activation recompute can be found [in this link](#). This parameter is common for all models supported in the library.

- **Type:** str
- **Possible Values:** selective, full, null
- **Required:** True

**activations\_checkpoint\_recompute** This config specifies which modules to recompute when using *selective* activation checkpointing. It accepts a list of module names as strings or *null*.

- **Type:** list[str] or *null*
- **Required:** False

### **fusions.softmax**

Setting this flag to True will replace the `torch.nn.Softmax` with a fused custom Softmax operator. This parameter is common for all models supported in the library.

- **Type:** bool
- **Required:** True

### **fusions.flash\_attention**

Setting this flag to True will insert the flash attention module for both forward and backward. This parameter is common for all models supported in the library.

- **Type:** bool

- **Required:** True

#### **fusions.ring\_attention**

Setting this flag to True will use the ring attention module for both forward and backward. This parameter must be true when `context_parallel_size` is greater than 1.

- **Type:** bool
- **Required:** False

#### **fusions.do\_layer\_norm\_weight\_decay**

Setting this flag to True will add layer norm weight decay. This parameter is common for all models supported in the library.

- **Type:** bool
- **Required:** True

#### **optim**

This is where the optimizers can be set. We can configure the optimizers supported by NeMo. All the optimizers can be configured according to the [parameters specified here](#).

- **Type:** config
- **Possible Values:** adamw, adamw\_fp32OptState, sgd, adam, adadelta, adamax, adagrad, rmsprop, rprop, novograd, adafactor
- **Required:** True

#### **optim.sched**

This is where the LR schedulers can be set. We can configure the schedulers supported by NeMo. All the schedulers can be configured according to the [parameters specified here](#).

- **Type:** config
- **Possible Values:** LinearAnnealingWithWarmUp, CosineAnnealing, WarmupPolicy, WarmupHoldPolicy, SquareAnnealing, NoamAnnealing, WarmupAnnealing, StepLR, rprop, ExponentialLR
- **Required:** True

### **Megatron Model**

The library enables a [megatron transformer](#) model which can be configured from the yaml file. The different available parameters are documented below after the following reference example.

```
# model architecture
encoder_seq_length: 4096
max_position_embeddings: ${.encoder_seq_length}
num_layers: 32
hidden_size: 4096
ffn_hidden_size: 11008
num_attention_heads: 32
num_kv_heads: 32
init_method_std: 0.021
hidden_dropout: 0
```

(continues on next page)

(continued from previous page)

```

attention_dropout: 0
ffn_dropout: 0
apply_query_key_layer_scaling: True
normalization: 'rmsnorm'
layernorm_epsilon: 1e-5
do_layer_norm_weight_decay: False # True means weight decay on all params
make_vocab_size_divisible_by: 8 # Pad the vocab size to be divisible by this value for
↳ computation efficiency.
persist_layer_norm: True # Use of persistent fused layer norm kernel.
share_embeddings_and_output_weights: False # Untie embedding and output layer weights.
position_embedding_type: 'rope' # Position embedding type. Options ['learned_absolute',
↳ 'rope']
rotary_percentage: 1 # If using position_embedding_type=rope, then the per head dim is
↳ multiplied by this.
activation: 'swiglu' # ['swiglu', 'gelu']
has_bias: False
# Miscellaneous
use_cpu_initialization: True

## Activation Checkpointing
activations_checkpoint_granularity: selective # 'selective' or 'full'

fusions:
  softmax: True
  flash_attention: False # Use NKI flash attention

optim:
  name: adamw
  lr: 3e-4
  weight_decay: 0.1
  capturable: True
  betas:
    - 0.9
    - 0.95
  sched:
    name: CosineAnnealing
    warmup_steps: 2000
    constant_steps: 0
    min_lr: 3.0e-5

```

**Note:** For common config, please refer to the HF Model section above.

### ffn\_hidden\_size

Transformer FFN hidden size.

- **Type:** int
- **Required:** True

### num\_attention\_heads

Number of Q attention heads.

- **Type:** int
- **Required:** True

**num\_kv\_heads**

Number of KV heads. This is where we can configure Q and KV differently to create GQA modules.

- **Type:** int
- **Required:** True

**init\_method\_std**

Standard deviation to use when we init layers of the transformer model.

- **Type:** float
- **Required:** True

**hidden\_dropout**

Dropout probability for hidden state transformer.

- **Type:** float
- **Required:** True

**attention\_dropout**

Dropout probability in the attention layer.

- **Type:** float
- **Required:** True

**ffn\_dropout**

Dropout probability in the feed-forward layer.

- **Type:** float
- **Required:** True

**apply\_query\_key\_layer\_scaling**

Scale  $Q * K^T$  by  $(1 / \text{layer-number})$ .

- **Type:** bool
- **Required:** True

**normalization**

Normalization layer to use.

- **Type:** str
- **Possible Values:** rmsnorm, layernorm
- **Required:** True

**layernorm\_epsilon**

Epsilon value for layernorm.

- **Type:** float
- **Required:** True



**share\_embeddings\_and\_output\_weights**

Setting this parameter to True will tie the vocab embedding weight with the final MLP weight.

- **Type:** bool
- **Required:** True

**make\_vocab\_size\_divisible\_by**

So lets say your vocab size is 31999 and you set this value to 4, the framework would pad the vocab-size such that it becomes divisible by 4. In this case the close divisible value is 32K.

- **Type:** int
- **Required:** True

**position\_embedding\_type**

Type of position embedding to be used.

- **Type:** str
- **Possible Values:** learned\_absolute, rope
- **Required:** True

**rotary\_percentage**

If using position\_embedding\_type=rope, then the per head dim is multiplied by this factor.

- **Type:** float
- **Required:** True

**activation**

Users can specify the activation function to be used in the model.

- **Type:** str
- **Possible Values:** swiglu, gelu
- **Required:** True

**has\_bias**

Setting this parameter to True will add bias to each of the linear layers in the model.

- **Type:** bool
- **Required:** True

**Precision**

This config can help to decide the dtype of the model/optimizer.

```
precision:
  type: 'mixed_precision' # ['bf16SR', 'fp32', 'autocast', 'mixed_precision', 'mixed_
↪precisionSR', 'manual']
  # Set the following only if precision type is manual, otherwise they will be_
↪automatically set.
  master_weights: False
  fp32_grad_acc: False
  xla_use_bf16: '@'
```

(continues on next page)

(continued from previous page)

```
xla_downcast_bf16: '0'
neuron_rt_stochastic_rounding_en: '0'
parallel_layers_reduce_dtype: 'bf16'
```

**Note:** Only if the precision type is manual, `master_weights`, `fp32_grad_acc`, `xla_use_bf16`, `xla_downcast_bf16`, `neuron_rt_stochastic_rounding_en` will be picked up from the config. These parameters are for more finer control of precision. It is recommended to use `mixed_precision` config for better accuracy.

## type

### **mixed\_precision**

The `mixed_precision` config uses the `zero1` optimizer. It performs grad accumulation, grad cc, and keeps the master copy of the weights in `fp32`. It also sets the `xla_downcast_bf16` environment variable to 1 and disables stochastic rounding.

### **mixed\_precisionSR**

`mixed_precisionSR` is a superset of the `mixed_precision` config with stochastic rounding enabled.

### **bf16SR**

`bf16SR` config will perform all operations in `bf16` and relies on stochastic rounding feature for accuracy gains.

### **autocast**

`autocast` config will follow the exact same precision strategy followed by `torch.autocast`.

**Note:** Autocast is not supported in this release.

### **manual**

To gain control of the different precision nobs, one can set the precision type to `manual` and control parameters like - `master_weights`, `fp32_grad_acc`, `xla_use_bf16`, `xla_downcast_bf16` and `neuron_rt_stochastic_rounding_en`.

## **parallel\_layers\_reduce\_dtype**

This config will perform reduce collectives (all-reduce and reduce-scatter) within parallel layers in the specified precision. If `fp32` precision type is used, then we implicitly set reduce dtype to `fp32`. Otherwise it will be defaulted to `bf16` in all other cases unless specified.

## **Model Alignment Specific**

You can configure fine-tuning (SFT) or model alignment (DPO/ORPO) through the YAML file, along with parameter-efficient fine-tuning using LoRA.

```
model_alignment_strategy:
  # DPO specific config
  dpo:
    kl_beta: 0.01
    loss_type: sigmoid
    max_prompt_length: 2048
    precompute_ref_log_probs: True
    truncation_mode: keep_start
```

(continues on next page)

(continued from previous page)

```

# Alternatively, you can also use SFT specific config
sft:
    packing: True

# Alternatively, can also use ORPO specific config
orpo:
    beta: 0.01
    max_prompt_length: 2048
    truncation_mode: keep_start

# Parameter-efficient finetuning - LoRA config
peft:
    lora_rank: 16
    lora_alpha: 32
    lora_dropout: 0.05
    lora_bias: "none"
    lora_verbose: True
    target_modules: ["qkv_proj"]

```

**model\_alignment\_strategy**

Set only when using finetuning specific algorithms (SFT, DPO, etc) and related hyperparameters DPO-specific parameters.

**dpo****kl\_beta**

KL-divergence beta to control divergence of policy model from reference model

- **Type:** float
- **Default:** 0.01
- **Required:** True

**loss\_type**

Currently support sigmoid version of optimized DPO loss

- **Type:** str
- **Default:** sigmoid
- **Required:** True

**max\_prompt\_length**

Set maximum length of prompt in the concatenated prompt and (chosen/rejected) response input

- **Type:** integer
- **Required:** True

**precompute\_ref\_log\_probs**

To enable precomputation of reference model log probabilities using pre-fit hook, False is not supported currently

- **Type:** bool

- **Required:** True

**truncation\_mode**

To define how to truncate if size (prompt+response) exceeds seq\_length options: ["keep\_start", "keep\_end"]

- **Type:** str
- **Default:** keep\_start`
- **Required:** True

SFT-specific parameters.

**sft****packing**

Appends multiple records in a single record until seq length supported by model, if false uses pad tokens to reach seq length. Setting it to True increases throughput but might impact accuracy.

- **Type:** bool
- **Default:** False
- **Required:** False

Odds Ratio Preference Optimization (ORPO) specific parameters.

**orpo****beta**

KL-divergence beta to control divergence of policy model from reference model

- **Type:** float
- **Default:** 0.01
- **Required:** True

**max\_prompt\_length**

Set maximum length of prompt in the concatenated prompt and (chosen/rejected) response input

- **Type:** integer
- **Required:** True

**truncation\_mode**

To define how to truncate if size (prompt+response) exceeds seq\_length options: ["keep\_start", "keep\_end"]

- **Type:** str
- **Default:** keep\_start`
- **Required:** True

**peft**

Configuration options for Parameter-Efficient Fine-Tuning (PEFT) methods, specifically LoRA settings.

**lora\_rank**

Rank of LoRA; determines the number of trainable parameters Higher rank allows for more expressive adaptations but increases memory usage

- **Type:** int
- **Default:** 16
- **Required:** True

#### **lora\_alpha**

Scaling factor for LoRA updates; affects the magnitude of LoRA adaptations.

- **Type:** int
- **Default:** 32
- **Required:** True

#### **lora\_dropout**

Dropout rate for LoRA layers to prevent overfitting.

- **Type:** float
- **Default:** 0.05
- **Required:** False

#### **lora\_bias**

Bias type for LoRA. Determines which biases are trainable. Can be 'none', 'all' or 'lora\_only'

- **Type:** str
- **Default:** "none"
- **Required:** False

#### **lora\_verbose**

Enables detailed LoRA-related logging during training.

- **Type:** bool
- **Default:** False
- **Required:** False

#### **target\_modules**

List of model layers to apply [LoRA](#).

- **Type:** list[str]
- **Default:** ["qkv\_proj"] (for Llama)
- **Required:** True

*This document is relevant for: Trn1, Trn2*

- [YAML Configuration Settings](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### 3.1.5 Developer Guide

This section will go over a variety of developer guides to help users get started with the Neuronx Distributed Training library.

*This document is relevant for: Trn1, Trn2*

#### Integrating a New Model

The NeuronX Distributed Training library is a modular framework that allows users to integrate their new modules with the framework while still utilizing the other modules provided by the library. In this section, we showcase how to integrate a new model with the library.

##### Table of contents

- *Model Building (torch.nn.Module)*
- *Model Integration*
  - *Build a Lightning Module*
  - *Plug into training.py*
  - *Create config file*
  - *Launching e2e training*

#### Model Building (torch.nn.Module)

Users can create a torch.nn.Module using the tensor-parallel APIs provided by the [NeuronxDistributed](#) library. Let's take an example of the [GPT-NeoX model built inside NxD examples](#). We can copy the model file and treat it as a new model to onboard using the framework.

---

**Note:** To understand more about how to build models using Tensor-parallel APIs check the [Developer guide here](#).

---

#### Model Integration

Once we have built the model, the next step is to integrate with the training framework. This can be done using the following steps:

##### Build a Lightning Module

Neuronx Distributed Training framework provides a `BaseModelModule` that implements the majority of the training APIs. Users can subclass this base module and implement few APIs that set up the model. Here is an example to setup the GPT-NeoX model example. Create a new file called `new_model_module.py` and add the following content.

```
from transformers import GPTNeoXConfig
import neuronx_distributed as nxd
from neuronx_distributed.parallel_layers.layer_norm import LayerNorm
from neuronx_distributed_training.lightning_modules.model.base import BaseModelModule
```

(continues on next page)

(continued from previous page)

```

from neuronx_distributed_training.utils.model_utils import get_param_groups_by_weight_
↳ decay
from modeling_gpt_neox_nxd import GPTNeoXForCausalLMNxN

class MyNewModel(BaseModelModule):

    def _get_model(self,):
        model_name = "EleutherAI/gpt-neox-20b"
        config = GPTNeoXConfig.from_pretrained(model_name)
        config.use_cache = False
        # Note: We can modify the model by reading parameters from self.config.model.
        # We would have to expose those config in the self.config.model accordingly.
        # Couple of examples are here, where we have exposed num_layers and hidden_size.
        if self.config.model.get('num_layers', -1) != -1:
            config.num_hidden_layers = self.config.model.get('num_layers')
        if self.config.model.get('hidden_size', -1) != -1:
            config.hidden_size = self.config.model.get('hidden_size')
        # This is because the GPT-Neox implementation requires this in the config.
        config.sequence_parallel_enabled = self.config.distributed_strategy.get(
↳ "sequence_parallel", False)
        return GPTNeoXForCausalLMNxN(config)

    def build_model(self):
        # This API is where we build the model object, and return the model.
        # However, in addition to returning the model, users need to
        # configure the nxd config too for pipeline parallelism and
        # activation checkpointing. Here is an example:
        if self.config.model.get("activations_checkpoint_granularity", None) ==
↳ "selective":
            # Here just to showcase how to recompute modules, we are using
            # GPTNeoXMLPNxD, users can add their own custom modules
            self.nxd_config["activation_checkpoint_config"] = GPTNeoXMLPNxD
        elif self.config.model.get("activations_checkpoint_granularity", None) == "full":
            self.nxd_config["activation_checkpoint_config"] = "full"

        # Read more about configuring pipeline parallel config here:
        # https://awsdocs-neuron.readthedocs-hosted.com/en/latest/libraries/neuronx-
↳ distributed/pp_developer_guide.html#pp-developer-guide
        self.nxd_config["pipeline_config"].update(
            {
                "transformer_layer_cls": GPTNeoXLayerNxN,
                "output_loss_value_spec": (True, False),
                "input_names": ["input_ids", "attention_mask", "labels"],
                "leaf_module_cls": [LayerNorm.__name__],
            }
        )
        return nxd.initialize_parallel_model(self.nxd_config, self._get_model)

    def setup_optimizer_param_groups(self):
        # Depending on what weight decay we need, users can configure
        # the params groups accordingly.
        no_decay = ["bias"]

```

(continues on next page)

(continued from previous page)

```

    if self.config.model.get("do_layer_norm_weight_decay", False):
        no_decay.append("LayerNorm")
    self._optimizer_param_groups = get_param_groups_by_weight_decay(self.model, no_
    →decay)

    def init_weights(self,):
        """
        This API is mainly to tell the framework how each layer needs
        to be initialized. This is required because Nx's PP API would
        use this to initialize the layers after model partition.
        Any layer that is unique to the model needs to be added here.
        """
        if isinstance(module, LayerNorm):
            module.weight.data.fill_(1.0)
        # The BaseModelModule already initializes the ColumnParallel, RowParallel
        # ParallelEmbedding layers.
        super().init_weights()

```

### Plug into training.py

Once the new model is created, we can then plug this into the training.py script under examples folder. We can modify the training.py script as follows:

```

...
# Assuming we are using the same DataModule we used for LLama example.
data_module = HFDataModule(cfg, trainer)
from new_model_module import MyNewModel
model = MyNewModel(cfg, trainer)

trainer.fit(model, datamodule=data_module)

```

The rest of the code can remain the same. The trainer will now use the MyNewModel for fetching the model code and run e2e training.

### Create config file

Next we can create a config file under conf to be used for this new model. We can start with a copy of hf\_llama\_7B\_config.yaml. Let's call this config file my\_new\_config.yaml. We can remove the key model.model\_config as we are not using it inside our MyNewModel. We can edit the distributed\_strategy config depending on what we need.

---

**Note:** For the dataset, we are using the same dataset that the llama example is using. To configure a new dataset, please check the [Integrating a new dataset/dataloader](#) section

---



## Launching e2e training

We can now launch training using the new model. This can be done using the following command:

```
CONF=my_new_config.yaml ./train.sh
```

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Integrating a new dataset/dataloader

In this section, we showcase how to integrate a new dataset/dataloader with the library.

### Table of contents

- *Building Dataset module*
- *Building DataModule*
  - *Plug into training.py*
  - *Create config file*
  - *Launching e2e training*

## Building Dataset module

One can use the guide on [PyTorch docs](#) to create a Dataset class.

## Building DataModule

To configure the dataloader, one needs to create a DataModule class. Neuronx Distributed Training library provides a BaseDataModule which one can use to implement their new DataModule. Create a new file called new\_data\_module.py and add the following content.

```
from neuronx_distributed_training.lightning_modules.data.base import BaseDataModule

class NewDataModule(BaseDataModule):
    def __init__(self, cfg, trainer):
        """
        DataModule class for configuring the dataset/dataloader

        Args:
            cfg: `data` cfg in the yaml file.
            trainer: PyTorch-Lightning trainer.
        """
        super().__init__(cfg, trainer)
        # Users can use the cfg argument to pass down
        # arguments from the yaml file to the DataModule.
```

(continues on next page)

(continued from previous page)

```

def get_batch_length(self, batch):
    """
    Returns the length of the batch.
    """
    return len(batch["input_ids"])

def process_global_batch(self, global_batch, global_batch_size=None):
    """ Any custom processing of batches can be done here.

    Args:
        global_batch: list of inputs, eg. [tokens, labels]
        global_batch_size: Length of tokens and labels
    """
    return global_batch

def train_dataloader(self):
    """
    This API should return a torch.utils.data.dataloader.DataLoader object
    """
    ...

def val_dataloader(self):
    """
    This API should return a torch.utils.data.dataloader.DataLoader object
    """
    ...

def test_dataloader(self):
    """
    This API should return a torch.utils.data.dataloader.DataLoader object
    """
    ...

```

## Plug into training.py

Once the new data module is created, we can then plug this into the `training.py` script under `examples` folder. We can modify the `training.py` script as follows:

```

...
# Assuming we are using the same ModelModule we used for LLama example.
from new_data_module import NewDataModule
data_module = NewDataModule(cfg, trainer)
model = HFLlamaModule(cfg, trainer)

trainer.fit(model, datamodule=data_module)

```

The rest of the code can remain the same. The trainer will now use the `NewDataModule` for fetching the `dataloader` and run e2e training.

## Create config file

Next, we can create a config file under `conf` to be used for this new dataloader. We can start with a copy of `hf_llama_7B_config.yaml`. Let's call this config file `my_new_config.yaml`. We can edit the `data` key to configure the `DataModule`

**Note:** For the model, we are using the same model that the llama example is using. To configure a new model, please check the [Integrating a New Model](#) section.

## Launching e2e training

We can now launch training using the new `data_module`. This can be done using the following command:

```
CONF=my_new_config.yaml ./train.sh
```

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Registering an optimizer and LR scheduler

A new optimizer or LR scheduler can be registered with the framework and enabled from the config.

### Table of contents

- [Setting up the optimizer](#)
- [Setting up the LR scheduler](#)

## Setting up the optimizer

One can write their own optimizer class. One such example is the `AdamW_FP32OptimParams`.

The inputs to the optimizer can be exposed in the config YAML file. To do this, we need to create a `Params` class as shown below:

```
from dataclasses import dataclass
from typing import Any, Dict, Optional, Tuple

from omegaconf import MISSING

@dataclass
class OptimizerParams:
    """
    All the params listed below can be configured from the YAML file
    """

    lr: Optional[float] = MISSING
    betas: Tuple[float, float] = (0.9, 0.999)
```

(continues on next page)

(continued from previous page)

```
eps: float = 1e-08
weight_decay: float = 0
amsgrad: bool = False
```

Once we create the optimizer and the optimizer params class, we can now register the optimizer with the framework using the following code:

```
from nemo.core.optim import register_optimizer

# `adamw_fp32OptState` would be the name in the optim config of the YAML file.
register_optimizer("adamw_fp32OptState", AdamW_FP32OptimParams, OptimizerParams)
```

This registration can be done inside the `training.py` file which resides in `examples` folder.

Once the registration is done, we can now expose the `OptimizerParams` under `optim` config of the YAML file.

### Setting up the LR scheduler

One can write their own LR scheduler and register with the framework. One such example of LR scheduler is shown below:

```
from functools import partial

from torch.optim.lr_scheduler import LambdaLR
from transformers.optimization import _get_linear_schedule_with_warmup_lr_lambda

class LinearAnnealingWithWarmUp(LambdaLR):
    def __init__(self, optimizer, warmup_steps, max_steps, last_epoch=-1):
        lr_lambda = partial(
            _get_linear_schedule_with_warmup_lr_lambda,
            num_warmup_steps=warmup_steps,
            num_training_steps=max_steps,
        )
        super().__init__(optimizer, lr_lambda, last_epoch)
```

Once we build this LR scheduler, we can expose the arguments to the config YAML file. Before that, we need to write up a `LRSchedulerParams` class. Here is an example for the same:

```
from nemo.core.config.schedulers import SchedulerParams

class LinearAnnealingWithWarmupParams(SchedulerParams):
    warmup_steps: int = 0
    max_steps: int = 0
```

Once the LR scheduler and the `SchedulerParams` class are set, we can now register the scheduler with the framework as below:

```
from nemo.core.optim.lr_scheduler import register_scheduler

# Here, `LinearAnnealingWithWarmUp` is the name of the scheduler we would use in the
```

(continues on next page)

(continued from previous page)

```
↪ config YAML file
register_scheduler("LinearAnnealingWithWarmUp", LinearAnnealingWithWarmUp, ↪
↪ LinearAnnealingWithWarmupParams)
```

This registration can be done inside the `training.py` file which resides under `examples` folder.

Once the registration is done, we can now expose the `LinearAnnealingWithWarmupParams` under `sched` config of the YAML file.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## Migrating from Neuron-NeMo-Megatron to Neuronx Distributed Training

In this section, we go over the changes one would have to make if they are migrating their training workload from [Neuronx-NeMo-Megatron \(NNM\)](#) to Neuronx Distributed Training (NxDT) framework.

### Table of contents

- [Config migration](#)
- [Model code](#)
- [Checkpointing Save/Load](#)
- [Config Mapping](#)

## Config migration

NxDT is a framework built on top of [NeMo](#) and [NeuronxDistributed \(NxD\)](#) and supports megatron-style model. The megatron model implementation is ported over from NNM. Hence, most of the config YAMLs from NNM can be migrated to use NxDT.

When building NxDT for the sake of modularity, we grouped certain parameters together, eg. `distributed_strategy` has all the configuration for model parallelism, `data` config now holds all the parameters required to configure the dataset.

At a high level, there are some differences with the NNM config, which are highlighted below:

1. The overall config structure has changed. For simplicity and ease of understanding, the config parameters are grouped according to their high level use case. For example, previously all the distributed config parameters used to reside inside `model` config, now it's been moved to a `distributed_config` of its own. Similarly data config is moved out to have clear separation between model and data.
2. Environment variables like `neuron_cc_flags` and `neuron_compile_cache_url` can be set from the config itself. There is no need to set the environment variables. The rationale is to avoid having to configure training scripts from multiple places.
3. **Activation Checkpointing:** NxDT only supports selective and full activation checkpointing. The selective checkpointing is done only for the `CoreAttention` block (in case of llama3-8K we recompute the MLP block, too) and full activation checkpointing is done only at a layer boundary. NxDT doesn't support config parameters like `activations_checkpoint_method`, `activations_checkpoint_num_layers`, `num_micro_batches_with_partial_activation_checkpoints`, `activations_checkpoint_layers_per_pipeline`, `disable_layer_norm_checkpointing`. Please remove these parameters from your config.yaml file.

---

**Note:** If you plan to add more modules that need to be recomputed, one would have to override the checkpointing config inside `ModelModule` (refer to `build_model` API at [Build a Lightning Module](#)) and add the modules that need to be recomputed.

---

4. **Tokenizer:** The tokenizer which used to reside under `model` is now moved to `data`. This is done so that all data related configuration can reside at one place.
5. **accumulate\_grad\_batches:** This param is removed since it should always be 1. Gradient accumulation is handled by setting the `global_batch_size` and `micro_batch_size` along with data-parallel degree.
6. **pre\_process** and **post\_process::** These two parameters were added to the model to decide if the embedding lookup needs to be added at the start and if a `pooler` layer needs to be added at the end. This has been set by default for all decoder models and hence the config param is no longer exposed.
7. **Mixed precision config:** `NxD` no longer exposes NeMo mixed precision parameters: `native_amp_init_scale`, `native_amp_growth_interval`, `hysteresis`, `fp32_residual_connection`, `fp16_lm_cross_entropy`. All these parameters are specific to the GPU mixed precision strategy, which Neuron doesn't support, or they are not applicable. Neuron has a different way to enable mixed precision training through `master_weights` and `fp32_grad_accumulation`.
8. **megatron\_amp\_o2:** This parameter is not supported.
9. **Fusions:** Neuron doesn't support fusion parameters like `grad_div_ar_fusion`, `gradient_accumulation_fusion`, `bias_activation_fusion`, `bias_dropout_add_fusion`, `masked_softmax_fusion`. All of these fusions are built for GPU and require CUDA kernels which cannot run on `Trn1`. Neuron would have its own set of kernels and when we support them, we would enable those parameters from the config.

---

**Note:** If there is a need to support these configs, please create a feature request with exact needs and we shall work on it.

---

For detailed mapping, please check the [Config Mapping](#).

## Model code

There are the following differences in the model code:

1. NNM used [Apex](#) to get all the distributed parallel layers and schedules. Since `NxD` uses `NxD` as the base library, all the `parallel layers/parallel state` are coming from `NxD`. Eg. `apex.parallel_state` is replaced with `nxd.parallel_layers.parallel_state`.
2. NNM explicitly creates a module for each pipeline-parallel (PP) rank, however, `NxD` uses `NxD` which does the partitioning under the hood. Hence, users no longer have to worry about creating a rank specific module. They can create one single model and `NxD's PP wrapper` takes care of sharding for each PP rank. Hence, all the code related to pipeline parallelism inside model code is removed. The model code assumes there is no PP and just uses TP layers from `NxD`.

---

**Note:** For the tracer to work efficiently, we configure the pipeline parallel config inside the `BaseModelModule` class inside `lightning_modules/model`.

---

3. In NNM, megatron module had to explicitly handle gradient reduction for shared weights across PP ranks. In `NxD`, since we are using `NxD's PP wrapper`, all that is handled for the user.

4. For activation checkpointing, NNM had explicit recompute functions which handled the [custom forward API](#). With NxDT, [NxDT's Activation Checkpoint wrapper](#) handles the recompute of the modules. Users just have to configure the `activation_checkpoint_config` inside `nxd_config` [here](#).

## Checkpointing Save/Load

NxDT supports all the checkpointing features which NNM supports. This includes async checkpointing, auto-resume, etc. There are some differences in the format of the checkpoint. This is because NxDT uses [NxDT's checkpoint api](#). The key differences are listed below:

1. NNM combines the model weights, optimizers and other state\_dicts into a single state\_dict and dump a file of the format: `tp_rank_0*_pp_rank_00*/model_optim_rng.ckpt`. However, with NxDT, we save the model state\_dict and the optimizer separately. The model state\_dict is saved in a folder of the form: `model/dp_rank_00_tp_rank_00_pp_rank_00.pt` and the optimizer is saved into a separate folder as: `optim/dp_rank_00_tp_rank_00_pp_rank_00.pt`. This is mainly done so that when we use [zero1](#), each DP rank can save its own optimizer shard.
2. In NNM, if we are using pipeline parallelism, each pipeline stage creates an independent model. So let's say we have a model with 32 layers and we use PP=4, then NNM would create 4 chunks with layers 0-7. So each PP rank would have model\_state\_dict with keys going from layer-0-7. However, in NxDT, the model is created as a whole and then sharded. So the layer numbers are preserved.
3. There are checkpoint conversion scripts provided under `examples/` of NxDT repository to convert the existing NNM checkpoints to NxDT format in case of migrating in the middle of training.

```
python nnm_nxd_t_ckpt_converter.py --tp 8 --pp 4 --n_layers 32 --nnm_ckpt_path {path_to_ckpt}/ckpt/nnm --nxd_ckpt_path {path_to_ckpt}/nnm-converted-nxd-ckpt/ --enable_parallel_processing True --num_parallel_processes 8
```

## Config Mapping

Here is a detailed mapping for all the parameters in the config file. For the below mapping, we chose the Llama-7B example across NNM and NxDT frameworks. The same mapping is also true for other models.

| NNM param            | NxDT param mapping                        | Comments   |
|----------------------|---|--|
| name                 | name                                      |  |
| restore_from_path    | Not supported                             | This config was not fully supported in NNM, either.  |
| <b>trainer</b>       |   |  |
| devices              | devices                                   |  |
| num_nodes            | num_nodes                                 |  |
| accelerator          | Not required                              | We made the default as TPU which maps to Neuron internally, so users no longer have to add it. |
| precision            | replaced by <code>precision_config</code> | There is a separate <i>precision</i> config to control the precision of model and optimizer.   |
| logger               | Replaced by default                       | We made the NNM logger default in NxDT.  |
| enable_checkpointing | Separate <code>exp_manager</code> config  | All checkpointing is controlled by <code>exp_manager</code> config.                            |
| replace_sampler_ddp  | Not supported                             | Had to be always False in NNM, made it default in NxDT. No setting required.                   |
| max_epochs           | max_epochs                                |  |
| max_steps            | max_steps                                 |  |

continues on next page

Table 3.1 – continued from previous page

| NNM param                            | NxDT param mapping  | Comments  |
|--------------------------------------|---|---|
| log_every_n_steps                    | log_every_n_steps   |   |
| val_check_interval                   | val_check_interval  |   |
| limit_val_batches                    | limit_val_batches   |   |
| limit_test_batches                   | limit_test_batches  |   |
| accumulate_grad_batches              | Removed   | This is automatically configured based on global_batchsize, micro-batchsize and distributed config.                 |
| gradient_clip_val                    | gradient_clip_val   |   |
| benchmark                            | Not supported   |   |
| enable_model_summary                 | Not supported   |   |
| <b>exp_manager</b>                   |   |   |
| log_local_rank_0_only                | log_local_rank_0_only                                     |   |
| create_tensorboard_logger            | create_tensorboard_logger                                 |   |
| explicit_log_dir                     | explicit_log_dir  |   |
| exp_dir                              | exp_dir   |   |
| name                                 | name  |   |
| create_wandb_logger                  | Not supported   | This was not supported under NNM, either. We have removed this argument from NxDT.                                  |
| wandb_logger_kwargs                  | Not supported   |   |
| resume_if_exists                     | resume_if_exists  |   |
| resume_ignore_no_checkpoint          | resume_ignore_no_checkpoint                               |   |
| create_checkpoint_callback           | create_checkpoint_callback                                |   |
| checkpoint_callback_params           | checkpoint_callback_params                                |   |
| <b>model</b>                         |   |   |
| tensor_model_parallel_size           | distributed_strategy.tensor_model_parallel_size           | All the parallelism config are moved to distributed_strategy config.  |
| pipeline_model_parallel_size         | distributed_strategy.pipeline_model_parallel_size         |   |
| virtual_pipeline_model_parallel_size | distributed_strategy.virtual_pipeline_model_parallel_size |   |
| sequence_parallel                    | distributed_strategy.sequence_parallel                    |   |
| wrap_with_zero                       | distributed_strategy.zero1                                |   |
| micro_batch_size                     | data.micro_batch_size                                     | All the dataset/dataloader/tokenizer configurations are now part of a separate config called data.                  |
| global_batch_size                    | data.global_batch_size                                    |   |
| tokenizer                            | data.tokenizer  |   |
| data                                 | Moved to data at the same level as model                  | The entire data key now controls a DataModule and is placed at the same level as model key in the config structure. |
| encoder_seq_length                   | encoder_seq_length  |   |
| max_position_embeddings              | max_position_embeddings                                   |   |
| make_vocab_size_divisible_by         | make_vocab_size_divisible_by                              |   |

continues on next page



Table 3.1 – continued from previous page

| NNM param                           | NxDT param mapping  | Comments   |
|-------------------------------------|---|--|
| pre_process                         | Not supported   | NxDT by default adds embedding layer at the start of the transformer block.  |
| post_process                        | Not supported   | NxDT by default adds a LM-head at the end of the transformer block.  |
| persist_layer_norm                  | persist_layer_norm  |  |
| share_embeddings_and_output_weights | share_embeddings_and_output_weights                       |  |
| position_embedding_type             | position_embedding_type                                   |  |
| rotary_percentage                   | rotary_percentage   |  |
| transformer_block_type              | transformer_block_type                                    |  |
| has_bias                            | has_bias  |  |
| native_amp_init_scale               | Not required  |  |
| native_amp_growth_interval          | Not required  | GPU optimizations which were not supported in NNM, have been removed from NxDT. Most of these fusion ops, the neuron compiler handles on its own. For Attention and Softmax, Neuron uses NKI kernels and custom ops to implement them. |
| hysteresis                          | Not required  |  |
| fp32_residual_connection            | Not required  |  |
| fp16_lm_cross_entropy               | Not required  |  |
| megatron_amp_O2                     | Not required  |  |
| grad_div_ar_fusion                  | Not required  |  |
| gradient_accumulation_fusion        | Not required  |  |
| bias_activation_fusion              | Not required  |  |
| bias_dropout_add_fusion             | Not required  |  |
| masked_softmax_fusion               | fusions.softmax   |  |
| seed                                | Seed is moved out of model and at the same level as model |  |
| resume_from_checkpoint              | exp_manager.resume_from_checkpoint                        |  |
| use_cpu_initialization              | use_cpu_initialization                                    |  |
| onnx_safe                           | Not supported   | This was not supported under NNM, either. We have removed this argument from NxDT.   |
| apex_transformer_log_level          | Not supported   |  |
| gradient_as_bucket_view             | Not supported   |  |
| sync_batch_comm                     | Not supported   |  |
| log_parameter_norm                  | exp_manager.log_gradient_norm                             |  |
| log_gradient_norm                   | exp_manager.log_gradient_norm                             |  |
| flexible_pipeline_parallel_stages   | Not supported   |  |
| activations_checkpoint_granularity  | activations_checkpoint_granularity                        | Currently, NxDT checkpoints the attention module in case of selective and a single layer in case of full checkpointing.  |
| activations_checkpoint_method       | Not supported   |  |
| activations_checkpoint_num_layers   | Not supported   |  |

continues on next page

Table 3.1 – continued from previous page

| NNM param   | NxDT param mapping                                   | Comments   |
|---|--|--|
| num_micro_batches_with_partial_activation_checkpoints | Not supported  |  |
| activations_checkpoint_layers_per_pipeline            | Not supported  |  |
| disable_layer_norm_checkpointing                      | Not supported  |  |
| zero_use_master_weight                                | Supported via precision config                       | See <i>manual precision config</i> .                         |
| zero_use_fp32_grad_accum                              | Supported via precision config                       | See <i>manual precision config</i> .                         |
| transformer_engine                                    | Not supported  | This is specifically built for NVIDIA GPUs.                  |
| fp8   | Not supported  | fp8 training is not supported on Neuron (both NNM and NxDT). |
| fp8_e4m3  | Not supported  | fp8 training is not supported on Neuron (both NNM and NxDT). |
| fp8_hybrid  | Not supported  | fp8 training is not supported on Neuron (both NNM and NxDT). |
| fp8_margin  | Not supported  | fp8 training is not supported on Neuron (both NNM and NxDT). |
| use_emha  | Not supported  | fp8 training is not supported on Neuron (both NNM and NxDT). |
| convert_to_hf   | Supported via separate script                        |  |
| nsys_profile  | Not supported  | This is specifically built for NVIDIA GPUs.                  |
| optim   | optim  |  |
| enable_recovery_time_instrumentation                  | exp_manager.<br>enable_recovery_time_instrumentation |  |
| async_checkpointing                                   | exp_manager.<br>async_checkpointing                  |  |

**Note:** For parameters that are not supported by NxDT, please create a feature request with specific use-case for the parameter, if needed.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## NxD Training Compatibility with NeMo

NxD Training (NxDT) is built on top of [NeMo-1.14](#). The framework reuses modules from NeMo and exposes them via similar config interface.

**Note:** At the moment, NxDT only allows running training of decoder LLM models.

This document goes over steps on how to run the NeMo training workloads inside NxDT.

### Table of contents

- *Model Integration*
- *Dataloader Integration*
- *Optimizer/LR Scheduler Integration*
- *Optimal Partitioning*
- *Fusions/kernels*
- *Checkpoint Saving/loading*
- *Config Mapping*

## Model Integration

### Model already Exists in NxDT Model Hub:

If the model you want to train is already included in the NxDT model hub, and the training workflow (e.g., pre-training, fine-tuning) is supported in NxDT, you need to modify NeMo YAML configuration file to the NxDT YAML file. Follow the mapping table in the [Config Mapping](#).

### Custom/New Model

If your model is not part of the NxDT model hub, please use the guide [Integrating a New Model](#).

## Dataloader Integration

### Dataloader already exposed via one of the NxDT configs

In this case, please map the NeMo YAML config parameters to NxDT config parameters using the mapping table provided here [Config Mapping](#).

### Custom/New Dataloader

If the dataloader is not part of the hub, please use the guide [Integrating a new dataset/dataloader](#).

## Optimizer/LR Scheduler Integration

Since NxDT is built on top of NeMo, all the optimizers/LR schedulers provided by NeMo can be enabled from the config.

## Optimal Partitioning

NxDT is built on top of [NxDT Core](#) primitives and exposes different model parallelism techniques. All of them can be configured using the `distributed_strategy` config.

## Fusions/kernels

All the fused kernels available inside the NeMo config are not available in NxDT. This is because fused kernels in NeMo are built specifically for GPUs. Neuron have a different set of kernels that can be enabled from the config. Also, since Neuron uses a graph based approach, the compiler can optimize some of the modules and do fusions wherever required.

## Checkpoint Saving/loading

1. NeMo combines the model weights, optimizers and other state\_dicts into a single state\_dict and dumps a file of the format: `tp_rank_0*_pp_rank_00*/model_optim_rng.ckpt`. However, with NxDT, we save the model state\_dict and the optimizer separately. The model statedict is saved in a folder of the form: `model/dp_rank_00_tp_rank_00_pp_rank_00.pt` and the optimizer is saved into a separate folder as: `optim/dp_rank_00_tp_rank_00_pp_rank_00.pt`. This is mainly done so that when we use [zero1](#), each DP rank can save its own optimizer shard.
2. NxDT doesn't support `.nemo` style checkpoint saving. If users have a `.nemo` checkpoint, they would have to unpack it themselves and build a checkpoint conversion script to load the checkpoint into NxDT.
3. In NeMo, if we are using pipeline parallel, each pipeline stage creates an independent model. So lets say we have a model with 32 layers and we use PP=4, then NeMo would create 4 chunks with layers 0-7. So each PP rank would have a `model_state_dict` with keys going from layer-0-7. However, in NxDT, the model is created as a whole and then sharded. So the layer numbers are preserved.
4. One would have to write up a checkpoint conversion script similar to the checkpoint conversion from NeMo to NxDT.

For a more detailed mapping of NeMo parameters to NxDT parameters, follow the guide [Config Mapping](#).

## Config Mapping

Here is a detailed mapping for all the parameters in the config file. For the below mapping, we chose the Llama example across both NeMo and NxDT frameworks. The same mapping is also true for other models.

| NeMo param              | NxDT param mapping                        | Comments   |
|-------------------------|---|--|
| name                    | name                                      |  |
| restore_from_path       | Not supported                             |  |
| <b>trainer</b>          |   |  |
| devices                 | devices                                   |  |
| num_nodes               | num_nodes                                 |  |
| accelerator             | Not required                              | We made the default as TPU which maps to Neuron internally, so users no longer have to add it. |
| precision               | replaced by <code>precision_config</code> | There is a separate <i>precision</i> config to control the precision of model and optimizer.   |
| logger                  | Not required                              | We set default value of logger to False.   |
| enable_checkpointing    | Separate <code>exp_manager</code> config  | All checkpointing is controlled by <code>exp_manager</code> config.                            |
| use_distributed_sampler | Not supported                             |  |
| max_epochs              | max_epochs                                |  |
| max_steps               | max_steps                                 |  |
| log_every_n_steps       | log_every_n_steps                         |  |
| val_check_interval      | val_check_interval                        |  |

continues on next page

Table 3.2 – continued from previous page

| NeMo param                           | NxDT param mapping  | Comments  |
|--------------------------------------|---|---|
| limit_val_batches                    | limit_val_batches   |   |
| limit_test_batches                   | limit_test_batches  |   |
| accumulate_grad_batches              | Removed   | This is automatically configured based on global_batchsize, micro-batchsize and distributed config.                 |
| gradient_clip_val                    | gradient_clip_val   |   |
| benchmark                            | Not supported   |   |
| enable_model_summary                 | Not supported   |   |
| <b>exp_manager</b>                   |   |   |
| log_local_rank_0_only                | log_local_rank_0_only                                     |   |
| create_tensorboard_logger            | create_tensorboard_logger                                 |   |
| explicit_log_dir                     | explicit_log_dir  |   |
| exp_dir                              | exp_dir   |   |
| name                                 | name  |   |
| create_wandb_logger                  | Not supported   | This was not supported under NNM, either. We have removed this argument from NxDT.                                  |
| wandb_logger_kwargs                  | Not supported   |   |
| resume_if_exists                     | resume_if_exists  |   |
| resume_ignore_no_checkpoint          | resume_ignore_no_checkpoint                               |   |
| create_checkpoint_callback           | create_checkpoint_callback                                |   |
| checkpoint_callback_params           | checkpoint_callback_params                                |   |
| <b>model</b>                         |   |   |
| mcgpt                                | Not supported   | NxDT has its own implementation of megatron_gpt_model which is based on v1.14 version of NeMo                       |
| tensor_model_parallel_size           | distributed_strategy.tensor_model_parallel_size           | All the parallelism config are moved to distributed_strategy config   |
| pipeline_model_parallel_size         | distributed_strategy.pipeline_model_parallel_size         |   |
| virtual_pipeline_model_parallel_size | distributed_strategy.virtual_pipeline_model_parallel_size |   |
| sequence_parallel                    | distributed_strategy.sequence_parallel                    |   |
| micro_batch_size                     | data.micro_batch_size                                     | All the dataset/dataloader/tokenizer configuration are now part of a separate config called data                    |
| global_batch_size                    | data.global_batch_size                                    |   |
| tokenizer                            | data.tokenizer  |   |
| data                                 | Moved to data at the same level as model                  | The entire data key now controls a DataModule and is placed at the same level as model key in the config structure. |
| encoder_seq_length                   | encoder_seq_length  |   |
| max_position_embeddings              | max_position_embeddings                                   |   |
| make_vocab_size_divisible_by         | make_vocab_size_divisible_by                              |   |
| pre_process                          | Not supported   | NxDT by default adds embedding layer at the start of the transformer block.   |

continues on next page

Table 3.2 – continued from previous page

| NeMo param  | NxDT param mapping  | Comments  |
|---|---|---|
| post_process  | Not supported   | NxDT by default adds a LM-head at the end of the transformer block.   |
| persist_layer_norm                                    | persist_layer_norm  |   |
| share_embeddings_and_output_weights                   | share_embeddings_and_output_weights                       |   |
| position_embedding_type                               | position_embedding_type                                   |   |
| rotary_percentage                                     | rotary_percentage   |   |
| transformer_block_type                                | transformer_block_type                                    |   |
| has_bias  | has_bias  |   |
| num_query_groups                                      | Not required  | query group attention can be configured using num_kv_heads parameter.   |
| native_amp_init_scale                                 | Not Required  |   |
| native_amp_growth_interval                            | Not Required  | GPU optimizations which were not supported in NNM, have been removed from NxDT. Most of these fusion ops, the neuron compiler handles on its own. For Attention and Softmax, Neuron uses NKI kernels and custom ops to implement them |
| hysteresis  | Not Required  |   |
| fp32_residual_connection                              | Not Required  |   |
| fp16_lm_cross_entropy                                 | Not Required  |   |
| megatron_amp_O2                                       | Not Required  |   |
| grad_div_ar_fusion                                    | Not Required  |   |
| gradient_accumulation_fusion                          | Not Required  |   |
| bias_activation_fusion                                | Not Required  |   |
| bias_dropout_add_fusion                               | Not Required  |   |
| masked_softmax_fusion                                 | fusions.softmax   |   |
| seed  | seed is moved out of model and at the same level as model |   |
| resume_from_checkpoint                                | exp_manager.resume_from_checkpoint                        |   |
| use_cpu_initialization                                | use_cpu_initialization                                    |   |
| onnx_safe   | Not supported   | This was not supported under NNM too, we have removed this argument from NxDT.  |
| apex_transformer_log_level                            | Not supported   |   |
| gradient_as_bucket_view                               | Not supported   |   |
| sync_batch_comm                                       | Not supported   |   |
| activations_checkpoint_granularity                    | activations_checkpoint_granularity                        | By default NxDT checkpoints attention module in case of selective and a single layer in case of full checkpointing.   |
| activations_checkpoint_method                         | Not supported   |   |
| activations_checkpoint_num_layers                     | Not supported   |   |
| num_micro_batches_with_partial_activation_checkpoints | Not supported   |   |
| activations_checkpoint_layers_per_pipeline            | Not supported   |   |
| disable_layer_norm_checkpointing                      | Not supported   |   |
| transformer_engine                                    | Not supported   | This is specifically built for NVIDIA GPUs.   |

continues on next page

Table 3.2 – continued from previous page

| NeMo param   | NxD param mapping | Comments  |
|--------------|-------------------|---|
| fp8          | Not supported     | fp8 training is not supported on Neuron (both NNM and NxD). |
| fp8_e4m3     | Not supported     |   |
| fp8_hybrid   | Not supported     |   |
| fp8_margin   | Not supported     |   |
| use_emha     | Not supported     |   |
| nsys_profile | Not supported     | This is specifically built for NVIDIA GPUs.                 |
| optim        | optim             |   |

**Note:** For parameters that are not supported by NxD, please create a feature request with specific use-case for the parameter, if needed.

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### 3.1.6 Tutorials

This section will go over tutorials to help users get started with NxD Training library.

*This document is relevant for: Trn1, Trn2*

#### Megatron GPT Pretraining

In this example, we will compile and train a Megatron GPT model on a single instance or on multiple instances using ParallelCluster with the NxD Training library. The example has the following main sections:

##### Table of contents

- *Setting up the environment*
  - *ParallelCluster Setup*
  - *Install Dependencies*
- *Download the dataset*
- *Pre-compile the model*
- *Training the model*
- *Monitoring Training*
  - *Tensorboard monitoring*
  - *neuron-top / neuron-monitor / neuron-ls*
- *Troubleshooting Guide*

## Setting up the environment

### ParallelCluster Setup

In this example, we will use 8 instances with ParallelCluster, please follow the instructions here to create a cluster: [Train your model on ParallelCluster](#)

ParallelCluster automates the creation of trn1 clusters, and provides the SLURM job management system for scheduling and managing distributed training jobs. Please note that the home directory on your ParallelCluster head node will be shared with all of the worker nodes via NFS.

### Install Dependencies

Once you have launched a trn1 instance or ParallelCluster, please follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxD Training and its dependencies. Please see the following installation guide for installing NxD Training: [NxD Training: NxDT Installation Guide](#)

### Download the dataset

This tutorial makes use of a preprocessed Wikipedia dataset that is stored in S3. The dataset can be downloaded to your cluster or instance by running the following commands on the head node or your trn1 instance:

```
export DATA_DIR=~/examples_datasets/gpt2
mkdir -p ${DATA_DIR} && cd ${DATA_DIR}
wget https://s3.amazonaws.com/models.huggingface.co/bert/gpt2-vocab.json
wget https://s3.amazonaws.com/models.huggingface.co/bert/gpt2-merges.txt
aws s3 cp s3://neuron-s3/training_datasets/gpt/wikipedia/my-gpt2_text_document.bin . --
↪no-sign-request
aws s3 cp s3://neuron-s3/training_datasets/gpt/wikipedia/my-gpt2_text_document.idx . --
↪no-sign-request
aws s3 cp s3://neuron-s3/training_datasets/gpt/wikipedia/license.txt . --no-sign-request
```

### Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

```
git clone https://github.com/aws-neuron/neuronx-distributed-training
cd neuronx-distributed-training/examples
```



Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `megatron_gpt_config`. The default config here is a 6.7B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

Next, run the following commands to launch an AOT pre-compilation job on your instance:

```
export COMPILE=1
./train.sh
```

The compile output and logs will be shown directly in the terminal and you will see a message similar to this:

```
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total graphs: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total successful compilations: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total failed compilations: 0
```

Then, you know your compilation has successfully completed.

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

If you are using `ParallelCluster`, then you will need to update the `conf/megatron_gpt_config.yaml` with

```
num_nodes: 8
```

Then to run the compile job:

```
export COMPILE=1
sbatch --exclusive \
  --nodes 8 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

Once you have launched the precompilation job, run the `squeue` command to view the SLURM job queue on your cluster. If you have not recently run a job on your cluster, it may take 4-5 minutes for the requested `trn1.32xlarge` nodes to be launched and initialized. Once the job is running, `squeue` should show output similar to the following:

| JOBID | PARTITION | NAME | USER   | ST | TIME | NODES | NODELIST(REASON)            |
|-------|-----------|------|--------|----|------|-------|-----------------------------|
| 10    | compute1  | wrap | ubuntu | R  | 5:11 | 8     | compute1-dy-queue1-i1-[0-7] |

You can view the output of the precompilation job by examining the file named `slurm-ZZ.out`, where `ZZ` represents the JOBID of your job in the `squeue` output above.

```
tail -f slurm-10.out
```

Once the precompilation job is complete, just like the above output you should see a message similar to the following in the logs:

```
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total graphs: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total successful compilations: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total failed compilations: 0
```

At this point, you can press `CTRL-C` to exit the `tail` command.

## Training the model

The pre-training job is launched almost exactly the same as the compile job. We now turn off the `COMPILE` environment variable and run the same training script to start pre-training.

On a single instance:

```
export COMPILE=0
./train.sh
```

If you are using `ParallelCluster`:

```
export COMPILE=0
sbatch --exclusive \
  --nodes 8 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

As outlined above, you can again use the `squeue` command to view the job queue, and also monitor the job in the same way with the `tail` command to see the training logs. Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 0%|          | 189/301501 [59:12<1573:03:24, 18.79s/it, loss=7.75, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=188.0, consumed_samples=24064.0]
Epoch 0: 0%|          | 190/301501 [59:30<1572:41:13, 18.79s/it, loss=7.74, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=189.0, consumed_samples=24192.0]
Epoch 0: 0%|          | 191/301501 [59:48<1572:21:28, 18.79s/it, loss=7.73, v_num=3-16,
↳ reduced_train_loss=7.910, global_step=190.0, consumed_samples=24320.0]
```

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use standard tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/megatron_gpt/`. Once you have identified the directory, `cd` into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/megatron_gpt/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```

## neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, first SSH into one of your compute nodes from the head node (if using `ParallelCluster`), and then run `neuron-top`:

```
ssh compute1-dy-queue1-i1-1 # to determine which compute nodes are in use, run the
↪ squeue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Troubleshooting Guide

For issues with NxD Training, please see: [NxD Training Known Issues](#)

For `ParallelCluster` issues see: [AWS ParallelCluster Troubleshooting](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## HuggingFace Llama3.1/Llama3-8B Pretraining

In this example, we will compile and train a HF Llama3.1/Llama3-8B model on a single instance with the NxD Training (NxD) library. The example has the following main sections:

### Table of contents

- *Setting up the environment*
  - *Install Dependencies*
- *Download the dataset*
- *Pre-compile the model*
- *Training the model*
- *Monitoring Training*
  - *Tensorboard monitoring*
  - *neuron-top / neuron-monitor / neuron-ls*
- *Troubleshooting Guide*

## Setting up the environment

### Install Dependencies

Once you have launched a Trn1 instance, please follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxDT and its dependencies. Please see the following installation guide for installing NxDT: [NxDT Installation Guide](#)

### Download the dataset

Let's download training-data scripts for our experiments

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪ training/llama/get_dataset.py
```

To tokenize the data, we must request the tokenizer from Hugging Face and Meta by following the instructions at the following link: [HuggingFace Llama 3 8B Model](#).

Use of the Llama models is governed by the Meta license. In order to download the model weights and tokenizer, please visit the above website and accept their License before requesting access. After access has been granted, you may use the following python3 script along with your own hugging face token to download and save the tokenizer.

```
from huggingface_hub import login
from transformers import AutoTokenizer

login(token='your_own_hugging_face_token')

tokenizer = AutoTokenizer.from_pretrained('meta-llama/Meta-Llama-3-8B')

tokenizer.save_pretrained(".")
```

For Llama3.1/Llama3, make sure your base directory has the following files:

```
'./tokenizer_config.json', './special_tokens_map.json', './tokenizer.json'
```

Next let's download and pre-process the dataset:

```
mkdir ~/examples_datasets/ && cd ~/examples_datasets/
python3 ~/get_dataset.py --llama-version 3
```

*Note:* In case you see an error of the following form when downloading data: `huggingface_hub.utils._validators.HFValidationError: Repo id must be in the form 'repo_name' or 'namespace/repo_name'. Use 'repo_type' argument if needed.` This could be because of a stale cache. Try deleting the cache using:

```
sudo rm -rf ~/.cache/
```

## Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

```
git clone https://github.com/aws-neuron/neuronx-distributed-training
cd neuronx-distributed-training/examples
```

Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `hf_llama3_8B_config`. The default config here is a 8B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

Next, run the following commands to launch an AOT pre-compilation job on your instance:

```
export COMPILE=1
./train.sh
```

The compile output and logs will be shown directly in the terminal and you will see a message similar to this:

```
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total graphs: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total successful compilations: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total failed compilations: 0
```

Then, you know your compilation has successfully completed.

---

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

---

## Training the model

The pre-training job is launched almost exactly the same as the compile job. We now turn off the `COMPILE` environment variable and run the same training script to start pre-training.

On a single instance:

```
export COMPILE=0
./train.sh
```

Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 0%|          | 189/301501 [59:12<1573:03:24, 18.79s/it, loss=7.75, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=188.0, consumed_samples=24064.0]
Epoch 0: 0%|          | 190/301501 [59:30<1572:41:13, 18.79s/it, loss=7.74, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=189.0, consumed_samples=24192.0]
Epoch 0: 0%|          | 191/301501 [59:48<1572:21:28, 18.79s/it, loss=7.73, v_num=3-16,
↳ reduced_train_loss=7.910, global_step=190.0, consumed_samples=24320.0]
```

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use standard tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/`. Once you have identified the directory, cd into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```

### neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, run `neuron-top`:

```
ssh compute1-dy-queue1-i1-1 # to determine which compute nodes are in use, run the
↳ squeue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Troubleshooting Guide

For issues with NxDT, please see: [NxDT Known Issues](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## HuggingFace Llama3.1/Llama3-8B Supervised Fine-tuning

In this example, we will compile and finetune pre-trained HF Llama3.1/Llama3-8B model on a single instance with the NxD Training library. The pre-trained Llama3-8B model serves as the foundation, and we will build upon this solid base by fine-tuning the model to adapt it to a specific task or dataset. The example has the following main sections:

### Table of contents

- *Setting up the environment*
  - *Install Dependencies*
- *SFT-YAML Configuration Overview*
- *Download the dataset*
- *Download pretrained model checkpoint and tokenizer*
  - *Checkpoint Conversion*
- *Pre-compile the model*
- *Training the model*
- *Monitoring Training*
  - *Tensorboard monitoring*
  - *neuron-top / neuron-monitor / neuron-ls*
- *Troubleshooting Guide*

## Setting up the environment

### Install Dependencies

Please follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxD Training and its dependencies. Please see the following installation guide for installing NxD Training: [NxD Training Installation Guide](#).

### SFT-YAML Configuration Overview

You can configure a variety of SFT-specific and model parameters for finetuning through the YAML file.

```
exp_manager:
  resume_from_checkpoint: /pretrained_ckpt

data:
  train_dir: /example_datasets/llama3_8b/training.jsonl
  val_dir: /example_datasets/llama3_8b/validation.json
  dev_choose_samples: 2250
  seq_length: 4096
  alignment_strategy:
    sft:
```

(continues on next page)

(continued from previous page)

```

        packing: True
tokenizer:
    type: /llama3_tokenizer
model:
    weight_init_only: True

```

**exp\_manager****resume\_from\_checkpoint**

Manually set the checkpoint file (pretrained checkpoint) to load from

- **Type:** str
- **Default:** /pretrained\_ckpt
- **Required:** True (start with pretrained checkpoint)

**data****train\_dir**

SFT training data - jsonl or arrow file

As for SFT we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** True

**val\_dir**

SFT validation data - jsonl or arrow file

As for SFT we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** False

**dev\_choose\_samples**

If set, will use that many number of records from the head of the dataset instead of using all.  
Set to null to use full dataset

- **Type:** integer
- **Default:** null
- **Required:** False

**seq\_length**

Set sequence length for the training job.

- **Type:** integer
- **Required:** True

**alignment\_strategy**

Set only when using finetuning specific algorithms (SFT, DPO, etc) and related hyperparameters SFT-specific parameters.



**sft****packing**

Appends multiple records in a single record until seq length supported by model, if false uses pad tokens to reach seq length. Setting it to True increases throughput but might impact accuracy.

- **Type:** bool
- **Default:** False
- **Required:** False

**tokenizer****type**

Set tokenizer path/type

- **Type:** str
- **Default:** /llama3\_tokenizer
- **Required:** True

**model****weight\_init\_only**

Load only model states and ignore the optim states from ckpt directory

- **Type:** bool
- **Default:** True

## Download the dataset

This tutorial makes use of a preprocessed version of *databricks-dolly* instruction-following dataset that is stored in S3. The dataset can be downloaded to your cluster or instance by running the following AWS CLI commands on the head node or your Trn1 instance:

```
export DATA_DIR=~/examples_datasets/llama3_8b
mkdir -p ${DATA_DIR} && cd ${DATA_DIR}
aws s3 cp s3://neuron-s3/training_datasets/llama/sft/training.jsonl . --no-sign-request
aws s3 cp s3://neuron-s3/training_datasets/llama/sft/validation.jsonl . --no-sign-
↪request
```

## Download pretrained model checkpoint and tokenizer

In this tutorial, we will use a pretrained Llama3-8B checkpoint from the original repository. Follow the steps to download tokenizer and model checkpoint from the pretraining stage: <https://llama.meta.com/llama-downloads/>

Alternatively, the model checkpoint and tokenizer can also be downloaded from HuggingFace by following this [guide](#). You can also directly download and convert the HuggingFace model checkpoint using *Direct HuggingFace Model Conversion*.

Create a folder llama3\_tokenizer and copy the tokenizer contents to it.

Modify the following paths in YAML file based on your specific directory configuration:

1. model.model\_config

2. `exp_manager.resume_from_checkpoint`
3. `tokenizer.type`
4. `train_dir` and `val_dir`

You can use your custom model, pretrained checkpoint and tokenizer by modifying the `hf_llama3_8B_SFT_config.yaml` file.

## Checkpoint Conversion

Follow this [Checkpoint Conversion Guide](#) to convert the HF-style Llama3-8B checkpoint to NxDt supported format and store it in `pretrained_ckpt` directory. Modify the config parameter `exp_manager.resume_from_checkpoint` path to the converted pretrained checkpoint path.

## Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

```
git clone https://github.com/aws-neuron/neuronx-distributed-training
cd neuronx-distributed-training/examples
```

Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `hf_llama3_8B_SFT_config`. The default config here is a 8B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

Next, run the following commands to launch an AOT pre-compilation job on your instance:

```
export COMPILE=1
./train.sh
```

The compile output and logs will be shown directly in the terminal and you will see logs similar to this:

```
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total graphs: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total successful compilations: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total failed compilations: 0
```

Then, you know your compilation has successfully completed.

---

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

---

## Training the model

The fine-tuning job is launched almost exactly in the same way as the compile job. We now turn off the `COMPILE` environment variable and run the same training script to start pre-training.

On a single instance:

```
export COMPILE=0
./train.sh
```

Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 0%|          | 189/301501 [59:12<1573:03:24, 18.79s/it, loss=7.75, v_num=3-16,
↪ reduced_train_loss=7.560, global_step=188.0, consumed_samples=24064.0]
Epoch 0: 0%|          | 190/301501 [59:30<1572:41:13, 18.79s/it, loss=7.74, v_num=3-16,
↪ reduced_train_loss=7.560, global_step=189.0, consumed_samples=24192.0]
Epoch 0: 0%|          | 191/301501 [59:48<1572:21:28, 18.79s/it, loss=7.73, v_num=3-16,
↪ reduced_train_loss=7.910, global_step=190.0, consumed_samples=24320.0]
```

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use standard tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/`. Once you have identified the directory, `cd` into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```

### neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, run `neuron-top`:

```
ssh compute1-dy-queue1-i1-1 # to determine which compute nodes are in use, run the
↪ queue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Troubleshooting Guide

For issues with NxD Training, please see: [NxD Training Known Issues](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## HuggingFace Llama3.1/Llama3-8B Efficient Supervised Fine-tuning with LoRA (Beta)

In this example, we will compile and finetune pre-trained HF Llama3.1/Llama3-8B model with LoRA adaptors on a single instance with the NxD Training (NxD) library. LoRA or Low Rank Adaptation allows for parameter-efficient fine-tuning (PEFT) by adding small trainable rank decomposition matrices to specified layer of the model, significantly reducing memory usage and training time compared to dense fine-tuning. The pre-trained Llama3-8B model serves as the foundation, and we will build upon this by fine-tuning the model to adapt it to a specific task or dataset. The example has the following main sections:

### Table of contents

- [Setting up the environment](#)
  - [Install Dependencies](#)
- [Download the dataset](#)
- [Download pretrained model checkpoint and tokenizer](#)
  - [Checkpoint Conversion](#)
- [LoRA SFT-YAML Configuration Overview](#)
- [Pre-compile the model](#)
- [Training the model](#)
- [Monitoring Training](#)
  - [Tensorboard monitoring](#)
  - [neuron-top / neuron-monitor / neuron-ls](#)
- [Troubleshooting Guide](#)

## Setting up the environment

### Install Dependencies

First, you can launch a Trn1 instance by following the Neuron DLAMI guide: [Neuron DLAMI User Guide](#).

Once you have launched a Trn1 instance, follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxD and its dependencies. Please see the following installation guide for installing NxD: [NxD Installation Guide](#).

## Download the dataset

This tutorial makes use of a preprocessed version of *databricks-dolly* instruction-following dataset that is stored in S3. The dataset can be downloaded to your cluster or instance by running the following AWS CLI commands on the head node or your Trn1 instance:

```
export DATA_DIR=~/examples_datasets/llama3_8b
mkdir -p ${DATA_DIR} && cd ${DATA_DIR}
aws s3 cp s3://neuron-s3/training_datasets/llama/sft/training.jsonl . --no-sign-request
aws s3 cp s3://neuron-s3/training_datasets/llama/sft/validation.jsonl . --no-sign-
↪request
```

Then, download the config.json file:

For Llama-3-8B:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪training/llama/tp_zero1_llama_hf_pretrain/8B_config_llama3/config.json ~/
```

## Download pretrained model checkpoint and tokenizer

In this tutorial, we will use a pretrained Llama3-8B checkpoint from the original repository. Follow the steps to download tokenizer and model checkpoint from the pretraining stage: <https://llama.meta.com/llama-downloads/>.

Alternatively, the model checkpoint and tokenizer can also be downloaded from HuggingFace by following this [guide](#).

You can also directly download and convert the HuggingFace model checkpoint using *Direct HuggingFace Model Conversion*.

If you choose to download the weights from HuggingFace with your own token, you can create a python script to run such as:

```
import transformers

tokenizer_path="llama3_tokenizer"
model_weights_path="llama3-8B_hf_weights"
model_id = "meta-llama/Meta-Llama-3-8B"

t = transformers.AutoTokenizer.from_pretrained(model_id)
t.save_pretrained(tokenizer_path)

m = transformers.AutoModelForCausalLM.from_pretrained(model_id)
m.save_pretrained(model_weights_path)
```

Create a folder llama3\_tokenizer and copy the tokenizer contents to it.

Modify the following paths in YAML file based on your specific directory configuration:

1. model.model\_config
2. exp\_manager.resume\_from\_checkpoint
3. tokenizer.type
4. train\_dir and val\_dir

You can use your custom model, pretrained checkpoint and tokenizer by modifying the hf\_llama3\_8B\_SFT\_lora\_config.yaml file.

## Checkpoint Conversion

Follow this [Checkpoint Conversion Guide](#) to convert the HF-style Llama3-8B checkpoint to NxDT supported format and store it in `pretrained_ckpt` directory. Modify the config parameter `exp_manager.resume_from_checkpoint` path to the converted pretrained checkpoint path.

## LoRA SFT-YAML Configuration Overview

You can configure a variety of SFT, DPO, PEFT-specific and model parameters for finetuning using the YAML file.

```
exp_manager:
  resume_from_checkpoint: /pretrained_ckpt

data:
  train_dir: /example_datasets/llama3_8b/training.jsonl
  val_dir: /example_datasets/llama3_8b/validation.json
  dev_choose_samples: 2250
  seq_length: 4096
  tokenizer:
    type: /llama3_tokenizer

model:
  weight_init_only: True

model_alignment_strategy:
  sft:
    packing: True
  peft:
    lora_rank: 16
    lora_alpha: 32
    lora_dropout: 0.05
    lora_bias: "none"
    lora_verbose: True
    target_modules: ["qkv_proj"]
```

### exp\_manager

#### resume\_from\_checkpoint

Manually set the checkpoint file (pretrained checkpoint) to load from

- **Type:** str
- **Default:** /pretrained\_ckpt
- **Required:** True (start with pretrained checkpoint)

### data

#### train\_dir

SFT training data - jsonl or arrow file

For SFT, we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** True

**val\_dir**

SFT validation data - jsonl or arrow file

For SFT, we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** False

**dev\_choose\_samples**

If set, will use that many number of records from the head of the dataset instead of using all.  
Set to null to use full dataset

- **Type:** integer
- **Default:** null
- **Required:** False

**seq\_length**

Set sequence length for the training job.

- **Type:** integer
- **Required:** True

**tokenizer  
type**

Set tokenizer path/type

- **Type:** str
- **Default:** /llama3\_tokenizer
- **Required:** True

**model****weight\_init\_only**

Load only model states and ignore the optim states from ckpt directory

- **Type:** bool
- **Default:** True

**model\_alignment\_strategy**

Set only when using finetuning specific algorithms (SFT, DPO, etc) and parameter-efficient fine-tuning methods like LoRA (Low-Rank Adaptation).

**sft**

Supervised Fine-Tuning (SFT) specific parameters.

**packing**

Appends multiple records in a single record until seq length supported by model, if false uses pad tokens to reach seq length. Setting it to True increases throughput but might impact accuracy.

- **Type:** bool
- **Default:** False
- **Required:** False

**peft**

Configuration options for Parameter-Efficient Fine-Tuning (PEFT) methods, specifically LoRA settings.

**lora\_rank**

Rank of LoRA; determines the number of trainable parameters. Higher rank allows for more expressive adaptations but increases memory usage.

- **Type:** int
- **Default:** 16
- **Required:** True

**lora\_alpha**

Scaling factor for LoRA updates; affects the magnitude of LoRA adaptations.

- **Type:** int
- **Default:** 32
- **Required:** True

**lora\_dropout**

Dropout rate for LoRA layers to prevent overfitting.

- **Type:** float
- **Default:** 0.05
- **Required:** False

**lora\_bias**

Bias type for LoRA. Determines which biases are trainable. Can be 'none', 'all' or 'lora\_only'.

- **Type:** str
- **Default:** "none"
- **Required:** False

**lora\_verbose**

Enables detailed LoRA-related logging during training.

- **Type:** bool
- **Default:** False
- **Required:** False

**target\_modules**

List of model layers to apply LoRA.

- **Type:** list[str]
- **Default:** ["qkv\_proj"] (for Llama)
- **Required:** True



## Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

```
git clone https://github.com/aws-neuron/neuronx-distributed-training
cd neuronx-distributed-training/examples
```

Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `hf_llama3_8B_SFT_lora_config`. The default config here is a 8B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

Next, run the following commands to launch an AOT pre-compilation job on your instance:

```
cd ~/neuronx-distributed-training/examples
export COMPILE=1
./train.sh
```

The compile output and logs will be shown directly in the terminal and you will see logs similar to this:

```
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total graphs: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total successful compilations: 22
2024-08-11 23:04:08.000738: INFO ||PARALLEL_COMPILE||: Total failed compilations: 0
```

Then, you know your compilation has successfully completed.

---

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

---

## Training the model

The fine-tuning job is launched almost exactly in the same way as the compile job. We now turn off the `COMPILE` environment variable and run the same training script to start pre-training.

On a single instance:

```
export COMPILE=0
./train.sh
```

Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 0%|          | 189/301501 [59:12<1573:03:24, 18.79s/it, loss=7.75, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=188.0, consumed_samples=24064.0]
Epoch 0: 0%|          | 190/301501 [59:30<1572:41:13, 18.79s/it, loss=7.74, v_num=3-16,
↳ reduced_train_loss=7.560, global_step=189.0, consumed_samples=24192.0]
Epoch 0: 0%|          | 191/301501 [59:48<1572:21:28, 18.79s/it, loss=7.73, v_num=3-16,
↳ reduced_train_loss=7.910, global_step=190.0, consumed_samples=24320.0]
```

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use standard tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/`. Once you have identified the directory, cd into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```

### neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, run `neuron-top`:

```
ssh compute1-dy-queue1-i1-1 # to determine which compute nodes are in use, run the
↳ squeue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Troubleshooting Guide

For issues with NxDT, please see: [NxDT Known Issues](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## HF Llama3.1/Llama3-8B Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO) based Fine-tuning (Beta)

In this example, we will show how to compile and finetune a pre-trained HF Llama3.1/Llama3-8B model on a single instance with the `NxD Training (NxDT)` library using [Direct Preference Optimization \(DPO\)](#) and [Odds Ratio Preference Optimization \(ORPO\)](#) based fine-tuning. The pre-trained Llama3-8B model serves as the foundation, and we will build upon this base by fine-tuning and aligning the model to adapt it to a specific task or dataset. The example has the following main sections:

### Table of contents

- *Setting up the environment*
  - *Install Dependencies*
- *DPO-YAML Configuration Overview*
- *ORPO-YAML Configuration Overview*
- *Download the dataset*
- *Convert data to DPO-specific Preference data format*
- *Download pretrained model checkpoint and tokenizer*
  - *Checkpoint Conversion*
- *Pre-compile the model*
- *Training the model*
- *Monitoring Training*
  - *Tensorboard monitoring*
  - *neuron-top / neuron-monitor / neuron-ls*
- *Troubleshooting Guide*

## Setting up the environment

### Install Dependencies

Once you have launched a Trn1 instance, Please follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxDT and its dependencies. Please see the following installation guide for installing NxDT: [NxDT Installation Guide](#).

We can download the `requirements_dpo.txt` and install using the command:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed-training/master/requirements_dpo.txt
pip install -r requirements_dpo.txt
```

## DPO-YAML Configuration Overview

You can configure a variety of DPO-specific and model parameters for finetuning through the YAML file.

```
exp_manager:
  resume_from_checkpoint: /pretrained_ckpt

data:
  train_dir: /example_datasets/llama3_8b/data_dpo.jsonl
  val_dir: null
  dev_choose_samples: null
  seq_length: 4096
  tokenizer:
    type: /llama3_tokenizer

model:
  weight_init_only: True

model_alignment_strategy:
  dpo:
    kl_beta: 0.01
    loss_type: sigmoid
    max_prompt_length: 2048
    precompute_ref_log_probs: True
    truncation_mode: keep_start
```

### exp\_manager

#### resume\_from\_checkpoint

Manually set the checkpoint file (pretrained/post-SFT checkpoint) to load from

- **Type:** str
- **Default:** /pretrained\_ckpt
- **Required:** True (start with pretrained checkpoint)

### data

#### train\_dir

DPO training data - jsonl or arrow file

As for DPO we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** True

#### val\_dir

DPO validation data - jsonl or arrow file

As for DPO we use HF style ModelAlignment dataloader, we also use HF style data file paths

- **Type:** str
- **Required:** False

#### dev\_choose\_samples

If set, will use that many number of records from the head of the dataset instead of using all. Set to null to use full dataset

- **Type:** integer
- **Default:** null
- **Required:** False

#### **seq\_length**

Set sequence length for the training job. For DPO, it is total sequence length of prompt and (chosen/rejected) response concatenated together

- **Type:** integer
- **Required:** True

#### **tokenizer type**

Set tokenizer path/type

- **Type:** str
- **Default:** /llama3\_tokenizer
- **Required:** True

#### **model**

##### **weight\_init\_only**

Load only model states and ignore the optim states from ckpt directory

- **Type:** bool
- **Default:** True

#### **model\_alignment\_strategy**

Set only when using finetuning specific algorithms (SFT, DPO, etc) and and parameter-efficient finetuning methods like LoRA (Low-Rank Adaptation).

##### **dpo**

Direct Preference Optimization (DPO) specific parameters.

##### **kl\_beta**

KL-divergence beta to control divergence of policy model from reference model

- **Type:** float
- **Default:** 0.01
- **Required:** True

##### **loss\_type**

Currently support sigmoid version of optimized DPO loss

- **Type:** str
- **Default:** sigmoid
- **Required:** True

**max\_prompt\_length**

Set maximum length of prompt in the concatenated prompt and (chosen/rejected) response input

- **Type:** integer
- **Required:** True

**precompute\_ref\_log\_probs**

To enable precomputation of reference model log probabilities using pre-fit hook, False is not supported currently

- **Type:** bool
- **Required:** True

**truncation\_mode**

To define how to truncate if size (prompt+response) exceeds seq\_length options: ["keep\_start", "keep\_end"]

- **Type:** str
- **Default:** keep\_start`
- **Required:** True

**ORPO-YAML Configuration Overview**

Here we show the ORPO-specific model parameters which can be configured for finetuning through the YAML file. And below we explain the parameters that are new as compared to DPO-specific parameters.

```
exp_manager:
  checkpoint_callback_params:
    every_n_train_steps: 10
    resume_from_checkpoint: /pretrained_ckpt

data:
  train_dir: /example_datasets/llama3_8b/data_orpo.jsonl
  val_dir: null
  dev_choose_samples: null
  seq_length: 4096
  tokenizer:
    type: /llama3_tokenizer

model:
  encoder_seq_len: 4096
  weight_init_only: True
  optim:
    lr: 1.5e-4
    sched:
      name: CosineAnnealing

model_alignment_strategy:
  orpo:
    beta: 0.1
```

(continues on next page)

(continued from previous page)

```
max_prompt_length: 2048
truncation_mode: keep_start
```

**exp\_manager****checkpoint\_callback\_params.every\_n\_train\_steps**

How often we want to checkpoint.

- **Type:** int
- **Required:** True

**model****encoder\_seq\_length**

Setting the sequence length for the training job. This parameter is common for all models supported in the library.

- **Type:** int
- **Required:** True

**optim.sched**

This is where the LR schedulers can be set. We can configure the schedulers supported by NeMo. All the schedulers can be configured according to the [parameters specified here](#).

- **Type:** config
- **Possible Values:** LinearAnnealingWithWarmUp, CosineAnnealing, WarmupPolicy, WarmupHoldPolicy, SquareAnnealing, NoamAnnealing, WarmupAnnealing,
- StepLR, rprop, ExponentialLR
- **Required:** True

**model\_alignment\_strategy**

Set only when using finetuning specific algorithms (SFT, DPO, ORPO, etc) and parameter-efficient fine-tuning methods like LoRA (Low-Rank Adaptation).

**orpo**

Odds Ratio Preference Optimization (ORPO) specific parameters.

**beta**

KL-divergence beta to control divergence of policy model from reference model

- **Type:** float
- **Default:** 0.01
- **Required:** True

## Download the dataset

The DPO (& ORPO) tutorial makes use of the same preprocessed version of *intel-orca\_dpo\_pairs* preference dataset that is stored in S3. The dataset can be downloaded to your cluster or instance by running the following AWS CLI commands on the head node or your Trn1 instance:

```
export DATA_DIR=~/examples_datasets/llama3_8b
mkdir -p ${DATA_DIR} && cd ${DATA_DIR}
aws s3 cp s3://neuron-s3/training_datasets/llama/dpo/data_dpo.jsonl . --no-sign-request
```

Then, download the config.json file:

For Llama-3.1-8B:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪ training/llama/tp_zero1_llama_hf_pretrain/8B_config_llama3.1/config.json ~/
```

For Llama-3-8B:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪ training/llama/tp_zero1_llama_hf_pretrain/8B_config_llama3/config.json ~/
```

## Convert data to DPO-specific Preference data format

If you directly downloaded the [Intel ORCA\\_dpo\\_pairs dataset](#), then you need to convert the data into preference dataset format using the script below.

**Note:** For different datasets with different field names, make necessary changes to the script accordingly.

```
from datasets import load_dataset
from transformers import AutoTokenizer

def preference_data_format(example):

    system = "<|im_start|>\n" + example['system'] + "<|im_end|>\n"

    # Format instruction
    prompt = "<|im_start|> " + example['question'] + "<|im_end|>\n<|im_start|>assistant\n"
    ↪ ""

    # Format chosen answer
    chosen = example['chosen'] + "<|im_end|>\n"

    # Format rejected answer
    rejected = example['rejected'] + "<|im_end|>\n"

    return {
        "prompt": system + prompt,
        "chosen": chosen,
        "rejected": rejected,
    }
```

(continues on next page)



(continued from previous page)

```
# Particular dataset with following fields: "system", "question", "chosen", "rejected"
dataset = load_dataset("json", data_files="orca_rlhf.jsonl", split="train")

# Save columns
original_columns = dataset.column_names

# Format dataset
dataset = dataset.map(
    preference_data_format,
    remove_columns=original_columns
)

# save converted preference dataset
dataset.to_json("data_dpo.jsonl")
```

## Download pretrained model checkpoint and tokenizer

In this tutorial, we will use a pretrained Llama3-8B checkpoint (post-SFT checkpoint preferred) from the original repository. Follow the steps to download tokenizer and model checkpoint from the pretraining stage: <https://llama.meta.com/llama-downloads/>

Alternatively, the model checkpoint and tokenizer can also be downloaded from HuggingFace by following this [guide](#). You can also directly download and convert the HuggingFace model checkpoint using *Direct HuggingFace Model Conversion*.

Create a folder llama3\_tokenizer and copy the tokenizer contents to it.

Modify the following paths in YAML file based on your specific directory configuration:

1. model.model\_config
2. exp\_manager.resume\_from\_checkpoint
3. tokenizer.type
4. train\_dir and val\_dir

You can use your Llama model, pretrained checkpoint and tokenizer by modifying the hf\_llama3\_8B\_<DPO/ORPO>\_config.yaml file.

## Checkpoint Conversion

Follow this [Checkpoint Conversion Guide](#) to convert the HF-style Llama3-8B checkpoint to NxDT supported format and store it in pretrained\_ckpt directory. Modify the config parameter exp\_manager.resume\_from\_checkpoint path to the converted pretrained checkpoint path.

## Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `hf_llama3_8B_<DPO/ORPO>_config.yaml`. The default config here is a 8B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

Next, run the following commands to launch an AOT pre-compilation job on your instance:

```
export COMPILE=1
export CONF_FILE=hf_llama3_8B_<DPO/ORPO>_config
./train.sh
```

The compile output and logs will be shown directly in the terminal and you will see logs similar to this:

```
2024-10-24 18:49:49.000950: INFO | | NEURON_PARALLEL_COMPILE | | : Total graphs: 32
2024-10-24 18:49:49.000950: INFO | | NEURON_PARALLEL_COMPILE | | : Total successful_
↳ compilations: 32
2024-10-24 18:49:49.000950: INFO | | NEURON_PARALLEL_COMPILE | | : Total failed compilations:_
↳ 0
```

Then, you know your compilation has successfully completed.

---

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

---

## Training the model

The fine-tuning job is launched almost exactly in the same way as the compile job. We now turn off the `COMPILE` environment variable and run the same training script to start pre-training.

On a single instance:

```
export COMPILE=0
export CONF_FILE=hf_llama3_8B_<DPO/ORPO>_config
./train.sh
```

Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 5%|^a- | 3/62 [02:59<58:44, 0.02it/s, v_num=8-06, reduced_train_
↳loss=6.930, chosen_rewards=-0.81, rejected_rewards=-0.675, lr=2.73e-5, parameter_
↳norm=1.95e+3, global_step=1.000, consumed_samples=32.00, throughput=0.108, throughput_
↳peak=0.0677, gradient_norm=8.600]
Epoch 0: 6%|^a- < | 4/62 [03:24<49:27, 0.02it/s, v_num=8-06, reduced_train_
↳loss=6.790, chosen_rewards=-0.628, rejected_rewards=-0.64, lr=5.45e-5, parameter_
↳norm=1.95e+3, global_step=3.000, consumed_samples=64.00, throughput=0.181, throughput_
↳peak=0.146, gradient_norm=6.590]
Epoch 0: 8%|^a- $ | 5/62 [03:50<43:42, 0.02it/s, v_num=8-06, reduced_train_
↳loss=6.790, chosen_rewards=-0.628, rejected_rewards=-0.64, lr=5.45e-5, parameter_
↳norm=1.95e+3, global_step=3.000, consumed_samples=64.00, throughput=0.181, throughput_
↳peak=0.146, gradient_norm=6.590]
```

**Note:** The values in the above logs will differ based on config used, package versions, models, and other factors. This is just an example.

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use standard tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/`. Once you have identified the directory, `cd` into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/hf_llama3_8B/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```

### neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, run `neuron-top`:

```
ssh compute1-dy-queue1-i1-1 # to determine which compute nodes are in use, run the
↳squeue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Troubleshooting Guide

For issues with NxDT, please see: [NxDT Known Issues](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

## HuggingFace Llama3.1/Llama3-70B Pretraining

In this example, we will compile and train a HuggingFace Llama3.1/Llama3-70B model on multiple trn1 or newly launched trn2 instances using ParallelCluster with the NxDT Training (NxDT) library. The example has the following main sections:

### Table of contents

- *Setting up the environment*
  - *ParallelCluster Setup*
  - *Install Dependencies*
- *Download the dataset*
- *Pre-compile the model*
- *Training the model*
- *Monitoring Training*
  - *Tensorboard monitoring*
  - *neuron-top / neuron-monitor / neuron-ls*
- *Continual Pre-training with Downloaded Meta Model Weights*
  - *Download the model and convert the state\_dict to NxDT checkpoint format*
  - *Start the continual training job by loading converted checkpoints*
- *Pretraining with Context Paralellism*
- *Troubleshooting Guide*

## Setting up the environment

### ParallelCluster Setup

In this example, we will use 16 trn1.32xlarge instances or 8 trn2.48xlarge instances with ParallelCluster. Please follow the instructions here to create a cluster: [Train your model on ParallelCluster](#)

ParallelCluster automates the creation of trainium clusters, and provides the Slurm job management system for scheduling and managing distributed training jobs. Please note that the home directory on your ParallelCluster head node will be shared with all of the worker nodes via NFS.

## Install Dependencies

Once you have launched ParallelCluster, please follow this guide on how to install the latest Neuron packages: [PyTorch Neuron Setup Guide](#).

Next, we will need to install NxDT and its dependencies. Please see the following installation guide for installing NxDT: [NxDT Installation Guide](#)

## Download the dataset

Let's download training-data scripts for our experiments

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪training/llama/get_dataset.py
```

Then download config.json file:

For Llama-3.1-70B:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪training/llama/tp_pp_llama_hf_pretrain/70B_config_llama3.1/config.json ~/
```

For Llama-3-70B:

```
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/
↪training/llama/tp_pp_llama_hf_pretrain/70B_config_llama3/config.json ~/
```

To tokenize the data, we must request the tokenizer from Hugging Face and Meta by following the instructions at the following link: [HuggingFace Llama 3.1 70B Model](#).

Use of the Llama models is governed by the Meta license. In order to download the model weights and tokenizer, please visit the above website and accept their License before requesting access. After access has been granted, you may use the following python3 script along with your own hugging face token to download and save the tokenizer.

```
from huggingface_hub import login
from transformers import AutoTokenizer

login(token='your_own_hugging_face_token')

tokenizer = AutoTokenizer.from_pretrained('meta-llama/Meta-Llama-3.1-70B')
# For llama3 uncomment line below
# tokenizer = AutoTokenizer.from_pretrained('meta-llama/Meta-Llama-3-70B')

tokenizer.save_pretrained(".")
```

For Llama3.1/Llama3, make sure your base directory has the following files:

```
'./tokenizer_config.json', './special_tokens_map.json', './tokenizer.json'
```

Next, let's download and pre-process the dataset:

```
mkdir ~/examples_datasets/
python3 get_dataset.py --llama-version 3
```

*Note:* In case you see an error of the following form when downloading data: `huggingface_hub.utils._validators.HFValidationError: Repo id must be in the form 'repo_name' or 'namespace/repo_name'. Use 'repo_type' argument if needed.` This could be because of a stale cache. Try deleting the cache using:

```
sudo rm -rf ~/.cache/
```

## Pre-compile the model

By default, PyTorch Neuron uses a just in time (JIT) compilation flow that sequentially compiles all of the neural network compute graphs as they are encountered during a training job. The compiled graphs are cached in a local compiler cache so that subsequent training jobs can leverage the compiled graphs and avoid compilation (so long as the graph signatures and Neuron version have not changed).

An alternative to the JIT flow is to use the included `neuron_parallel_compile` command to perform ahead of time (AOT) compilation. In the AOT compilation flow, the compute graphs are first identified and extracted during a short simulated training run, and the extracted graphs are then compiled and cached using parallel compilation, which is considerably faster than the JIT flow.

First, clone the open-source `neuronx-distributed-training` library

```
git clone https://github.com/aws-neuron/neuronx-distributed-training
cd neuronx-distributed-training/examples
```

Now, ensure that you are using the proper config file in the `conf/` directory. In the `train.sh` file, ensure that the `CONF_FILE` variable is properly set to the config for the model you want to use. In our case, it will be `hf_llama3_70B_config.yaml` for training on `trn1` cluster, and `hf_llama3_70B_trn2_config.yaml` for `trn2`.

In this tutorial, we will train Llama3-70B model on multiple compute nodes. For training on `trn1`, please make sure `hf_llama3_70B_config` has the right configuration:

```
trainer:
  devices: 32
  num_nodes: 16
```

For pretraining on `trn2`, `hf_llama3_70B_trn2_config` would contain:

```
trainer:
  devices: 64
  lnc: 2 # default for trn2 workloads
  num_nodes: 8
```

On `trn2` instances, the configuration `lnc: 2` indicates that there is a 2-to-1 mapping between logical Neuron Core (`lnc`) and physical Neuron Core. Another supported configuration is `lnc: 1`, in which case each node would expose 128 logical devices.

The default config here is a 70B parameter model, but users can also add their own `conf/*.yaml` files and run different configs and hyperparameters if desired. Please see [Config Overview](#) for examples and usage for the `.yaml` config files.

On `trn1` cluster, run the following commands to launch an AOT pre-compilation job on your instance:

```
export COMPILER=1
export CONF_FILE=hf_llama3_70B_config
sbatch --exclusive \
  --nodes 16 \
```

(continues on next page)

(continued from previous page)

```
--cpus-per-task 128 \
--wrap="srun ./train.sh"
```

On trn2 cluster, run the following:

```
export COMPILE=1
export CONF_FILE=hf_llama3_70B_trn2_config
sbatch --exclusive \
  --nodes 8 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

Once you have launched the precompilation job, run the `squeue` command to view the Slurm job queue on your cluster. If you have not recently run a job on your cluster, it may take 4-5 minutes for the requested trn1.32xlarge or trn2.48xlarge nodes to be launched and initialized. Once the job is running, `squeue` should show output similar to the following:

| JOBID | PARTITION | NAME | USER   | ST | TIME | NODES | ODELIST (REASON)             |
|-------|-----------|------|--------|----|------|-------|------------------------------|
| 7     | compute1  | wrap | ubuntu | R  | 5:11 | 16    | compute1-st-queue1-i1-[1-16] |

You can view the output of the precompilation job by examining the file named `slurm-ZZ.out`, where `ZZ` represents the JOBID of your job in the `squeue` output above.

```
tail -f slurm-7.out
```

Once the precompilation job is complete, just like the above output you should see a message similar to the following in the logs:

```
2024-11-07 09:57:13.000144: 39810 INFO | NEURON_PARALLEL_COMPILE | : Total graphs: 36
2024-11-07 09:57:13.000144: 39810 INFO | NEURON_PARALLEL_COMPILE | : Total successful_
↪compilations: 36
2024-11-07 09:57:13.000144: 39810 INFO | NEURON_PARALLEL_COMPILE | : Total failed_
↪compilations: 0
```

At this point, you can press CTRL-C to exit the `tail` command.

---

**Note:** The number of graphs will differ based on package versions, models, and other factors. This is just an example.

---

## Training the model

You can launch pre-training job similar to compilation by using the same training script but now turning off the `COMPILE` environment variable

On trn1 ParallelCluster:

```
export COMPILE=0
export CONF_FILE=hf_llama3_70B_config
sbatch --exclusive \
  --nodes 16 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

On trn2 ParallelCluster:

```
export COMPILE=0
export CONF_FILE=hf_llama3_70B_trn2_config
sbatch --exclusive \
    --nodes 8 \
    --cpus-per-task 128 \
    --wrap="srun ./train.sh"
```

As outlined above, you can again use the `squeue` command to view the job queue, and also monitor the job in the same way with the `tail` command to see the training logs. Once the model is loaded onto the Trainium accelerators and training has commenced, you will begin to see output indicating the job progress:

Example:

```
Epoch 0: 3%|          | 3/91 [16:05<7:52:06, 321.89s/it, loss=6.7, v_num=2, reduced_
↪train_loss=13.40, lr=7.5e-9, parameter_norm=5536.0, global_step=1.000, consumed_
↪samples=2048.0]
Epoch 0: 3%|          | 3/91 [16:05<7:52:06, 321.89s/it, loss=4.47, v_num=2, reduced_
↪train_loss=13.40, lr=7.5e-9, parameter_norm=5536.0, global_step=2.000, consumed_
↪samples=3072.0]
Epoch 0: 4%|          | 4/91 [21:20<7:44:18, 320.22s/it, loss=4.47, v_num=2, reduced_
↪train_loss=13.40, lr=7.5e-9, parameter_norm=5536.0, global_step=2.000, consumed_
↪samples=3072.0]
Epoch 0: 4%|          | 4/91 [21:20<7:44:18, 320.22s/it, loss=3.35, v_num=2, reduced_
↪train_loss=13.40, lr=7.5e-9, parameter_norm=5536.0, global_step=3.000, consumed_
↪samples=4096.0]
```

---

**Note:** The convergence is for demonstration and would differ based on instance type, model, and other factors.

---

## Monitoring Training

### Tensorboard monitoring

In addition to the text-based job monitoring described in the previous section, you can also use tools such as TensorBoard to monitor training job progress. To view an ongoing training job in TensorBoard, you first need to identify the experiment directory associated with your ongoing job. This will typically be the most recently created directory under `~/neuronx-distributed-training/examples/nemo_experiments/hf_llama/`. Once you have identified the directory, `cd` into it, and then launch TensorBoard:

```
cd ~/neuronx-distributed-training/examples/nemo_experiments/hf_llama/8/
tensorboard --logdir ./
```

With TensorBoard running, you can then view the TensorBoard dashboard by browsing to `http://localhost:6006` on your local machine. If you cannot access TensorBoard at this address, please make sure that you have port-forwarded TCP port 6006 when SSH'ing into the head node,

```
ssh -i YOUR_KEY.pem ubuntu@HEAD_NODE_IP_ADDRESS -L 6006:127.0.0.1:6006
```



## neuron-top / neuron-monitor / neuron-ls

The `neuron-top` tool can be used to view useful information about NeuronCore utilization, vCPU and RAM utilization, and loaded graphs on a per-node basis. To use `neuron-top` during an ongoing training job, run `neuron-top`:

```
ssh compute1-st-queue1-i1-1 # to determine which compute nodes are in use, run the
↪squeue command
neuron-top
```

Similarly, once you are logged into one of the active compute nodes, you can also use other Neuron tools such as `neuron-monitor` and `neuron-ls` to capture performance and utilization statistics and to understand NeuronCore allocation.

## Continual Pre-training with Downloaded Meta Model Weights

If you want to perform continual pre-training using the model weights provided by Meta, follow these steps:

Ensure you have the `config.json` file, which should have been downloaded as described in the [Download the dataset](#) section.

### Download the model and convert the state\_dict to NxDT checkpoint format

Get the conversion scripts described in the [Checkpoint Conversion](#). Mention the `hf_model_name` argument to specify the HuggingFace model identifier for the model you want to download and convert the checkpoint to NxDT format.

Run the following to download the model and convert the `state_dict` to NxDT sharded checkpoint.

On `trn1` cluster:

```
python3 ./checkpoint_converter_scripts/checkpoint_converter.py \
  --model_style hf \
  --hf_model_name meta-llama/Meta-Llama-3-70B \
  --hw_backend trn1 \
  --tp_size 32 --pp_size 8 --n_layers 80 \
  --output_dir /fsx/pretrained_weight/ \
  --convert_from_full_state --save_xser True \
  --kv_size_multiplier 4 --qkv_linear True \
  --config ~/config.json
```

On `trn2` cluster:

```
python3 ./checkpoint_converter_scripts/checkpoint_converter.py \
  --model_style hf \
  --hf_model_name meta-llama/Meta-Llama-3-70B \
  --hw_backend trn2 \
  --tp_size 32 --pp_size 4 --n_layers 80 \
  --output_dir /fsx/pretrained_weight/ \
  --convert_from_full_state --save_xser True \
  --kv_size_multiplier 4 --qkv_linear True \
  --config ~/config.json
```

**Note:** This conversion process requires larger host memory. Please run it on a `trn1.32xlarge` or `trn2.48xlarge` compute node. In this example, the converted model is stored on FSx for Lustre to be accessed by all compute nodes.

## Start the continual training job by loading converted checkpoints

In order to start the continual training job with loading this converted model as initial weights, please update the config file (hf\_llama3\_70B\_config.yaml or hf\_llama3\_70B\_trn2\_config.yaml) as below:

```
exp_manager:
.
.
  resume_from_checkpoint: /fsx/pretrained_weight/ # manually set the checkpoint file to_
↪load from
.
.
model:
  # Miscellaneous
  use_cpu_initialization: False # Init weights on the CPU (slow for large models)
  weight_init_only: True
```

Compared to initial pre-training loss value, you should see lower loss value when the training starts with Meta's model weights. Logs for one such sample run look like below.

```
Epoch 0: 3%|          | 3/91 [16:09<7:53:59, 323.17s/it, loss=0.834, v_num=7, reduced_
↪train_loss=1.670, lr=7.5e-9, parameter_norm=4736.0, global_step=1.000, consumed_
↪samples=2048.0]
Epoch 0: 3%|          | 3/91 [16:09<7:53:59, 323.17s/it, loss=0.556, v_num=7, reduced_
↪train_loss=1.670, lr=7.5e-9, parameter_norm=4736.0, global_step=2.000, consumed_
↪samples=3072.0]
Epoch 0: 4%|          | 4/91 [21:25<7:46:02, 321.41s/it, loss=0.556, v_num=7, reduced_
↪train_loss=1.670, lr=7.5e-9, parameter_norm=4736.0, global_step=2.000, consumed_
↪samples=3072.0]
Epoch 0: 4%|          | 4/91 [21:25<7:46:02, 321.41s/it, loss=0.417, v_num=7, reduced_
↪train_loss=1.670, lr=7.5e-9, parameter_norm=4736.0, global_step=3.000, consumed_
↪samples=4096.0]
```

## Pretraining with Context Parallelism

To run pretraining with context parallelism, use the following yaml config file: hf\_llama3\_70B\_CP\_config.yaml. This YAML file has the following changes to enable context parallelism:

```
distributed_strategy:
  context_parallel_size: 2

fusions:
  flash_attention: False
  ring_attention: True
```

**distributed\_strategy**  
**context\_parallel\_size**

Context parallel degree to be used for sharding sequence.

- **Type:** int
- **Required:** False
- **Default:** 1

**fusions****ring\_attention**

Setting this flag to True will use the ring attention module for both forward and backward. This parameter must be true when context parallel is `context\_parallel\_size` is greater than 1.

- **Type:** bool
- **Required:** False

In the config file, `context_parallel_size` is set to the desired degree, and as context parallelism leverages ring attention instead of flash attention, we set `ring_attention: True`, and `flash_attention: False`.

Context parallelism currently supports sequence lengths up to 32k and is supported on TRN1.

Compile with:

```
export COMPILE=1
export CONF_FILE=hf_llama3_70B_CP_config
sbatch --exclusive \
  --nodes 16 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

and pre-training with:

```
export COMPILE=0
export CONF_FILE=hf_llama3_70B_CP_config
sbatch --exclusive \
  --nodes 16 \
  --cpus-per-task 128 \
  --wrap="srun ./train.sh"
```

**Troubleshooting Guide**

For issues with NxDT, please see: [NxDT Known Issues](#)

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

**Checkpoint Conversion****Table of Contents**

- *Supported Model Architectures*
- *Conversion Scenarios and Usage*
  - *Key Arguments*
- *Conversion Example*
- *Direct HuggingFace Model Conversion*
  - *Troubleshooting*

The NxD Training library provides a versatile checkpoint conversion functionality, allowing seamless transition between different model styles. This tutorial aims to provide a comprehensive guide through the various use cases and demonstrate how to perform the checkpoint conversions.

## Supported Model Architectures

The checkpoint conversion functionality supports conversion of the following model styles to/from NxDT checkpoints:

1. **HuggingFace (HF) style models**
2. **Megatron style models**

Extends support for both GQA (Llama-3) and non-GQA models (Llama-2).

## Conversion Scenarios and Usage

The tool supports the following conversion scenarios. It internally uses `NeuronxDistributed` (NxD) to convert to/from checkpoints. Run the following commands from the `/examples/checkpoint_conversion_scripts/` directory:

---

### Note:

1. **Important:** You must set the `--hw_backend` argument correctly for your hardware. The sample commands below use `trn1`.
  - Set `--hw_backend trn1` for Trainium (Trn1) hardware
  - Set `--hw_backend trn2` for Trainium 2 (Trn2) hardware

All example commands in this tutorial use `trn1`. If you're using Trn2, remember to replace `trn1` with `trn2` in every command.

2. Ensure that the model configuration `config.json` file is present, as it is required for checkpoint conversions. It is suggested to use specific json files like [examples](#) . If not present, you will need to create it.
  3. If your HF/custom checkpoint has multiple `.bin` or `.pt` or `.pth` files then merge and convert to a single file before conversion.
- 

For conversion of non-GQA based models (e.g. Llama2), just set the `--qkv_linear` argument to `False`.

1. **HF style model:**

- a. **HF to NxDT checkpoint:**

**Command:**

```
python3 checkpoint_converter.py --model_style hf --hw_backend trn1 --input_dir /
↪home/ubuntu/pretrained_llama_3_8B_hf/pytorch_model.bin --output_dir /home/
↪ubuntu/converted_hf_style_hf_to_nxdt_tp8pp4/ --save_xser True --config /home/
↪ubuntu/pretrained_llama_3_8B_hf/config.json --tp_size 8 --pp_size 4 --n_
↪layers 32 --kv_size_multiplier 1 --qkv_linear True --convert_from_full_state
```

This converts an HF-style checkpoint to an NxDT checkpoint.

- b. **NxDT to HF checkpoint:**

**Command:**

```
python3 checkpoint_converter.py --model_style hf --hw_backend trn1 --input_dir ~/examples/nemo_experiments/hf_llama3_8B_SFT/2024-07-19_23-07-40/checkpoints/hf_llama3_8B--step=5-consumed_samples=160.0.ckpt/model --output_dir ~/converted_hf_style_nxdt_to_hf_tp8pp4/ --load_xser True --config ~/config.json --tp_size 8 --pp_size 4 --kv_size_multiplier 1 --qkv_linear True --convert_to_full_state
```

This converts an NxDT checkpoint to an HF-style checkpoint.

## 2. Megatron style model (non-GQA models: e.g., Llama-2, and GQA models: e.g., Llama-3):

### a. HF to NxDT Megatron checkpoint:

**Command:**

```
python3 checkpoint_converter.py --model_style megatron --hw_backend trn1 --input_dir ~/megatron-tp8pp4-nxdt-to-hf4/checkpoint.pt --output_dir ~/meg_nxdt_hf3_nxdt3 --config ~/llama_gqa/config.json --save_xser True --tp_size 8 --pp_size 4 --n_layers 32 --kv_size_multiplier 1 --qkv_linear True --convert_from_full_state
```

This converts an HF-style checkpoint to an NxDT Megatron-style checkpoint.

### b. NxDT Megatron checkpoint to HF:

**Command:**

```
python3 checkpoint_converter.py --model_style megatron --hw_backend trn1 --input_dir ~/examples/nemo_experiments/megatron_llama/2024-07-23_21-07-30/checkpoints/megatron_llama--step=5-consumed_samples=5120.0.ckpt/model --output_dir ~/megatron-tp8pp4-nxdt-to-hf4 --load_xser True --config ~/llama_gqa/config.json --tp_size 8 --pp_size 4 --kv_size_multiplier 1 --qkv_linear True --convert_to_full_state
```

This converts an NxDT Megatron-style checkpoint to an HF-style checkpoint (GQA-based model, see: `--qkv_linear` set to True).

## Key Arguments

The `checkpoint_converter.py` script supports the following key arguments:

- `--model_style`: Specifies the model style, either *hf* (HuggingFace: default) or *megatron*
- `--hw_backend`: (required) Specifies the hardware backend either *trn1* or *trn2*
- `--input_dir`: (required) directory containing the input checkpoint
- `--hf_model_name`: (optional) HuggingFace model identifier for directly converting models hosted on HuggingFace
- `--output_dir`: (required) directory to save the converted checkpoint directory
- `--save_xser`: Saves the checkpoint with `torch_xla` serialization
- `--load_xser`: Loads the checkpoint with `torch_xla` serialization
- `--convert_from_full_state`: Converts full model checkpoint to sharded model checkpoint
- `--convert_to_full_state`: Converts sharded model checkpoint to full model checkpoint

- `--config`: path to the model configuration file (create *json* file if not present)
- `--tp_size`: tensor parallelism degree
- `--pp_size`: pipeline parallelism degree
- `--n_layers`: number of layers in the model
- `--kv_size_multiplier`: key-value size multiplier
- `--qkv_linear`: boolean to specify GQA/non-GQA models
- `--fuse_qkv`: boolean to specify fused QKV in GQA models

We recommend enabling xser for significantly faster save and load times. Note that if the checkpoint is saved with xser, it can only be loaded with xser, and vice versa.

### Conversion Example

Assuming you have a pre-trained HF-style Llama3-8B model checkpoint looking similar to:

`input_dir: /hf/checkpoint/pytorch_model.bin`

```
$ ls /hf/checkpoint
-rw-r--r-- 1 user group 123 Aug 27 2024 pytorch_model.bin
```

Convert the HF-style checkpoint to an NxDT checkpoint on a single instance:

```
python3 checkpoint_converter.py --model_style hf --hw_backend trn1 --input_dir /hf/
↪checkpoint/pytorch_model.bin --output_dir /nxd/checkpoint --save_xser True --convert_
↪from_full_state --config /path/to/config.json --tp_size 8 --pp_size 4 --n_layers 32 --
↪kv_size_multiplier 1 --qkv_linear True --convert_from_full_state
```

This command will create an NxDT checkpoint in `output_dir: /nxd/checkpoint` and it will be sharded with (tp=8, pp=4) like:

```
$ ls /nxd/checkpoint/model
-rw-r--r-- 1 user group 123 Aug 27 2024 dp_rank_00_tp_rank_00_pp_rank_00.pt
-rw-r--r-- 1 user group 456 Aug 27 2024 dp_rank_00_tp_rank_01_pp_rank_00.pt
.....
-rw-r--r-- 1 user group 789 Aug 27 2024 dp_rank_00_tp_rank_07_pp_rank_02.pt
-rw-r--r-- 1 user group 122 Aug 27 2024 dp_rank_00_tp_rank_07_pp_rank_03.pt
```

### Direct HuggingFace Model Conversion

Using the `--hf_model_name` argument allows users to directly convert checkpoint files hosted on HuggingFace without the need for manual downloading or merging of checkpoint files.

To use this feature, you can specify the HuggingFace model identifier using the `--hf_model_name` argument. The script will then download the model and convert it directly to the NxDT format.

---

#### Note:

1. When using `--hf_model_name`, do not specify `--input_dir`. These arguments are mutually exclusive.

2. If both `--hf_model_name` and `--input_dir` are specified, the script will prioritize `--input_dir` and ignore `--hf_model_name`
3. You will be prompted to enter your HuggingFace API token. If you don't have one, you can create it at <https://huggingface.co/settings/tokens>.
4. Ensure you have sufficient disk space to download and process the model.

Example usage:

```
python3 checkpoint_converter.py --model_style hf --hw_backend trn1 --hf_model_name "meta-llama/Llama-2-7b-hf" --output_dir /path/to/output --save_xser True --config /path/to/config.json --tp_size 8 --pp_size 4 --n_layers 32 --kv_size_multiplier 1 --qkv_linear_ False --convert_from_full_state
```

This command will download the Llama-2-7b model from HuggingFace. Convert it to NxD format, and save it in the specified output directory.

## Troubleshooting

- If you encounter an error related to HuggingFace authentication, ensure you're using a valid API token.
- If the download fails, check your internet connection and verify that the model identifier is correct.

*This document is relevant for: Trn1, Trn2*

- *Megatron GPT Pretraining*
- *HuggingFace Llama3.1/Llama3-8B Pretraining*
- *HuggingFace Llama3.1/Llama3-8B Supervised Fine-tuning*
- *HuggingFace Llama3.1/Llama3-8B Efficient Supervised Fine-tuning with LoRA (Beta)*
- *HF Llama3.1/Llama3-8B Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO) based Fine-tuning (Beta)*
- *HuggingFace Llama3.1/Llama3-70B Pretraining*
- *Checkpoint Conversion*

*This document is relevant for: Trn1, Trn2*

*This document is relevant for: Trn1, Trn2*

### 3.1.7 Misc.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

**NxD Training Release Notes (neuronx-distributed-training)****Table of contents**

- *NeuronX Distributed Training [1.4.1]*
- *NeuronX Distributed Training [1.4.0]*
- *NeuronX Distributed Training [1.3.0]*
- *NeuronX Distributed Training [1.2.0]*
- *NeuronX Distributed Training [1.1.1]*
- *NeuronX Distributed Training [1.1.0]*
- *NeuronX Distributed Training [1.0.1]*
- *NeuronX Distributed Training [1.0.0]*

This document lists the release notes for Neuronx Distributed Training library.

**NeuronX Distributed Training [1.4.1]**

Date: 06/30/2025

Features in this release

- Fixed an installation issue.

**NeuronX Distributed Training [1.4.0]**

Date: 06/24/2025

Features in this release

- Added support for PyTorch 2.7.

**NeuronX Distributed Training [1.3.0]**

Date: 05/16/2025

Features in this release

- (Beta release) Added autocast for HF based Llama3 8B and Llama3 70B models
- (Beta release) Added support for context parallel sequence lengths up to 32k on TRN1
- Added support for ORPO
- Added support for nemo-toolkit 2.1
- Added support for Transformers 4.48.0
- Added support for PyTorch-Lightning 2.5.0
- Added support for PyTorch 2.6



## NeuronX Distributed Training [1.2.0]

Date: 04/03/2025

Features in this release

- Added support for LoRA supervised fine-tuning.
- Added option to configure collectives data types.
- Minor fixes to reduce the amount of logs during training.
- Removes `-llm-training` flag by default in all configs, except llama2. Note: this flag should not be enabled when using the Neuron Kernel Interface.

## NeuronX Distributed Training [1.1.1]

Date: 1/14/2025

Features in this release

- Added a flag in Llama3/3.1 70B config to control the dtype of reduce-scatter operations in Column/Row Parallel linear layers.

## NeuronX Distributed Training [1.1.0]

Date: 12/20/2024

Features in this release

- Added support for HuggingFace Llama3/3.1 70B with trn2 instances
- Added support for custom pipeline parallel cuts in HuggingFace Llama3
- Added support for PyTorch 2.5
- Added support for DPO post-training model alignment
- Added support for Mixtral 8x7B Megatron and HuggingFace models
- Added option in checkpoint converter to download and convert checkpoints using HuggingFace model identifier
- Fix the validation loss to properly compute the average loss across the validation epoch
- Minor bug fixes for error logging and imports

## Known Issues and Limitations

- Autocast option may not properly cast all inputs to bf16, recommended to use mixed precision option (currently is default) in configs for best results
- With PT2.5, some of the key workloads like Llama3-8B training may show a reduced performance when using `-llm-training` compiler flag as compared to PT2.1.

In such a case, try removing `-llm-training` flag from `compiler_flags` in the config.yaml

## NeuronX Distributed Training [1.0.1]

Date: 11/20/2024

Features in this release

- Added support for transformers 4.36.0

## NeuronX Distributed Training [1.0.0]

Date: 09/16/2024

Features in this release

This is the first release of NxDT Training (NxDT), NxDT is a PyTorch-based library that adds support for user-friendly distributed training experience through a YAML configuration file compatible with NeMo, allowing users to easily set up their training workflows. At the same time, NxDT maintains flexibility, enabling users to choose between using the YAML configuration file, PyTorch Lightning Trainer, or writing their own custom training script using the NxDT Core. The library supports PyTorch model classes including Hugging Face and Megatron-LM. Additionally, it leverages NeMo's data engineering and data science modules enabling end-to-end training workflows on NxDT, and providing a compatibility with NeMo through minimal changes to the YAML configuration file for models that are already supported in NxDT. Furthermore, the functionality of the Neuron NeMo Megatron (NNM) library is now part of NxDT, ensuring a smooth migration path from NNM to NxDT.

This release of NxDT includes:

- Installation through *neuronx-distributed-training* package.
- Open Source Github repository: <https://github.com/aws-neuron/neuronx-distributed-training>
- Support for YAML based interface allowing users to configure training from a config file.
- Support for 3D-parallelism, sequence-parallelism and zero1.
- Support for megatron-model and hugging-face based Llama model.
- Support flash attention kernel.
- Support for async checkpointing and s3 checkpointing.
- Examples to pretrain and fine-tune Llama model

## Known Issues and Limitations

- Model checkpointing saves sharded checkpoints. Users will have to write a script to combine the shards
- Validation/Evaluation with interleaved pipeline feature is not supported.
- NxDT shows slightly higher memory utilization as compared to NxDT based examples.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Trn1, Trn2

## Known Issues and Workarounds

This section covers the common failures that one can see while working with Neuronx Distributed Training library. Some of the failures regarding installation have been documented in [Common failures during installation](#).

### Table of contents

- *Shared weights error*
- *HOST OOM issues*
  - *During checkpoint saving*
  - *During async\_checkpointing*
  - *During Dataloading*
- *ImportError: helpers*
- *Matplotlib error*
- *Flash Attention not supported for megatron-style models*

### Shared weights error

Tieing weights is not supported when using pipeline parallelism. This means currently, the `share_embeddings_and_output_weights` parameter is not supported when using pipeline parallelism. It would produce an error that looks like this

```
File "/home/ubuntu/aws_neuron_venv_pytorch/lib/python3.8/site-packages/neuronx_
distributed/pipeline/model.py", line 625, in _reduce_shared_weights
assert p.grad is not None, f"Found shared weight {n} has None grad"
AssertionError: Found shared weight language_model_embedding_word_embeddings.weight has_
None grad
```

Please set this flag to False when using pipeline parallelism.

### HOST OOM issues

You would see an error log that looks like this without any other error above it.

```
WARNING:torch.distributed.elastic.agent.server.api:Received 15 death signal, shutting_
down workers
WARNING:torch.distributed.elastic.multiprocessing.api:Sending process 3721028 closing_
signal SIGTERM
WARNING:torch.distributed.elastic.multiprocessing.api:Sending process 3721029 closing_
signal SIGTERM
WARNING:torch.distributed.elastic.multiprocessing.api:Sending process 3721030 closing_
signal SIGTERM
WARNING:torch.distributed.elastic.multiprocessing.api:Sending process 3721031 closing_
signal SIGTERM
```

You can confirm HOST OOM by checking `sudo dmesg` on the Trn1 node. HOST OOM can occur because of multiple reasons:

## During checkpoint saving

If you see the above error immediately after a checkpoint saving log, this indicates that the entire checkpoint is copied to CPU. In this case, please check if the `save_xser` parameter is set to `True`. This mode will ensure each worker saves only one tensor at a time to disk. Setting this to `False` will make all the workers copy the entire checkpoint to CPU and can result in `HOST OOM`.

## During `async_checkpointing`

`async_checkpointing` when used with a low number of nodes can cause `HOST OOM` as it increases memory pressure per node. When we use more nodes, the memory pressure gets divided among the nodes and hence you would get an `OOM`.

On a high level, `async_checkpointing` copies data from device memory to host memory, then launch a new process to save host memory to storage, and let the main process continue with the training. Since we launch a new process, it requires a lot more extra host memory, because the launched process has the exact copy of memory space of the parent process. Let's use the following example to demonstrate how much memory we would need. For a llama2 70b training using `tp32` on 32 nodes, we launch 32 processes on each node. As baseline, each process uses 5 GB of host memory. There is also the XRT server, which uses 110 GB of host memory, so in total 270 GB host memory is used ( $5 \times 32 + 110$ ). If we enable `async_checkpointing` on this setting, the final memory usage can reach as high as 482 GB because of the following reasons:

1. Each training process needs to allocate memory to hold the model. The model weights for llama2 70B would require 280GB of memory to store the weights. The optimizer state would require twice as much memory. So total amount of host memory is 840 GB. Because we used all ranks for saving, the 840GB of data was evenly distributed among 1,024 processes ( $32 \times 32$ ), which means 0.84 GB of memory per process, or 26 GB of memory per instance. So each process's host memory usage is 5.8GB.
2. Second, each training process will fork a process for saving. The forked process will have a copy of parent's memory. In practice, linux uses a Copy-On-Write mechanism to save memory usage, but still in theory the actual memory usage of the child process can reach 5.8 GB combined. When `async_checkpointing` is enabled, we have 64 processes each using 5.8 GB of memory, and the XRT server uses 110 GB of memory. Therefore the total memory usage will be 482GB ( $64 \times 5.8 + 110$ ).

Hence with 32 nodes, we are already on the edge (each Trn1 node has 512GB of host memory) and we could `OOM` at 32 nodes. For a more stable run, enabling `async_checkpointing` at 64 nodes is recommended.

## During Dataloading

Another common reason for `HOST OOM` is loading too much data onto CPU. For pipeline-parallel processing, the library loads the entire global batch onto CPU and then moves it one-by-one to device. If we have a large batchsize with each batch taking space, it can lead to `HOST OOM`.

## ImportError: helpers

If you see an error that looks like:

```
ImportError: cannot import name 'helpers' from 'nemo.collections.nlp.data.language_
↪ modeling.megatron' (/usr/local/lib/python3.8/dist-packages/nemo/collections/nlp/data/
↪ language_modeling/megatron/__init__.py)
```

This could be because the `helpers.cpp` didn't get built correctly at the time of execution. We can pre-built it by running the following code:

```

import sys
import types

import torch

if torch.__version__.startswith("2"):
    string_classes = str
    inf = torch.inf
else:
    string_classes = None
    inf = None

# conditionally modify the import
def modify_torch_six_import():
    if string_classes is not None:
        try:
            if "torch._six" not in sys.modules:
                # Create and add dummy module to sys.modules
                six_module = types.ModuleType("torch._six")
                six_module.string_classes = string_classes
                six_module.inf = inf
                sys.modules["torch._six"] = six_module
        except Exception as e:
            raise RuntimeError(f"Failed to override torch._six import: {e}")

modify_torch_six_import()
from nemo.collections.nlp.data.language_modeling.megatron.dataset_utils import compile_
    ↪ helper
compile_helper()

```

Alternatively, if you see

```

ImportError: /shared/username/aws_neuron_venv_pytorch/lib/python3.10/site-packages/nemo/
    ↪ collections/nlp/data/language_modeling/megatron/helpers.cpython-310-x86_64-linux-gnu.
    ↪ so: file too short

```

A current workaround for this case is to delete the .so file and run the above snippet explicitly.

## Matplotlib error

If you see an error that looks like:

```

TimeoutError: Lock error: Matplotlib failed to acquire the following lock file

```

It means there is some contention in compute/worker nodes to access the matplotlib cache, and hence the lock error. To resolve this add or run `python -c 'import matplotlib.pyplot as plt'` as part of your setup. This will create a matplotlib cache and avoid the race condition.

## Flash Attention not supported for megatron-style models

Flash attention kernel is supported only for HF-style models and will be added for megatron-style models in one of the future releases.

*This document is relevant for: Trn1, Trn2*

- *NxD Training Release Notes (neuronx-distributed-training)*
- *Known Issues and Workarounds*

*This document is relevant for: Trn1, Trn2*

NxD Training is a PyTorch library for end-to-end distributed training.

## NxD Training Overview & Setup

- *Overview*
- *Setup*

## API Reference Guide

- *YAML Configuration Settings*

## Developer Guide

- *Integrating a New Model*
- *Integrating a new dataset/dataloader*
- *Registering an optimizer and LR scheduler*
- *Migrating from Neuron-NeMo-Megatron to Neuronx Distributed Training*
- *NxD Training Compatibility with NeMo*

## Tutorials

- *Megatron GPT Pretraining*
- *HuggingFace Llama3.1/Llama3-8B Pretraining*
- *HuggingFace Llama3.1/Llama3-8B Supervised Fine-tuning*
- *HuggingFace Llama3.1/Llama3-8B Efficient Supervised Fine-tuning with LoRA (Beta)*
- *HF Llama3.1/Llama3-8B Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO) based Fine-tuning (Beta)*
- *HuggingFace Llama3.1/Llama3-70B Pretraining*
- *Checkpoint Conversion*

## App Notes

- [Introducing NxD Training](#)
- [Tensor Parallelism Overview](#)
- [Pipeline Parallelism Overview](#)
- [Activation Memory Reduction](#)

## Misc

- [Known Issues and Workarounds](#)

This document is relevant for: Trn1, Trn2

## 3.2 NxD Inference

### 3.2.1 Overview

- [NxD Inference Overview](#)

## NxD Inference Overview

### Table of contents

- [Overview](#)
- [Using NxD Inference Library](#)

## Overview

NxD Inference (where NxD stands for NeuronX Distributed) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. Introduced with Neuron SDK 2.21 release, it offers advanced inference capabilities, including features such as continuous batching and speculative decoding for high performance inference. Additionally, it supports inference engine for vLLM for seamless integration with the majority of customers' production deployment systems. ML developers can use NxD Inference library at different levels of abstraction that fits their inference use case.

NxD Inference(NxDI) library offers the following benefits:

- **Production ready models:** NxD Inference provides production ready models like Llama-3.1, DBRX, and Mixtral that you can quickly deploy for high performance inference.
- **LLM Inference Features:** NxD Inference provides support for various LLM inference features like KV Cache, Multi-Head Attention (MHA), Grouped Query Attention (GQA), Flash Attention, Quantization, MoE , Continuous Batching and Speculative Decoding enabling high performance inference.
- **Modular Design:** Inference features in NxDI like KV Caching are implemented with a modular design, allowing developers to easily incorporate them into new models or customize and extend them.

- **Distributed Strategies:** NxID Inference enables distributing inference workload of large models across multiple NeuronCores in a single instance using Tensor parallelism and Sequence Parallelism. Pipeline parallelism and multi-node inference will be supported in future Neuron releases.
- **Support for NKI Kernels:** NxID Inference provides support for integrating custom NKI kernels on Trainium and Inferentia instances.
- **Open Source and SW Release:** NxID Inference library is provided as pip wheel and corresponding source code is made available on [GitHub](#) . We encourage developers to contribute new model implementations or feature optimizations to the NxDI library by submitting a pull request.

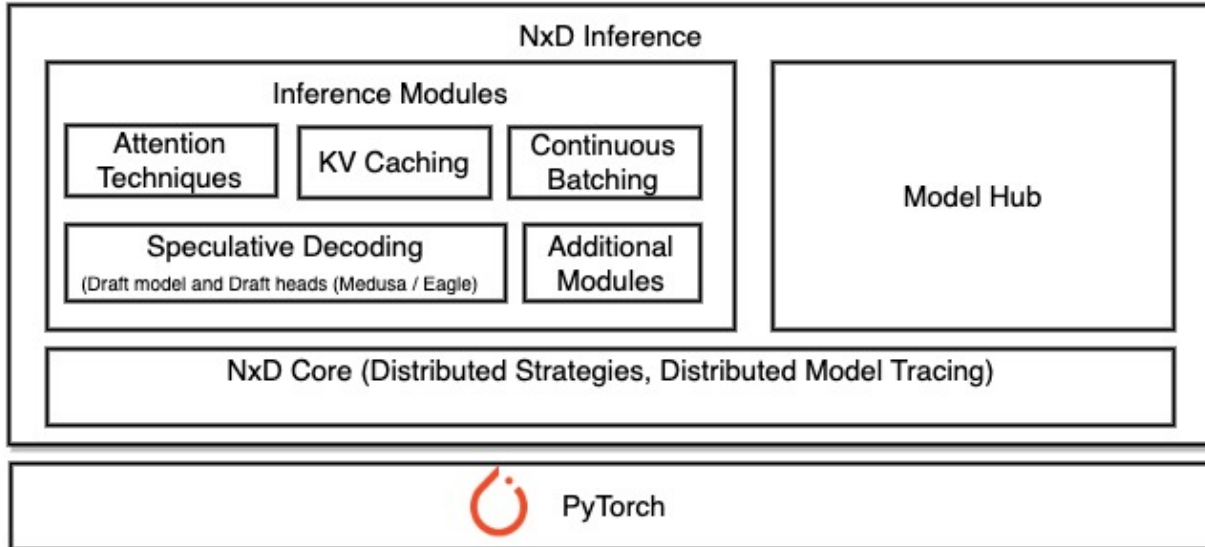


Fig. 3.4: NxID Inference High level Overview

### Using NxID Inference Library

ML developers can use NxID Inference library at different levels of abstraction. As shown in the below diagram Fig. 3.5, developers can use NxDI library in 3 different ways.

- **Deploy production ready models with vLLM:** NxDI supports production ready models like Llama-3.1, DBRX and Mixtral that can be easily deployed directly through vLLM. Customers can integrate their inference scripts directly with vLLM API.
- **Deploy production ready models with NxDI:** For customers who are not using vLLM, they can integrate with NxDI models directly for use cases such as static batching. For continuous batching, customers can also integrate with NxDI API to implement a custom model server with scheduler (other than vLLM) . See Fig. 3.5 b) for reference.
- **Integrate with Inference modules and NxID Core primitives:** As described in Fig. 3.5 c), customers who are looking to onboard new models which are not in NxDI model hub can integrate with inference modules and NxID Core primitives. In addition, customers who are looking to integrate with model servers other than vLLM can also integrate directly with NxID Inference modules and NxID core primitives.



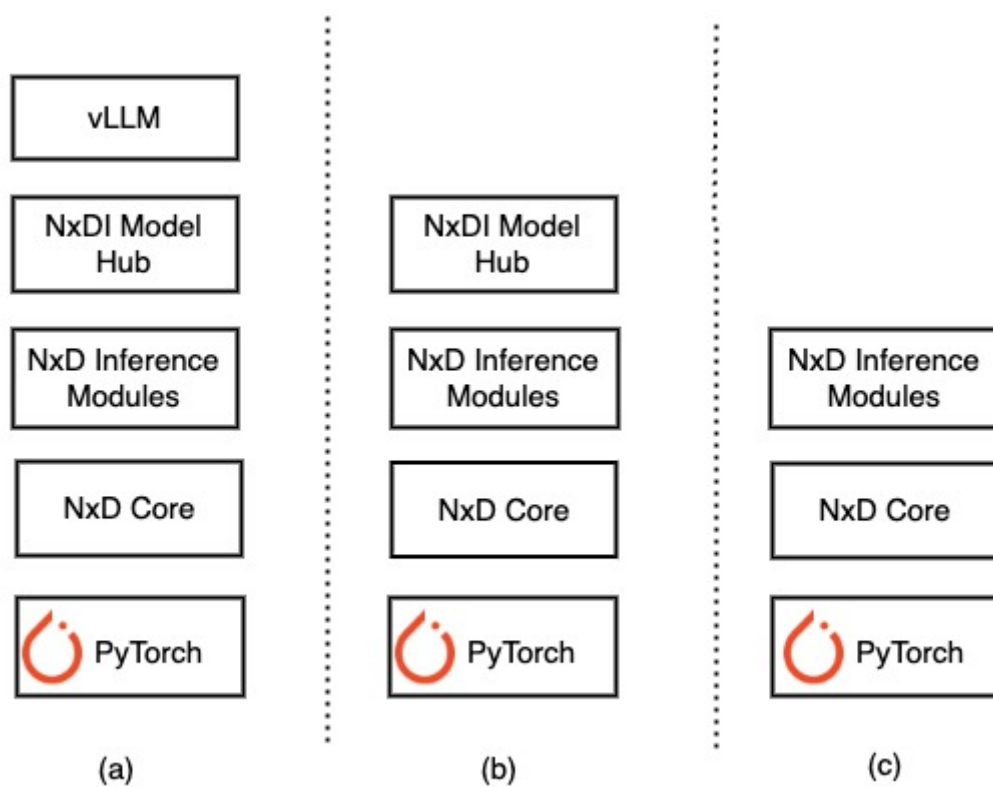


Fig. 3.5: Using Nx inference through various abstractions

### 3.2.2 Nx D Inference Setup Guide

The NeuronX Distributed (Nx D) Inference framework is built on top of *Nx D Core*. Follow the steps in this guide to set up your environment to run inference using the Nx D Inference framework.

#### Table of contents

- *Option 1: Launch an instance using a Neuron DLAMI*
- *Option 2: Use a Neuron Deep Learning Container (DLC)*
- *Option 3: Manually Install Nx D Inference*
  - *Setup a Neuron Environment*
  - *Install Nx D Inference*
- *Verify Nx D Inference Installation*

#### Option 1: Launch an instance using a Neuron DLAMI

Neuron Deep Learning AMIs (DLAMIs) are Amazon Machine Images (AMIs) that come with the Neuron SDK pre-installed. To quickly get started with Nx D Inference, you can launch an EC2 instance with the multi-framework DLAMI, which includes Nx D Inference and its dependencies. For more information, see the Neuron Multi-Framework DLAMI Guide and *Neuron DLAMI User Guide*.

After you launch an instance, you can run the following command to activate the Nx D Inference virtual environment.

```
source /opt/aws_neuronx_venv_pytorch_2_6_nxd_inference/bin/activate
```

#### Option 2: Use a Neuron Deep Learning Container (DLC)

Neuron Deep Learning Containers (DLCs) are Docker images that come with the Neuron SDK pre-installed. To run Nx D Inference in a Docker container, use the *Neuronx PyTorch Inference Containers*. For more information, see *Neuron Containers*.

#### Option 3: Manually Install Nx D Inference

Follow these instructions to manually install Nx D Inference on an instance.

---

**Note:** For information about which Python versions are compatible with the Neuron SDK, see Release Artifacts.

---

## Setup a Neuron Environment

Before you install NxD Inference, you must install the Neuron SDK and its dependencies, including PyTorch Neuron (`torch-neuronx`). Follow instructions for one of the following operating systems:

- PyTorch NeuronX Setup on Ubuntu 22
- PyTorch NeuronX Setup on Amazon Linux 2023

## Install NxD Inference

Run this command to install NxD Inference.

```
source aws_neuron_venv_pytorch/bin/activate
pip install -U pip
pip install --upgrade neuronx-cc==2.* neuronx-distributed-inference --extra-index-url https://pip.repos.neuron.amazonaws.com
```

## Verify NxD Inference Installation

To verify that NxD Inference installed successfully, check that you can run the `inference_demo` console script.

```
inference_demo --help
```

## 3.2.3 API Reference Guides

### NxD Inference API Reference

NeuronX Distributed (NxD) Inference (`neuronx-distributed-inference`) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. Neuronx Distributed Inference includes a model hub and modules that users can reference to implement their own models on Neuron.

This API guide describes API and configuration functions and parameters that you can use when you directly interact with the NxD Inference library.

**Note:** NxD Inference also supports integration with vLLM. When you use vLLM, you can use the `override_neuron_config` attribute to override defaults using the *NeuronConfig parameters* described in this API guide. For more information about vLLM integration, see *vLLM User Guide for NxD Inference*.

#### Table of contents

- *Configuration*
  - *NeuronConfig*
  - *InferenceConfig*
  - *MoENeuronConfig*
  - *FusedSpecNeuronConfig*

- *Generation*
  - *HuggingFaceGenerationAdapter*
- *Models*
  - *NeuronApplicationBase*
  - *NeuronBaseForCausalLM*
  - *NeuronBaseModel*
  - *ModelWrapper*

## Configuration

NxD Inference defines configuration objects that enable you to control how a model is compiled and used for inference. When you compile a model, its configuration is serialized to a JSON file in the compiled checkpoint, so you can distribute the compiled checkpoint to additional Neuron instances without needing to compile on each instance.

NxD Inference supports loading HuggingFace model checkpoints and configurations. When you run a model from a HuggingFace checkpoint, NxD Inference loads the model configuration from the model's `PretrainedConfig`.

## NeuronConfig

`NeuronConfig` contains compile-time configuration options for inference on Neuron.

## Initialization

Pass the `NeuronConfig` attributes as keyword args.

## Functions

- `NeuronConfig(**kwargs)` - Initializes a `NeuronConfig` with attributes from `kwargs`.

## Attributes

- General configuration
  - `batch_size` - The number of inputs to process in a single request. Defaults to 1.
  - `padding_side` - The padding side. Defaults to `right`.
  - `seq_len` - The sequence length, which is typically the sum of `max_context_length` and `max_new_tokens`. This value is the maximum sequence size that the model can process in a single request. Defaults to 128.
  - `max_context_length` - The maximum context length. Default to the `seq_len`.
  - `max_new_tokens` - The maximum number of tokens to generate in a single request. Default to the difference between `seq_len` and `max_context_length`. If the difference is zero, then `max_new_tokens` is set to `None`.
  - `max_length` - The maximum length to process. Default to the `seq_len`.

- `n_active_tokens` - The number of active tokens to track. Defaults to the `seq_len`.
  - `n_positions` - The number of positions to track. Defaults to the `seq_len`.
  - `torch_dtype` - The torch data type to use for computation. Choose from the following options. Defaults to `torch.bfloat16`.
    - \* `torch.bfloat16`
    - \* `torch.float16`
    - \* `torch.float32`
  - `rpl_reduce_dtype` - The torch data type to use for `all_reduce` operations in `RowParallelLinear` layers. Defaults to `None` and does the reduction in the input tensors dtype.
  - `cast_type` - The type of casting strategy to use when loading model parameters. Can be set to `config` (default) which casts all parameters to `torch_dtype`, or `as-declared` which casts all parameters to the dtype they were defined with.
  - `async_mode` - Whether to use asynchronous mode for inference. Defaults to `false`.
  - `save_sharded_checkpoint` - Whether to save the sharded weights in the compiled checkpoint. If this option is disabled, NxD Inference shards the weights during model load. Defaults to `true`.
  - `logical_nc_config` - The Logical NeuronCore Configuration (LNC). On Trn1 and Inf2, this defaults to 1. On Trn2, this defaults to 2. You can also configure LNC with the `NEURON_LOGICAL_NC_CONFIG` environment variable. For more information about LNC, see [Logical NeuronCore configuration](#).
    - \* Note: If you use Trn2 with NxD Inference v0.1 (Neuron 2.21), you must specify LNC=2 by setting `logical_neuron_cores=2` in `NeuronConfig`. The `logical_neuron_cores` attribute is deprecated in NxD Inference v0.2 and later.
  - `skip_sharding` - Whether to skip weight sharding during compilation. You can use this option if the compiled checkpoint path already includes sharded weights for the model. Defaults to `false`.
  - `weights_to_skip_layout_optimization` - The list of weight names to skip during weight layout optimization.
  - `skip_warmup` - Whether to skip warmup during model load. To improve the performance of the first request sent to a model, NxD Inference warms up the model during load. Defaults to `false`.
- Distributed configuration
    - `tp_degree` - The number of Neuron cores to parallelize across using tensor parallelism. Defaults to 1.
      - \* The number of attention heads needs to be divisible by the tensor-parallelism degree.
      - \* The total data size of model weights and key-value caches needs to be smaller than the tensor-parallelism degree multiplied by the amount of HBM memory per Neuron core.
        - On trn2, each Neuron core has 24GB of memory (with `logical_nc_config` set to 2).
        - On inf2/trn1, each Neuron core has 16GB of memory.
      - \* The Neuron runtime supports the following tensor-parallelism degrees:
        - trn2: 1, 2, 4, 8, 16, 32, and 64 (with `logical_nc_config` set to 2)
        - inf2: 1, 2, 4, 8, and 24
        - trn1: 1, 2, 8, 16, and 32
  - Attention
    - `flash_decoding_enabled` - Whether to enable flash decoding. Defaults to `false`.

- `fused_qkv` - Whether to fuse the query (Q), key (K), and value (V) weights in the models attention layers. This option improves performance by using larger matrices. Defaults to `false`.
- `sequence_parallel_enabled` - Whether to use sequence parallelism, which splits tensors along the sequence dimension. Defaults to `false`. Sequence parallel requires context sequence length to be divisible with tensor parallelism degree. Once enabled, sequence parallelism is only applied to context encoding.
- `qk_layernorm` - Whether to enable QK layer normalization. Defaults to `false`.
- `attention_dtype` - The torch data type to use for all operations in attention. Defaults to `None` and infers the dtype based on the dtype of the hidden\_states passed to attention.
- On-device sampling
  - `on_device_sampling_config` - The on-device sampling configuration to use. Specify this config to enable on-device sampling. This config is an `OnDeviceSamplingConfig`, which has the following attributes:
    - \* `do_sample` - Whether to use multinomial sampling (true) or greedy sampling (false). Defaults to `false`.
    - \* `top_k` - The top-k value to use for sampling. Defaults to 1.
    - \* `dynamic` - Whether to enable dynamic sampling. With dynamic sampling, you can pass different `top_k`, `top_p`, and `temperature` values to the `forward` call to configure sampling for each input in a batch. Defaults to `false`.
    - \* `deterministic` - Whether to enable deterministic sampling. Defaults to `false`.
    - \* `global_topk` - The global topK value to use. Defaults to 256.
- Bucketing
  - `enable_bucketing` - Whether to enable bucketing. Defaults to `false`. You can specify the buckets to use with the `context_encoding_buckets` and `token_generation_buckets` attributes. If you don't specify the buckets to use, NxDI automatically selects buckets based on the following logic.
    - \* Context encoding: Powers of two between 128 and the max context length.
      - Note: Max context length is equivalent to sequence length by default.
    - \* Token generation: Powers of two between 128 and the maximum sequence length.
  - `context_encoding_buckets` - The list of bucket sizes to use for the context encoding model.
  - `token_generation_buckets` - The list of bucket sizes to use for the token generation model.
- Quantization
  - `quantized` - Whether the model weights are quantized. Defaults to `false`.
  - `quantized_checkpoints_path` - The path to the quantized checkpoint. To quantize the model and save it to this path, use `NeuronApplicationBase`'s `save_quantized_state_dict` function. Specify one of the following:
    - \* A folder path. During quantization, NxD Inference saves the quantized model in safetensors format to this folder. To use a quantized model from a folder, it can be in safetensors or pickle format.
    - \* A file path to a quantized model file in pickle format.
  - `quantization_dtype` - The data type to use for quantization. Choose from the following options. Defaults to `int8`.
    - \* `int8` - 8 bit int.
    - \* `f8e4m3` - 8-bit float with greater precision and less range.

- Important: To use `f8e4m3` for quantization, you must set the `XLA_HANDLE_SPECIAL_SCALAR` environment variable to 1.
  - \* `f8e5m2` - 8-bit float with greater range and less precision.
- `quantization_type` - The type of quantization to use. Choose from the following options. Defaults to `per_tensor_symmetric`.
  - \* `per_tensor_symmetric`
  - \* `per_channel_symmetric`
- `modules_to_not_convert` - Specify a list of modules to be not quantized. Also, required when running inference on custom quantized models(using external libraries) where certain layers are left in full precision. Example: `["lm_head", "layers.0.self_attn", "layers.1.mlp", ...]`. Defaults to `None` (meaning all modules will be quantized)
- `draft_model_modules_to_not_convert` - Specify a list of modules in full precision when working with fused speculation. If no layers are required, add all layers in the list. Example: `["lm_head", "layers.0.self_attn", "layers.1.mlp", ...]`. This is only required in the case of fused speculation.
- KV cache quantization
  - `kv_cache_quant` - Whether to quantize the KV cache. When enabled, the model quantizes the KV cache to the `torch.float8_e4m3fn` data type. Defaults to `false`.
    - \* Important: To use `kv_cache_quant`, you must set the `XLA_HANDLE_SPECIAL_SCALAR` environment variable to 1.
- Kernels
  - `attn_kernel_enabled` - Whether to enable the flash attention kernel when supported. Defaults to `false`. Flash attention is automatically enabled by default for certain conditions, see `NeuronAttentionBase.get_flash_attention_strategy` in `neuronx_distributed_inference.modules.attention.attention_base`. Even explicitly enabled flash attention with `NeuronConfig(attn_kernel_enabled=True)` will be disabled for use cases where enabling it would be less efficient.
  - `qkv_kernel_enabled` - Whether to enable the fused QKV kernel. To use this option, you must set `fused_qkv` to `true` and `torch_dtype` to `torch.bfloat16`. Defaults to `false`.
  - `mlp_kernel_enabled` - Whether to enable the MLP kernel. To use this option, you must set `torch_dtype` to `torch.bfloat16`. Defaults to `false`.
  - `quantized_mlp_kernel_enabled` - Whether to enable the quantized MLP kernel, which uses FP8 compute to improve performance. To use this option, you must set `mlp_kernel_enabled` to `true`. Defaults to `false`.
  - `rmsnorm_quantize_kernel_enabled` - Whether to enable the quantized RMS norm kernel. Defaults to `false`.
- Continuous batching
  - `is_continuous_batching` - Whether to enable continuous batching. Defaults to `false`.
  - `max_batch_size` - The maximum batch size to use for continuous batching. Defaults to `batch_size`.
  - `ctx_batch_size` - The maximum batch size to use for the context encoding model in continuous batching. Defaults to `batch_size`.
  - `tkg_batch_size` - The maximum batch size to use for the token generation model in continuous batching. Defaults to `batch_size`.
- Speculative decoding

- `speculation_length` - The number of tokens to generate with the draft model before checking work with the primary model. Set this value to a positive integer to enable speculation. Defaults to 0.
- `spec_batch_size` - The batch size to use for speculation. Defaults to `batch_size`.
- `enable_eagle_speculation` - Whether to enable EAGLE speculation, where the previous hidden state is passed to a specialized target model to improve performance. Defaults to `false`.
- `enable_eagle_draft_input_norm` - Whether to perform input normalization in the EAGLE draft model. Defaults to `false`.
- `enable_fused_speculation` - Whether to enable fused speculation, where the target and draft model are fused into a single compiled model to improve performance. Fused speculation is enabled by default if `enable_eagle_speculation` is true. Otherwise, this defaults to `false`.
- Medusa decoding - Medusa is a speculation method that uses multiple smaller LM heads to perform speculation.
  - `is_medusa` - Whether to use Medusa decoding. Defaults to `false`
  - `medusa_speculation_length` - The number of tokens to generate with the Medusa heads before checking work with the primary model. Set this value to a positive integer. Defaults to 0.
  - `num_medusa_heads` - The number of LM heads to use for Medusa. Defaults to 0.
  - `medusa_tree` - The Medusa tree to use. For an example, see `medusa_mc_sim_7b_63.json` in the `examples` folder.
- Multi-LoRA serving
  - `lora_config` - The multi-lora serving configuration to use. Defaults to `none`. Specify this config to enable multi-LoRA serving. This config is `LoraServingConfig`, which has the following attributes:
    - \* `max_loras` - The maximum number of concurrent LoRA adapters in device memory. Defaults to 1.
    - \* `lora_ckpt_paths` - The checkpoint paths for LoRA adapters with key-value pairs. The key is the adapter ID and the value is the local path of the LoRA adapter checkpoint.
    - \* `lora_memory_transpose` - Transpose memory layout to optimize inference performance. Defaults to `True`.
    - \* `lora_shard_linear_layer` - Shard the linear layer across TP group to reduce memory consumption at the cost of communication overhead. Defaults to `False`.
- Compilation configuration
  - `cc_pipeline_tiling_factor` - The pipeline tiling factor to use for collectives. Defaults to 2.
- Debugging
  - `output_logits` - Whether to return model logits from the Neuron device when using on-device sampling. With on-device sampling, the model samples the logits on-device to return a singular token, and the model output includes only the tokens (without the logits) to improve performance. The `output_logits` feature enables you to output the logits alongside the token, which enables you to run logit validation and investigate the model output. Note: This feature impacts performance and shouldn't be used in production; this should only be used for testing and debugging model logits.



## InferenceConfig

InferenceConfig contains a NeuronConfig and model configuration attributes.

## Initialization

You can pass attributes through keyword args, or provide a `load_config` hook that is called during initialization to load the configuration attributes.

InferenceConfig is compatible with HuggingFace transformers. To use a model from HuggingFace transformers, you can populate an InferenceConfig with the attributes from the model's PretrainedConfig, which is stored in `config.json` in the model checkpoint.

```
from neuronx_distributed_inference.models.llama import (
    LlamaInferenceConfig,
    LlamaNeuronConfig
)
from neuronx_distributed_inference.utils.hf_adapter import load_pretrained_config

model_path = "/home/ubuntu/models/Meta-Llama-3.1-8B"

neuron_config = LlamaNeuronConfig()
config = LlamaInferenceConfig(
    neuron_config,
    load_config=load_pretrained_config(model_path),
)
```

## Attributes

An InferenceConfig includes `neuron_config` and any other attributes that you set during initialization.

- `neuron_config` - The NeuronConfig for this inference config.
- `fused_spec_config` - The FusedSpecNeuronConfig for this inference config. Provide a fused spec config if using fused speculation.
- `load_config` - The `load_config` hook to run during initialization. You can provide a load config hook to load configuration attributes from another source. To load from a HuggingFace PretrainedConfig, pass the load config hook returned by `load_pretrained_config`. The `load_pretrained_config` hook provider takes the model path as its argument.

InferenceConfig also supports an attribute map, which lets you configure additional names or aliases for attributes. When you get or set an attribute by an alias, you retrieve or modify the value of the original attribute. When you initialize an InferenceConfig from a HuggingFace PretrainedConfig, it automatically inherits the attribute map from that PretrainedConfig.

## Functions

- `InferenceConfig(neuron_config, load_config=None, **kwargs)` - Initializes an `InferenceConfig`.
- `load_config(self)` - Loads the config attributes. This function does nothing by default; subclasses can override it to provide a model-specific implementation. This function is called during initialization unless a `load_config` hook is provided.
- `get_required_attributes(self)` - Returns the list of attribute names that must be present in this config for it to validate during initialization. This function returns an empty list by default; subclasses can override it to require model-specific attributes to be present.
- `validate_config(self)` - Checks that the config is valid. This function is called during initialization. By default, this function checks that the attributes returned by `get_required_attributes` are present. Subclasses can override this function to implement model-specific validation.
- `save(self, model_path)` - Serializes the config to a JSON file, `neuron_config.json` in the given model path.
- `to_json_file(self, json_file)` - Serializes the config to the given JSON file.
- `to_json_string(self)` - Serializes the config to a string in JSON format.
- `load(cls, model_path, **kwargs)` - Loads the config from the `neuron_config.json` file in the given model path. You can specify `kwargs` to override attributes in the config.
- `from_json_file(cls, json_file, **kwargs)` - Loads the config from the given JSON file. You can specify `kwargs` to override attributes in the config.
- `from_json_string(cls, json_string, **kwargs)` - Loads the config from the given JSON string. You can specify `kwargs` to override attributes in the config.
- `get_neuron_config_cls(cls)` - Returns the `NeuronConfig` class type to use for this `InferenceConfig`. This function returns `NeuronConfig` by default; subclasses can override this function to configure a specific `NeuronConfig` subclass to use.

## MoENeuronConfig

A `NeuronConfig` subclass for mixture-of-experts (MoE) models. This config includes attributes specific to MoE models. MoE model configurations, such as `DbrxNeuronConfig`, are subclasses of `MoENeuronConfig`.

## Initialization

Pass the attributes as keyword args.

## Functions

- `MoENeuronConfig(**kwargs)` - Initializes an `MoENeuronConfig` with attributes from `kwargs`.

## Attributes

- `capacity_factor` - The capacity factor to use when allocating tokens across experts. When an expert is at capacity, tokens allocated to that expert are dropped until that expert has capacity again. Defaults to `None`, which means that NxDI waits until an expert has capacity, and no tokens are dropped.
- `glu_mlp` - Whether to use a Gated Linear Unit in the MLP. Defaults to `false`.

## FusedSpecNeuronConfig

A configuration for a model that uses fused speculation, which is a speculative decoding feature where the target and draft models are compiled into a combined model to improve performance. For more information, see [Fused Speculation](#).

## Attributes

- `worker_cls` - The model class to use for fused speculation. This class should be a subclass of `NeuronBaseModel`.
- `draft_config` - The `InferenceConfig` for the draft model.
- `draft_model_path` - The path to the draft model checkpoint.

## Generation

### HuggingFaceGenerationAdapter

NxD Inference supports running inference with the HuggingFace `generate` inference. To use HuggingFace-style generation, create a `HuggingFaceGenerationAdapter` that wraps a Neuron application model. Then, you can call `generate` on the adapted model.

```
generation_model = HuggingFaceGenerationAdapter(neuron_model)
outputs = generation_model.generate(
    inputs.input_ids,
    attention_mask=inputs.attention_mask,
    generation_config=generation_config
)
```

## Models

NxD Inference provides a [model hub](#) with production ready models. You can use these existing models to run inference, or use them as reference implementations when you develop your own models on Neuron. All model inherit from base classes that provide a basic set of functionality that is common to all models.

## NeuronApplicationBase

NeuronApplicationBase is the base class for all application models, including NeuronBaseForCausalLM. NeuronApplicationBase provides functions to compile and load models. This class extends `torch.nn.Module`. Application models are the entry point to running inference with NxD Inference. You can extend this class to define new application models that implement use cases in addition to causal LM.

### Attributes

- `config` - The InferenceConfig for this model.
- `neuron_config` - The NeuronConfig for this model.
- `model_path` - The model path for this model.
- `models` - The list of models that make up this application model. These models are instances of ModelWrapper. Add models to this list to compile them with `compile`.
- `is_compiled` - Whether this model is compiled.
- `is_loaded_to_neuron` - Whether this model is loaded to the Neuron device.

### Functions

- `NeuronApplicationBase(self, model_path, config=None, neuron_config=None)` - Initializes an application model from the given model path, and optionally the given InferenceConfig (`config`) and NeuronConfig (`neuron_config`). If no InferenceConfig is provided, this function loads the config from the given model path.
- `compile(self, compiled_model_path, debug=False)` - Compiles this model for Neuron and saves the compiled model to the given path. This function compiles all models added to `self.models`. This function also shards the weights for the model. To produce HLO files that have source annotations enabled for debugging, set `debug` to `True`. When `debug` is enabled, HLOs contain following attributes for each computation: `op_type`, `op_name`, `source_file`, and `source_line`.
- `load(self, compiled_model_path)` - Loads the compiled model from the given path to the Neuron device. This function also loads the model weights to the Neuron device.
- `load_weights(self, compiled_model_path)` - Loads the model weights from the given path to the Neuron device. You can call this function to load new weights without reloading the entire model.
- `shard_weights(self, compiled_model_path)` - Shards the model's weights and serializes the sharded weights to the given path.
- `forward(self, **kwargs)` - The forward function for this application model. This function must be implemented by subclasses.
- `validate_config(cls, config)` - Checks whether the config is valid for this model. By default, this function requires that `neuron_config` is present. This function can be implemented by subclasses to provide model-specific validation.
- `get_compiler_args(self)` - Returns the Neuron compiler arguments to use when compiling this model. By default, this returns no compiler arguments. This function can be implemented by subclasses to use model-specific compiler args.
- `to_cpu(self)` - Allows inference to be run entirely on CPU. Use this in place of the `compile` and `load` functions. Note that CPU inference doesn't currently work for configurations that use kernels.

- `get_state_dict(cls, model_path, config)` - Gets the state dict for this model. By default, this function loads the state dict from the given model path. This function calls the class' `convert_hf_to_neuron_state_dict` function to convert the state dict according to the specific model. Subclasses can override this function to provide custom state dict loading.
  - When loading the state dict, this function replaces keys that start with the class' `_STATE_DICT_MODEL_PREFIX` value with the class' `_NEW_STATE_DICT_MODEL_PREFIX` value. Subclasses can set these values to update the state dict keys accordingly.
- `convert_hf_to_neuron_state_dict` - Converts a state dict from HF format to the format expected by Neuron. By default, this function returns the state dict without modifying it; subclasses can override this to provide custom conversion for each model.
- `save_quantized_state_dict(cls, model_path, config)` - Quantizes the model's state dict and saves the quantized checkpoint to the `quantized_checkpoint_path` from the given config's `NeuronConfig`.
- `generate_quantized_state_dict(cls, model_path, config)` - Generates the quantized state dict for this model. This function loads the HuggingFace model from the given model path in order to quantize the model. Then, this function passes the quantized model to `prepare_quantized_state_dict` to generate the state dict. Subclasses can override this function to customize quantization.
- `prepare_quantized_state_dict(cls, hf_model_quant)` - Prepares the quantized state dict for the model. By default, this function converts the state dict from qint8 to int8. Subclasses can override this function to customize quantization.
- `load_hf_model(model_path)` - Loads the equivalent HuggingFace model from the given model path. Subclasses must implement this function to use quantization or to generate expected outputs when evaluating accuracy with `accuracy.py`.
- `reset(self)` - Resets the model state. By default, this function does nothing; subclasses can implement it to provide custom behavior.

## NeuronBaseForCausalLM

`NeuronBaseForCausalLM` is the base application class that you use to generate text with causal language models. This class extends `NeuronApplicationBase`. You can extend this class to run text generation in custom models.

### Attributes

- `kv_cache_populated` - Whether the KV cache is populated.

### Functions

- `NeuronBaseForCausalLM(self, *args, **kwargs)` - Initializes the `NeuronApplicationBase` and configures the models used in this LM application, including context encoding, token gen, and others, based on the given `NeuronConfig`.
- `forward(self, input_ids=None, seq_ids=None, attention_mask=None, position_ids=None, sampling_params=None, prev_hidden=None, past_key_values=None, inputs_embeds=None, labels=None, use_cache=None, output_attentions=None, output_hidden_states=None, medusa_args=None, return_dict=None, input_capture_hook=None)` - The forward function for causal LM. This function routes the forward pass to the correct sub-model (such as context encoding or token generation) based on the current model state. If an `input_capture_hook` function is provided, the forward function calls the hook with the model inputs as arguments.

- `reset(self)` - Resets the model for a new batch of inference. After the model is reset, a subsequent run will invoke the context encoding model.
- `reset_kv_cache(self)` - Resets the KV cache by replacing its key values with zeroes.

## NeuronBaseModel

NeuronBaseModel is the base class for all models. This class extends `torch.nn.Module`. In instances of NeuronBaseModel, you define the modules, such as attention, MLP, and decoder layers, that make up a model. You can extend this class to define custom decoder models.

## Attributes

- `sampler` - The sampler to use for on-device sampling.
- `kv_mgr` - The KV cache manager to use to manage the KV cache.
- `sequence_dimension` - The dimension for sequence parallelism.

## Functions

- `NeuronBaseModel(config, optimize_inference=True)` - Initializes the Neuron model from the given config. If `optimize_inference` is true, then this initializes a KV cache manager and sampler (if on-device sampling).
- `setup_attr_for_model(self, config)` - Initializes the following attributes for the model. These attributes are used by modules within the model. Subclasses must implement this function to set these attributes from the config.
  - `on_device_sampling`
  - `tp_degree`
  - `hidden_size`
  - `num_attention_heads`
  - `num_key_value_heads`
  - `max_batch_size`
  - `buckets`
- `init_model(self, config)` - Initializes the following modules for the model. Subclasses must implement this function.
  - `embed_tokens`
  - `layers`
  - `norm`
  - `lm_head`
- `forward(self, input_ids, attention_mask, position_ids, seq_ids, accepted_indices=None, current_length=None, medusa_mask=None, scatter_index=None)`
  - The forward function for this model.

## ModelWrapper

Wraps a model to prepare it for compilation. Neuron applications, such as `NeuronBaseForCausalLM`, use this class to prepare a model for compilation. `ModelWrapper` defines the inputs to use when tracing the model during compilation.

To define a custom model with additional model inputs, you can extend `ModelWrapper` and override the `input_generator` function, which defines the inputs for tracing.

## Functions

- `ModelWrapper(config, model_cls, tag, compiler_args)` - Initializes a model wrapper from a given config and model class. This model class is used to compile the model with the given compiler args. The tag is used to identify the compiled model in the application.
- `input_generator(self)` - Returns a list of input tensors to use to trace the model for compilation. When you trace and compile a model, the trace captures only the code paths that are run with these inputs. To support different inputs and code paths based on configuration options, provide configuration-specific inputs in `input_generator`.

Use the NxD Inference (`neuronx-distributed-inference`) API Reference Guides to learn how to use NxD Inference.

- [NxD Inference API Reference](#)

## 3.2.4 Developer Guides

### NxD Inference Features Configuration Guide

NxD Inference (`neuronx-distributed-inference`) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. Neuronx Distributed Inference includes a model hub and modules that users can reference to implement their own models on Neuron.

#### Table of contents

- [Checkpoint compatibility with HuggingFace Transformers](#)
- [Checkpoint support](#)
- [Compiling models](#)
- [Neuron Persistent Cache](#)
- [Serialization support](#)
- [Logical NeuronCore Configuration \(LNC\) support](#)
- [Tensor-parallelism support](#)
  - [Examples](#)
- [Sequence Parallelism](#)
- [Compile-time Configurations](#)
- [Hugging Face generate\(\) API support](#)
  - [Example](#)

- *On-device Sampling Support*
  - *Dynamic Sampling*
  - *Greedy Sampling*
  - *Multinomial (Top-K) Sampling*
  - *Top-P Support in On-Device Sampling*
  - *Temperature Support in On-Device Sampling*
- *QKV Weight Fusion*
- *Bucketing*
  - *Automatic Bucketing*
  - *Configuring Specific Buckets*
- *Quantization*
  - *Model Weight Quantization*
  - *KV Cache Quantization*
- *Speculative Decoding*
  - *Speculative Decoding with a Draft model*
  - *Medusa Speculative Decoding*
  - *EAGLE Speculative Decoding*
- *MoE model architecture support*
- *Grouped-query attention (GQA) support*
- *Asynchronous Runtime Support*
- *Prefix Caching Support*
  - *Bucketing with Prefix Caching*
- *Multi-LoRA Serving*
  - *Enable multi-LoRA serving*
  - *Maximum number of LoRA adapters supported*
- *Disaggregated Inference [BETA]*

## Checkpoint compatibility with HuggingFace Transformers

Models included in the NxD Inference model hub are checkpoint-compatible with HuggingFace Transformers. Supporting other checkpoint formats in NxD Inference is possible through converting the obtained checkpoint to the standard HuggingFace Transformers checkpoint format.



## Checkpoint support

NxD Inference supports older PyTorch binary checkpoints and newer [safetensors](#) checkpoints. For improved load speed and reduced host memory consumption, we recommend to always use safetensors by default. Both regular and sharded variants of checkpoints are supported.

NxD Inference supports weights stored in the model path in the following formats:

| Format      | Sharded | File name                    |
|-------------|---------|------------------------------|
| Safetensors | No      | model.safetensors            |
| Safetensors | Yes     | model.safetensors.index.json |
| Pickle      | No      | pytorch_model.bin            |
| Pickle      | Yes     | pytorch_model.bin.index.json |

If your weights are in another format, you must convert them to one of these formats before you can compile and load the model to Neuron. See the following references for more information about these formats:

- Safetensors:
  - <https://github.com/huggingface/safetensors>
  - <https://huggingface.co/docs/safetensors/en/convert-weights>
- Pickle:
  - <https://docs.python.org/3/library/pickle.html>

## Compiling models

To run a model on Neuron with NxD Inference, you compile Python code into a NEFF file (Neuron Executable File Format), which you can load to Neuron devices using the Neuron Runtime.

When you call `compile()`, NxD Inference does the following:

1. Trace the Python code to produce an HLO file.
2. Use the Neuron Compiler to compile the HLO file into a NEFF.

During the trace process, the model code is traced based on a given sample tensor for each input. As a result, model code should avoid dynamic logic that depends on the input values in a tensor, because NxD Inference compiles only the code path that is traced for the sample input tensor.

```
# Configure, initialize, and compile a model.
model = NeuronLlamaForCausalLM(model_path, config)
model.compile(compiled_model_path)
```

## Neuron Persistent Cache

The Neuron Persistent Cache is enabled by default for NxD Inference library. Model artifacts which have been compiled once will be cached and reused on successive runs when possible. Model artifacts will only be reused when compiling with the same compiler version (`neuronx-cc`), model configurations, and compiler flags. Neuron Persistent Cache also includes other features, such as using an S3 bucket as the cache backend. For more detailed information, see the [Persistent cache documentation](#)

## Serialization support

When you compile a model with NxD Inference, the library serializes the model to a given folder. After you have a serialized model, you can load it directly to a Neuron device without needing to compile again.

The compile function does not serialize sharded weights by default, and you can enable this functionality with the `save_sharded_checkpoint` flag in `NeuronConfig`. For more information on weights sharding, see [NxD Inference Weights Sharding Guide](#).

## Logical NeuronCore Configuration (LNC) support

On Trn2 instances, Neuron supports Logical NeuronCore (LNC) configuration, which combines multiple physical NeuronCores into a single logical NeuronCore. On Trn2 instances, the Neuron SDK is optimized for LNC=2, which means each NeuronCore visible to the Neuron SDK is two physical NeuronCores.

NxD Inference automatically chooses the correct LNC configuration based on the target platform. To override the default LNC configuration, you can set the `NEURON_LOGICAL_NC_CONFIG` environment variable, or set the `logical_nc_config` flag in `NeuronConfig`.

```
neuron_config = NeuronConfig(logical_nc_config=2)
```

For more information about logical NeuronCore support, see [Logical NeuronCore configuration](#).

## Tensor-parallelism support

For transformer decoders used in large language models, tensor-parallelism is necessary as it provides a way to shard the models' large weight matrices onto multiple NeuronCores, and having NeuronCores working on the same matrix multiply operation collaboratively. `neuronx-distributed-inference`'s tensor-parallelism support makes heavy use of collective operations such as all-reduce, which is supported natively by the Neuron runtime.

There are some principles for setting tensor-parallelism degree (number of NeuronCores participating in sharded matrix multiply operations) for Neuron-optimized transformer decoder models.

1. The number of attention heads needs to be divisible by the tensor-parallelism degree.
2. The total data size of model weights and key-value caches needs to be smaller than the tensor-parallelism degree multiplied by the amount of memory per Neuron core.
  1. On Trn2, each Neuron core has 24GB of memory (with LNC2).
  2. On Inf2/Trn1, each Neuron core has 16GB of memory.
3. The Neuron runtime supports the following tensor-parallelism degrees:
  1. Trn2: 1, 2, 4, 8, 16, 32, and 64 (with LNC2)
  2. Inf2: 1, 2, 4, 8, and 24
  3. Trn1: 1, 2, 8, 16, and 32

## Examples

1. `meta-llama/Meta-Llama-3.1-8B` has 32 attention heads, and when running at batch size 1 and `bfloat16` precision, the model requires about 16GB memory. Therefore, a `trn1.2xlarge` with 32GB device memory is sufficient.
2. `meta-llama/Meta-Llama-3.1-70B` has 64 attention heads, and when running at batch size 1 and `bfloat16` precision, the model requires about 148GB memory. Therefore, it can run on 16 NeuronCores on one `trn1.32xlarge` using 256GB device memory.

## Sequence Parallelism

Sequence parallelism splits tensors across the sequence dimension to improve performance. You can enable sequence parallelism by setting `sequence_parallel_enabled=True` in `NeuronConfig`.

```
neuron_config = NeuronConfig(sequence_parallel_enabled=True)
```

## Compile-time Configurations

NxD Inference models support a variety of compile-time configurations you can use to tune model performance. For more information, see the [NxD Inference API Reference](#).

## Hugging Face `generate()` API support

NxD Inference models support the HuggingFace `generate()` API via the `HuggingFaceGenerationAdapter` class. This adapter wraps a Neuron model to provide the HuggingFace generation interface.

NxD Inference's supports the following HuggingFace generation modes:

- Greedy decoding — `num_beams=1` and `do_sample=False`.
- Multinomial sampling — `num_beams=1` and `do_sample=True`.
- Assisted (speculative) decoding — `assistant_model` or `prompt_lookup_num_tokens` are specified.

NxD Inference doesn't currently support other HuggingFace generation modes such beam-search sampling.

Note: When you call `generate`, the number of prompts must match the `batch_size` for the model, which is an attribute of `NeuronConfig`.

```
neuron_config = NeuronConfig(batch_size=2)
```

## Example

The following example demonstrates how to wrap a model with `HuggingFaceGenerationAdapter` to call `generate()`.

```
from neuronx_distributed_inference.utils.hf_adapter import HuggingFaceGenerationAdapter

# Init Neuron model, HuggingFace tokenizer, HuggingFace and generation config.

# Run generation with HuggingFaceGenerationAdapter.
```

(continues on next page)

(continued from previous page)

```

generation_model = HuggingFaceGenerationAdapter(model)
inputs = tokenizer(prompts, padding=True, return_tensors="pt")
outputs = generation_model.generate(
    inputs.input_ids,
    generation_config=generation_config,
    attention_mask=inputs.attention_mask,
    max_length=model.neuron_config.max_length,
    **kwargs,
)

output_tokens = tokenizer.batch_decode(
    outputs, skip_special_tokens=True, clean_up_tokenization_spaces=False
)

print("Generated outputs:")
for i, output_token in enumerate(output_tokens):
    print(f"Output {i}: {output_token}")

```

## On-device Sampling Support

On-device sampling performs sampling logic on the Neuron device (rather than on the CPU) to achieve better performance. To enable on device sampling, provide an `OnDeviceSamplingConfig` for the `on_device_sampling_config` attribute in `NeuronConfig`.

```

on_device_sampling_config = OnDeviceSamplingConfig(global_topk=256)
neuron_config = NeuronConfig(on_device_sampling_config=on_device_sampling_config)

```

## Dynamic Sampling

With dynamic sampling, you can pass different `top_k`, `top_p`, and `temperature` values to the forward call to configure sampling for each input in a batch. To enable dynamic sampling, provide an `OnDeviceSamplingConfig` with `dynamic=True`.

```

on_device_sampling_config = OnDeviceSamplingConfig(dynamic=True)
neuron_config = NeuronConfig(on_device_sampling_config=on_device_sampling_config)

```

To use dynamic sampling, pass a `sampling_params` tensor to the forward function of the model. The `sampling_params` tensor has shape `[batch_size, 3]`, where the three values per batch are `top_k`, `top_p`, and `temperature`.

The following example demonstrates how to create `sampling_params` for a batch with two inputs. In the first input, `top_k=50`, `top_p=0.5`, and `temperature=0.75`. In the second input, `top_k=5`, `top_p=1.0`, and `temperature=1.0`.

```

sampling_params = torch.tensor([[50, 0.5, 0.75], [5, 1.0, 1.0]])

```

## Greedy Sampling

By default, on-device sampling uses greedy sampling, where the model picks the highest scoring token.

## Multinomial (Top-K) Sampling

With multinomial (top-k) sampling, the model picks one of the top  $k$ -highest scoring tokens. To use on-device multinomial sampling, you must enable dynamic sampling. You can configure the default `top_k` attribute in the `OnDeviceSamplingConfig`, or you can specify the `top_k` value in each call to the model's `forward` function.

```
on_device_sampling_config = OnDeviceSamplingConfig(top_k=5)
```

## Top-P Support in On-Device Sampling

To use top-p in on-device sampling, enable dynamic sampling, and specify `top_p` values in the `sampling_params`.

## Temperature Support in On-Device Sampling

To adjust temperature in on-device sampling, enable dynamic sampling, and specify `temperature` values in the `sampling_params`.

## QKV Weight Fusion

QKV weight fusion concatenates a model's query, key and value weight matrices to achieve better performance, because larger matrices allow for more efficient data movement and compute. You can enable QKV weight fusion by setting `fused_qkv=True` in the `NeuronConfig`.

```
neuron_config = NeuronConfig(fused_qkv=True)
```

## Bucketing

LLM inference is a generation process that can produce variable length sequences. This poses a problem since the Neuron compiler produces executables which expect statically shaped inputs and outputs. To make LLMs work with different shapes, NxD Inference supports buckets and applies padding wherever it is required. When you run inference, NxD Inference automatically chooses the smallest bucket that fits the input for optimal performance. For more information about bucketing, see [torch-neuronx-autobucketing-devguide](#).

## Automatic Bucketing

When automatic bucketing is enabled, NxD Inference automatically chooses buckets for each model according to the following logic:

- Context encoding: Powers of two between 128 and the max context length.
  - Note: Max context length is equivalent to sequence length by default.
- Token generation: Powers of two between 128 and the maximum sequence length.

To enable automatic bucketing, set `enable_bucketing=True` in `NeuronConfig`.

```
neuron_config = NeuronConfig(enable_bucketing=True)
```

### Configuring Specific Buckets

You can configure specific buckets to further optimize inference based on the input and output length distribution that you expect to process with your model. In `NeuronConfig`, set `enable_bucketing=True`, and provide a list of bucket sizes in `context_encoding_buckets` and/or `token_generation_buckets`.

```
neuron_config = NeuronConfig(
    enable_bucketing=True,
    context_encoding_buckets=[1024, 2048, 4096],
    token_generation_buckets=[8192]
)
```

### Quantization

NxD Inference supports quantization, where model weights and data are converted to a smaller data type to reduce memory bandwidth usage, which improves model performance.

Note: Quantization slightly reduces accuracy due to using data types with lower precision and/or lower range.

### Model Weight Quantization

NxD Inference supports quantizing model weights to the following data types:

- INT8 (int8) - 8 bit int.
- FP8 - 8 bit float.
  - f8e4m3 - 8-bit float with greater precision and less range.
    - \* Important: To use f8e4m3 for quantization, you must set the `XLA_HANDLE_SPECIAL_SCALAR` environment variable to 1.
  - f8e5m2 - 8-bit float with greater range and less precision.

NxD Inference supports the following quantization methods, which you specify with `quantization_type` in `NeuronConfig`:

- *per\_tensor\_symmetric*
- *per\_channel\_symmetric*

### Example

The following example demonstrates how to quantize a model to INT8. To quantize a model to a different data type, change the `quantization_dtype` config attribute in `NeuronConfig`.

```
from neuronx_distributed_inference.models.config import NeuronConfig
from neuronx_distributed_inference.models.llama.modeling_llama import (
    LlamaInferenceConfig,
    NeuronLlamaForCausalLM
```

(continues on next page)

(continued from previous page)

```

)
from neuronx_distributed_inference.utils.hf_adapter import load_pretrained_config

model_path = "/home/ubuntu/models/Llama-3.1-8B"
quantized_model_path = "/home/ubuntu/models/Llama-3.1-8B-quantized"

neuron_config = NeuronConfig(
    quantized=True,
    quantized_checkpoints_path=quantized_model_path,
    quantization_dtype="int8",
    quantization_type="per_tensor_symmetric"
)

config = LlamaInferenceConfig(
    neuron_config,
    load_config=load_pretrained_config(model_path)
)

# Quantize the model and save it to `quantized_checkpoints_path`.
NeuronLlamaForCausalLM.save_quantized_state_dict(model_path, config)

# Compile, load, and use the model.
model = NeuronLlamaForCausalLM(model_path, config)

```

## KV Cache Quantization

NxD Inference supports KV cache quantization, where the model's KV cache is quantized to a smaller data type. When enabled, the model quantizes the KV cache to the `torch.float8_e4m3fn` data type. Before using the KV cache, the model dequantizes the KV cache to the data type specified by `torch_dtype` in `NeuronConfig`.

To enable KV cache quantization, set `kv_cache_quant=True` in `NeuronConfig`.

```
neuron_config = NeuronConfig(kv_cache_quant=True)
```

- Important: To use KV cache quantization, you must set the `XLA_HANDLE_SPECIAL_SCALAR` environment variable to 1.

## Speculative Decoding

Speculative decoding is a performance optimization technique where a smaller *draft* LLM model predicts the next tokens, and the larger *target* LLM model verifies those predictions. NxD Inference supports the following speculative decoding implementations:

1. *Speculative decoding with a draft model*, where a separate draft model predicts the next  $n$  tokens for the target model. Each model is compiled independently.
2. *Medusa speculative decoding*, where several small model heads predict next tokens, and the target model verifies all predictions at the same time.
3. *EAGLE speculative decoding*, where the draft model uses additional context from the target model to improve generation efficiency. NxD Inference supports EAGLE v1 with a flat draft structure.

## Speculative Decoding with a Draft model

To use speculative decoding with a draft model, you configure, compile, and load a draft model in addition to the main target model. To enable speculative decoding with a draft model, set `speculation_length` and `trace_tokengen_model=False` in the target model's `NeuronConfig`. The draft model's `NeuronConfig` should use the same configuration but with these additional attributes reset to their defaults.

Speculative decoding with a draft model currently supports only batch sizes of 1.

### Example

The following example demonstrates using Llama-3.2 3B as a draft model for Llama-3.1 70B. The speculation length is set to 5 tokens.

```
import copy

from transformers import AutoTokenizer, GenerationConfig

from neuronx_distributed_inference.models.config import NeuronConfig
from neuronx_distributed_inference.models.llama.modeling_llama import (
    LlamaInferenceConfig,
    NeuronLlamaForCausalLM
)
from neuronx_distributed_inference.utils.accuracy import get_generate_outputs
from neuronx_distributed_inference.utils.hf_adapter import load_pretrained_config

prompts = ["I believe the meaning of life is"]

model_path = "/home/ubuntu/models/Llama-3.1-70B"
draft_model_path = "/home/ubuntu/models/Llama-3.2-3B"
compiled_model_path = "/home/ubuntu/neuron_models/Llama-3.1-70B"
compiled_draft_model_path = "/home/ubuntu/neuron_models/Llama-3.2-3B"

# Initialize target model.
neuron_config = NeuronConfig(
    speculation_length=5,
    trace_tokengen_model=False
)
config = LlamaInferenceConfig(
    neuron_config,
    load_config=load_pretrained_config(model_path)
)
model = NeuronLlamaForCausalLM(model_path, config)

# Initialize draft model.
draft_neuron_config = copy.deepcopy(neuron_config)
draft_neuron_config.speculation_length **= 0
draft_neuron_config.trace_tokengen_model **= True
draft_config = LlamaInferenceConfig(
    draft_neuron_config,
    load_config=load_pretrained_config(draft_model_path)
)
draft_model = NeuronLlamaForCausalLM(draft_model_path, draft_config)
```

(continues on next page)



(continued from previous page)

```

# Compile and save models.
model.compile(compiled_model_path)
draft_model.compile(compiled_draft_model_path)

# Load models to the Neuron device.
model.load(compiled_model_path)
draft_model.load(compiled_draft_model_path)

# Load tokenizer and generation config.
tokenizer **=** AutoTokenizer.from_pretrained(model_path, padding_side**=**neuron_config.
↳padding_side)
generation_config = GenerationConfig.from_pretrained(model_path)

# Run generation.
_, output_tokens = get_generate_outputs(
    model,
    prompts,
    tokenizer,
    is_hf=False,
    draft_model=draft_model,
    generation_config=generation_config
)

print("Generated outputs:")
for i, output_token in enumerate(output_tokens):
    print(f"Output {i}: {output_token}")

```

## Medusa Speculative Decoding

To use Medusa speculative decoding, you must use a model that is specifically fine-tuned for Medusa speculation, such as `text-generation-inference/Mistral-7B-Instruct-v0.2-medusa`. You must also provide a Medusa tree. For an example Medusa tree, see `medusa_mc_sim_7b_63.json` in the `examples` folder in NeuronX Distributed Inference.

To enable Medusa, set `is_medusa=True`, set the `medusa_speculation_length`, set the `num_medusa_heads`, and specify the `medusa_tree`.

```

def load_json_file(json_path):
    with open(json_path, "r") as f:
        return json.load(f)

medusa_tree = load_json_file("medusa_mc_sim_7b_63.json")

neuron_config = NeuronConfig(
    is_medusa=True,
    medusa_speculation_length=64,
    num_medusa_heads=4,
    medusa_tree=medusa_tree
)

```

To run generation with a Medusa model and the HuggingFace `generate()` API, set the `assistant_model` to the target model.

For more information about Medusa speculative decoding, see the official implementation on GitHub: <https://github.com/FasterDecoding/Medusa>.

Medusa speculative decoding currently supports only batch sizes of 1.

## EAGLE Speculative Decoding

NxD Inference supports EAGLE v1 speculative decoding with a flat draft structure.

## EAGLE Checkpoint Compatibility

To use EAGLE speculative decoding, you must use a draft model that is specifically fine-tuned for EAGLE speculation. Additionally, to use EAGLE with NxD Inference, the draft model must include the LM head weights from the target model. These weights are shared between the draft and target model.

Because NxD Inference uses a flat draft structure, it predicts only one token per draft iteration. Although NxD Inference doesn't support EAGLE with a tree structure, you can train an EAGLE checkpoint in the same way. Note that depending on your use case and dataset, you might see lower acceptance rate with the flat draft structure compared with using a tree structure.

NxD Inference supports EAGLE models with or without input normalization. By default, NxD Inference expects that the EAGLE model doesn't use input normalization. To use an EAGLE model with input normalization, set `enable_eagle_draft_input_norm` to `True` in `NeuronConfig`.

You can find links to pretrained EAGLE draft model checkpoints for various popular models in the official EAGLE repository on GitHub: <https://github.com/SafeAILab/EAGLE>. However, these pretrained EAGLE model checkpoints don't include the LM head weights from the target model. To use these pretrained checkpoints with NxD Inference, you must first copy the LM head weights from the target to the draft model.

The following code demonstrates how to perform this operation for a `Llama-3.1-70B-Instruct` target model and the corresponding `EAGLE draft`:

```
import json
import os

import torch
from safetensors import safe_open
from safetensors.torch import save_file

target_model_path = "Meta-Llama-3.1-70B-Instruct"
draft_model_path = "Llama-3.1-70B-Instruct-EAGLE-Draft"

DRAFT_MODEL_SAFETENSORS_NAME = "model.safetensors"
LM_HEAD_WEIGHT_TENSOR_NAME = "lm_head.weight"
TARGET_MODEL_SAFETENSORS_INDEX_NAME = "model.safetensors.index.json"

def find_lm_head_safetensors_location(model_dir):
    model_index_location_path = os.path.join(model_dir, TARGET_MODEL_SAFETENSORS_INDEX_
↪NAME)

    with open(model_index_location_path, 'r') as f:
        model_index_locations = json.load(f)

    lm_head_safetensors_name = model_index_locations["weight_map"][LM_HEAD_WEIGHT_TENSOR_
```

(continues on next page)

(continued from previous page)

```

↪NAME]

    return lm_head_safetensors_name

# Find the target model `lm_head.weight` location in safetensors
target_lm_head_safetensors_name = find_lm_head_safetensors_location(target_model_path)
target_lm_head_safetensors_path = os.path.join(target_model_path, target_lm_head_
↪safetensors_name)

# Open the target model.safetensor containing `lm_head.weight`
with safe_open(target_lm_head_safetensors_path, framework="pt") as f:
    target_lm_head = f.get_tensor(LM_HEAD_WEIGHT_TENSOR_NAME)

# Collect all tensors in the draft model
draft_model_safetensors_path = os.path.join(draft_model_path, DRAFT_MODEL_SAFETENSORS_
↪NAME)
tensors = {}
with safe_open(draft_model_safetensors_path, framework="pt") as f:
    for key in f.keys():
        tensors[key] = f.get_tensor(key)

# Add the LM head weights and save out the new draft model.safetensors file
tensors[LM_HEAD_WEIGHT_TENSOR_NAME] = target_lm_head.type(torch.float16)
save_file(tensors, draft_model_safetensors_path)

```

## Fused Speculation

EAGLE speculation uses a feature called *fused speculation*, where the draft model and target model are fused into a single compiled model to improve performance. Fused speculation uses a different config called `FusedSpecNeuronConfig`, which specifies the model class, draft config, and draft model path to fuse with the target model.

## Example

```

import copy

from neuronx_distributed_inference.models.config import (
    FusedSpecNeuronConfig,
    NeuronConfig,
    OnDeviceSamplingConfig
)
from neuronx_distributed_inference.models.llama.modeling_llama import (
    NeuronLlamaForCausalLM,
    NeuronLlamaModel
)
from neuronx_distributed_inference.utils.accuracy import get_generate_outputs
from neuronx_distributed_inference.utils.hf_adapter import load_pretrained_config
from transformers import AutoTokenizer, GenerationConfig

prompt = "The future of AI is"

```

(continues on next page)

(continued from previous page)

```

model_path = "/home/ubuntu/models/Llama-3.1-70B-Instruct"
draft_model_path = "/home/ubuntu/models/Llama-3.1-70B-Instruct-EAGLE-Draft"
compiled_model_path = "/home/ubuntu/neuron_models/Llama-3.1-70B-Instruct-EAGLE"
max_sequence_length = 1024

# Initialize on-device sampling configuration.
on_device_sampling_config = OnDeviceSamplingConfig(
    temperature=0.7,
    top_k=50,
    top_p=1.0,
)

# Initialize model configuration.
neuron_config = NeuronConfig(
    # Neuron supports EAGLE batch sizes greater than 1.
    # We set batch size to 1 in this tutorial due to a
    # limitation in the transformers library for
    # generation with speculative decoding.
    # For more information, see: https://github.com/huggingface/transformers/issues/32165
    batch_size = 1,
    enable_eagle_speculation=True,
    enable_fused_speculation=True,
    max_context_length=max_sequence_length,
    max_length=max_sequence_length,
    on_device_sampling_config=on_device_sampling_config,
    seq_len=max_sequence_length,
    speculation_length=5,
    # For best performance, set to the maximum tensor
    # parallelism of your Neuron instance type.
    tp_degree=32,
    trace_tokengen_model=False
)

config = NeuronLlamaForCausalLM.get_config_cls()(
    neuron_config, load_config=load_pretrained_config(model_path)
)

# Initialize draft model configuration and set EAGLE-specific values.
draft_neuron_config = copy.deepcopy(neuron_config)
draft_neuron_config.trace_tokengen_model = True
draft_neuron_config.enable_fused_speculation = False
draft_neuron_config.is_eagle_draft = True
draft_neuron_config.sequence_parallel_enabled = False

draft_config = NeuronLlamaForCausalLM.get_config_cls()(
    draft_neuron_config, load_config=load_pretrained_config(draft_model_path))

# Initialize fused speculation configuration.
fused_spec_config = FusedSpecNeuronConfig(
    NeuronLlamaForCausalLM._model_cls,
    draft_config=draft_config,

```

(continues on next page)

(continued from previous page)

```

    draft_model_path=draft_model_path,
)
config.fused_spec_config = fused_spec_config

# Initialize model from configuration.
model = NeuronLlamaForCausalLM(model_path, config)

# Compile and save model.
model.compile(compiled_model_path)

# Load model to the Neuron device.
model.load(compiled_model_path)

# Load tokenizer and generation config.
tokenizer = AutoTokenizer.from_pretrained(model_path, padding_side=neuron_config.padding_
↪side)
generation_config = GenerationConfig.from_pretrained(model_path)
generation_config.max_length = 1024
# pad_token_id is required for Hugging Face assisted sampling.
generation_config.pad_token_id = tokenizer.eos_token_id

# Run generation and print outputs.
_, output_tokens = get_generate_outputs(
    model,
    [prompt],
    tokenizer,
    is_hf=False,
    # draft_model is not set here due to fused speculation.
    draft_model=None,
    generation_config=generation_config
)

print("Generated output:")
for _, output in enumerate(output_tokens):
    print(output)

```

## MoE model architecture support

NxD Inference supports mixture-of-experts (MoE) models. The library includes ready-to-use modeling code for Mixtral and DBRX. These models are built using reusable MoE modules from NeuronX Distributed Core: RouterTopK, ExpertMLPs, and MoE. You can use these modules to onboard additional MoE models.

NxD Inference also provides a helper function, `initialize_moe_module`, which you can use to initialize an MoE model's MLP module from these MoE modules. For examples of how to use this helper function, see the decoder layer module implementation in the [Mixtral](#) and [DBRX](#) modeling code.

## Grouped-query attention (GQA) support

NxD Inference provides a reusable attention module, `NeuronAttentionBase`, which you can use when onboarding models. This module is also used in NxD Inference modeling code like Llama and Mixtral.

NxD Inference supports the following sharding strategies for the KV cache used in the attention module:

- `CONVERT_TO_MHA` — Transforms a GQA attention mechanism into a traditional MHA mechanism by replicating the K/V heads to evenly match the corresponding Q heads. This consumes more memory than would otherwise be used with other sharding mechanisms but works in all cases.
- `REPLICATE_TO_TP_DEGREE` — Transforms a GQA attention mechanism such that there is exactly one K/V head per `tp_degree` through replication e.g. 8 K/V heads with `tp_degree=32` results in 32 K/V heads. This is more memory efficient but does not work for all configurations. Q heads are padded interleaved to retain correct alignment between Q and K/V heads.

The `NeuronAttentionBase` module uses `REPLICATE_TO_TP_DEGREE` by default. If the TP degree isn't divisible by the number of KV heads, `NeuronAttentionBase` uses `CONVERT_TO_MHA`.

## Asynchronous Runtime Support

NxD Inference offers certain model configurations to be run with Asynchronous Runtime Mode (Async mode). Async mode allows NxD Inference to parallelize CPU logic with Neuron (NEFF) logic. As a result, any CPU overheads within NxDI that exist between sequential model executions (ex. autoregressive loop in LLMs) are almost fully eliminated. This reduces latency anywhere from 5% to 20% based on the model configuration.

This feature can be enabled with by setting `async_mode` to `True` in `NeuronConfig`.

To use Async mode, a model configuration must meet the following prerequisites: - On-device sampling is enabled. - If speculation is enabled, fused speculation must also be enabled.

It is highly recommended to set `async_mode` to `True` for every other case, since it offers a latency reduction. Furthermore, this feature is a purely runtime feature, so if you have a previously compiled model, and its configuration doesn't fall under the unsupported case, `async_mode` will likely be able to improve performance.

---

**Note:** If you are using vLLM, this feature works independently of vLLM's Async Engine. As a result, `async_mode` can be enabled whether vLLM is used or not.

---

## Prefix Caching Support

Prefix caching is a performance optimization technique where prompts in multiple requests sharing the same prefix can reuse the previously computed KV cache. When context encoding a prompt that starts with a previously computed prefix, the encoding of the prefix tokens will be skipped and the corresponding KV Cache will be fetched and used for encoding the rest of the tokens (suffix). The performance benefit comes from the time saved by re-using the KV Cache instead of re-encoding the prefix tokens. NxD Inference supports prefix caching during context encoding. To store KV cache and match to prefix efficiently, NxD Inference uses block KV Cache layout for prefix caching. NxD Inference does not implement its own cache eviction, memory management, or prefix hashing for matches. Instead, it requires external management of the block KV cache and expects active block tables and slot mappings to be provided with each generation request. This feature integrates with vLLM by enabling automatic prefix caching, which manages the block tables, handles automatic prefix matching across prompts, and performs cache evictions. More on automatic prefix caching support on vLLM can be found [here](#).

To enable prefix caching with NxD Inference, set `is_prefix_caching=True` in `NeuronConfig` along with configurations for block KV cache layout.

```
neuron_config = NeuronConfig(
    is_prefix_caching=True,
    is_block_kv_layout=True,
    pa_num_blocks=1024,
    pa_block_size=32,
)
```

`is_block_kv_layout=True` and `is_prefix_caching=True` are set in `NeuronConfig` to enable block KV Cache layout and enable prefix caching. The first two dimensions of the KV cache are set to the number of blocks and block size, respectively. These configurations are specified using `pa_num_blocks` and `pa_block_size` in `NeuronConfig`. For optimal performance with Neuron, it's recommended to set `pa_block_size=32`. The minimum required `pa_num_blocks` to be initialized is  $(\text{batch\_size} * \text{max\_model\_len}) / \text{pa\_block\_size}$ . However, it is recommended to initialize more blocks than the required minimum to accommodate caching of common prefixes. The higher the number of blocks, the greater the likelihood of cache hits, as fewer cache evictions will occur. NxD Inference does not currently provide an automated solution to determine the maximum number of KV Cache blocks that can be initialized in HBM without exceeding available memory space. Customers are advised to experiment with increasing the number of blocks that balances the cache hit rate and memory taken. Any memory taken by increasing the cache will impact the batch sizes and sequence lengths that can be supported, so customers are suggested to pick the correct number of blocks considering these trade offs and the specific inference workload they plan to run in production.

NxD Inference does not use paged attention for prefix caching. Instead, it follows a different process: first gathering the block KV cache using the block table, then converting it to a flat KV cache layout, computing attention, and finally scattering the computed cache back to the block KV cache layout. This approach introduces overhead during token generation requests due to layout conversions, which can negatively impact performance as the `max_model_len` increases.

## Bucketing with Prefix Caching

Prefix caching handles both the prefix (cache hit) and suffix (no cache hit) portions of input prompts during context encoding. A two-dimensional bucketing system has been introduced to support context encoding when prefix caching is enabled. This system uses separate dimensions corresponding to the prefix and suffix (non cache-hit portion) of the input prompts. In contrast, token generation still uses one-dimensional bucketing based on the maximum sequence length.

When bucketing is enabled, NxD Inference creates prefill (suffix) buckets (covering suffix portion) starting with powers of 2, ranging from 512 up to the maximum context length. The prefix buckets mirror the prefill buckets, with one key difference: a special prefix bucket of size 0 is added to handle requests with no cache hits. NxD Inference then creates a two-dimensional grid of all prefill and prefix bucket combinations, which represents the effective set of buckets during context encoding. During request processing, NxD Inference first identifies the smallest prefill bucket that can accommodate the largest suffix portion of the input prompts. If prefill padding is needed, NxD Inference prioritizes moving tokens from the prefix's end to the prefill bucket before adding padding. It then determines the smallest prefix bucket that can fit the largest prefix across prompts. These two dimensions together determine the final (prefill, prefix) bucket combination used to serve the context encoding request.

You can configure specific buckets to optimize inference based on the expected distribution of prefix lengths, input lengths, and output lengths for your model. In `NeuronConfig`, set `enable_bucketing=True`, and provide a list of bucket sizes in `context_encoding_buckets`, `prefix_buckets` and/or `token_generation_buckets`. `context_encoding_buckets` corresponds to prefill buckets when prefix caching is enabled.

```
neuron_config = NeuronConfig(
    enable_bucketing=True,
    context_encoding_buckets=[512, 1024, 2048],
    prefix_buckets=[512, 1024]
```

(continues on next page)

(continued from previous page)

```
token_generation_buckets=[2048]
)
```

## Examples

For `context_encoding_buckets=[512, 1024, 2048]` and `prefix_buckets=[512, 1024]`

For requests with:

- Input prompt of size 1000 with no prefix, NxDI uses prefill bucket as 1024 and prefix bucket as 0.
- Input prompt of size 800 with 128 as the prefix size, and remaining 672 as the suffix size, NxDI first selects 1024 as the prefill bucket. Remaining 352 prefill slots are filled up by moving entire prefix to the suffix part. So prefill bucket of 1024 and prefix bucket as 0 is used here.
- Input prompt of size 900 with 640 as the prefix size, and remaining 260 as the suffix size, NxDI first selects 512 as the prefill bucket. Remaining 252 prefill slots are filled up by moving 252 tokens from the end of prefix to the suffix part. Effective prefix length now becomes 388, so prefill bucket of 512 and prefix bucket as 512 is used.
- Input prompt of size 1600 with 1280 as the prefix size and remaining 320 as the suffix size, NxDI selects 512 as the prefill bucket. Remaining 192 prefill slots are filled up by moving 192 tokens from the end of prefix to the suffix part. Effective prefix length now becomes 1088 which is larger than the largest prefix bucket of 1024. This leads to exception during the request processing.

The two-dimensional bucketing system exponentially increases the number of context encoding buckets. Therefore, users should exercise caution when using auto-bucketing with large context lengths. It is recommended to limit the granularity of prefix buckets based on your specific workload requirements.

For detailed examples of prefix caching with NxI Inference and vLLM, see [nxdi-trn2-llama3.3-70b-apc-tutorial](#).

## Multi-LoRA Serving

NxD Inference supports serving with multiple LoRA adapters and users can specify different LoRA adapters for their requests at runtime. It also supports multi-LoRA serving with vLLM as the frontend. NxI Inference currently supports loading of LoRA adapters at server startup for dense models like Llama-3.3-70B. Dynamic loading of LoRA adapters at runtime is not currently supported and will be supported in a future Neuron release.

### Enable multi-LoRA serving

To enable multi-LoRA serving, provide a `LoraServingConfig` for `lora_config` attribute in `NeuronConfig`.

```
lora_config = LoraServingConfig(
    max_loras=max_loras,
    lora_ckpt_paths=lora_ckpt_paths,
)
neuron_config = NeuronConfig(lora_config=lora_config)
```

Refer to [NxI Inference API Reference](#) for more details of `LoraServingConfig`.

NxD Inference primarily supports the format of LoRA adapters from [Huggingface PEFT](#). Each checkpoint path is a folder that contains a checkpoint file (.safetensors, .bin, or .pt) and a configuration json file (.json). In addition, NxI inference also supports LoRA adapters trained from [NxI LoRA finetuning](#). Each checkpoint path is a checkpoint file (.pt) that includes both LoRA adapter weights and the configuration.



NxD Inference assumes all the LoRA adapters for multi-LoRA serving are available locally during compilation and their weights are loaded on neuron devices during serving. When uploading a LoRA adapter checkpoint to NxDI for multi-LoRA serving, the user is required to name the adapter with a unique adapter ID, which will be used by users to specify the LoRA adapter for serving at runtime and by NxDI for model compilation.

The number of the multiple LoRA adapters for serving is specified by `max_loras`. The set of LoRA adapters in NxD Inference are specified by `lora_ckpt_paths`, which is a dictionary with a key-value pair for each LoRA adapter. The key is the adapter ID named by the user and the value is the local path of the LoRA adapter checkpoint. For detailed examples of multi-LoRA serving in NxDI, see [Tutorial: Multi-LoRA serving for Llama-3.1-8B on Trn2 instances](#).

## Maximum number of LoRA adapters supported

The LoRA adapter size is much smaller than the base model, but its weights still consumes non-negligible on-device memory. The maximum number of LoRA adapters that can be concurrently supported in NxD Inference depends on the base model, the LoRA rank, the reserved memory size for LoRA adapters, and how the LoRA adapters are sharded across TP groups.

Suppose Trn1 instance is used for multi-LoRA serving and the reserved memory size on each neuron core for LoRA adapters is 2GB. Each LoRA adapter has two parts, LoRA A and LoRA B, and only one of them can be partitioned with tensor parallelism and the other is just Linear layer. We analyze the maximum number of LoRA adapters supported in NxD inference under two cases: the linear layer is duplicated and the linear layer is sharded. These two cases can be specified by `lora_shard_linear_layer` in `LoraServingConfig`.

### When the linear layer is duplicated

The weight size of a LoRA adapter on each device is around half of the total LoRA adapter size in this case. When the base model is Llama3.1 8B, the LoRA adapter checkpoint size with LoRA rank 16 in BF16 is around 170MB. Because  $2GB / (170MB / 2) = 23$ , the maximum number of concurrent LoRA adapters is 23. When the base model is Llama3.3 70B, the LoRA adapter checkpoint size with LoRA rank 16 in BF16 is around 830MB and we can set `max_loras=4`. We analyze the maximum number of LoRA adapters supported in NxD inference under two cases: the linear layer is duplicated and the linear layer is sharded. These two cases can be specified by `lora_shard_linear_layer` in `LoraServingConfig`.

| Model        | Reserved Memory size | LoRA rank | Maximum LoRAs |
|--------------|----------------------|-----------|---------------|
| Llama3.1 8B  | 2GB                  | 16        | 23            |
| Llama3.1 8B  | 2GB                  | 32        | 12            |
| Llama3.3 70B | 2GB                  | 16        | 4             |
| Llama3.3 70B | 2GB                  | 32        | 2             |

### When the linear layer is sharded

The linear layer in a LoRA adapter is sharded across neuron cores in a TP group at the cost of Allgather communication overhead in this case. The weight size of a LoRA adapter on each device is  $1/TP\_DEGREE$  of the total LoRA adapter size.

| Model        | Reserved Memory size | LoRA rank | TP degree | Maximum LoRAs |
|--------------|----------------------|-----------|-----------|---------------|
| Llama3.1 8B  | 2GB                  | 16        | 32        | 376           |
| Llama3.1 8B  | 2GB                  | 32        | 32        | 188           |
| Llama3.3 70B | 2GB                  | 16        | 32        | 77            |
| Llama3.3 70B | 2GB                  | 32        | 32        | 38            |

## Disaggregated Inference [BETA]

Disaggregated Inference is an LLM serving architecture separates the prefill and decode phases of inference onto different hardware resources. Separating the compute intensive prefill phase from the memory bandwidth intensive decode phase can improve the LLM serving experience by

1. Removing prefill interruptions to decode from continuous batching to reduce inter token latency (ITL). These gains can be used to achieve higher throughput by running with a higher decode batch size while staying under Service Level Objectives (SLO).
2. Adapt to changing traffic patterns while still remaining under application SLOs.
3. Enable independent scaling of resources and parallelism strategies for prefill (compute bound) and decode (memory bound).

See the [Disaggregated Inference Developer Guide](#) and the [Disaggregated Inference Tutorial](#)

## NxD Inference - Production Ready Models

Neuronx Distributed Inference provides production ready models that you can directly use for seamless deployment. You can view the source code for all supported models in the [NxD Inference GitHub repository](#).

---

**Note:** If you are looking to deploy a custom model integration, you can follow the [model onboarding guide](#). You can refer to the source code for supported models in the [NxD Inference GitHub repository](#) and make custom changes required for your use case.

---

### Table of contents

- [Using Models to Run Inference](#)
  - [Using vLLM](#)
  - [Integrating Directly with NxD Inference](#)
- [Supported Model Architectures](#)
  - [Llama \(Text\)](#)
  - [Llama \(Multimodal\)](#)
  - [Mixtral](#)
  - [DBRX](#)
  - [Qwen2.5](#)

## Using Models to Run Inference

You can run models through vLLM or integrate directly with NxD Inference.

### Using vLLM

If you are using vLLM for production deployment, we recommend that you use the vLLM API to integrate with NxD Inference. The vLLM API automatically chooses the correct model and config classes based on the model's config file. For more information, refer to the [vLLM User Guide for NxD Inference](#).

### Integrating Directly with NxD Inference

To use NxD Inference directly, you construct model and configuration classes. For more information about which model and configuration classes to use for each model, see [Supported Model Architectures](#). To see an example of how to run inference directly with NxD Inference, see the [generation\\_demo.py](#) script.

## Supported Model Architectures

NxD Inference currently provides support for the following model architectures.

### Llama (Text)

NxD Inference supports Llama text models. The Llama model architecture supports all Llama text models, including Llama 2, Llama 3, Llama 3.1, Llama 3.2, and Llama 3.3. You can also use the Llama model architecture to run any model based on Llama, such as Mistral.

### Neuron Classes

- Neuron config class: `MoENeuronConfig`
- Inference config class: `MixtralInferenceConfig`
- Causal LM model class: `NeuronMixtralForCausalLM`

### Compatible Checkpoint Examples

- <https://huggingface.co/meta-llama/Llama-3.1-405B-Instruct> (requires Trn2)
- <https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct>
- <https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct>
- <https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>
- <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

## **Llama (Multimodal)**

NxD Inference supports Llama 3.2 multimodal models. You can use HuggingFace checkpoints or the original Meta checkpoints. To use the Meta checkpoint, you must first convert the checkpoint to Neuron format. For more information about how to run Llama3.2 multimodal inference, and for details about how to convert the original Meta checkpoints to run on NxD Inference, see *[Tutorial: Deploying Llama3.2 Multimodal Models](#)*.

### **Neuron Classes**

- Neuron config class: `MultimodalVisionNeuronConfig`
- Inference config class: `MllamaInferenceConfig`
- Causal LM model class: `NeuronMllamaForCausalLM`

### **Compatible Checkpoint Examples**

- <https://huggingface.co/meta-llama/Llama-3.2-11B-Vision-Instruct>
- <https://huggingface.co/meta-llama/Llama-3.2-90B-Vision-Instruct>

## **Mixtral**

NxD Inference supports models based on the Mixtral model architecture, which uses mixture-of-experts (MoE) architecture.

### **Neuron Classes**

- Neuron config class: `MoENeuronConfig`
- Inference config class: `MixtralInferenceConfig`
- Causal LM model class: `NeuronMixtralForCausalLM`

### **Compatible Checkpoint Examples**

- <https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

## **DBRX**

NxD Inference supports models based on the DBRX model architecture, which uses mixture-of-experts (MoE) architecture.

## Neuron Classes

- Neuron config class: `DbrxNeuronConfig`
- Inference config class: `DbrxInferenceConfig`
- Causal LM model class: `NeuronDbrxForCausalLM`

## Compatible Checkpoint Examples

- <https://huggingface.co/databricks/dbrx-instruct>

## Qwen2.5

NxD Inference supports models based on the Qwen2.5 model architecture.

## Neuron Classes

- Neuron config class: `Qwen2NeuronConfig`
- Inference config class: `Qwen2InferenceConfig`
- Causal LM model class: `NeuronQwen2ForCausalLM`

## Compatible Checkpoint Examples

- <https://huggingface.co/Qwen/Qwen2.5-72B-Instruct>
- <https://huggingface.co/Qwen/Qwen2.5-32B-Instruct>
- <https://huggingface.co/Qwen/Qwen2.5-14B-Instruct> (Not tested, but expected to work out of the box)
- <https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>
- <https://huggingface.co/Qwen/Qwen2.5-3B-Instruct> (Not tested, but expected to work out of the box)
- <https://huggingface.co/Qwen/Qwen2.5-1.5B-Instruct> (Not tested, but expected to work out of the box)
- <https://huggingface.co/Qwen/Qwen2.5-0.5B-Instruct>

## Onboarding models to run on NxD Inference

This guide covers how to onboard a model to get it to run on NxD Inference for the first time. To learn more about how to optimize a model on Neuron, see the *NxD Inference Features Configuration Guide*.

### Table of contents

- *Overview*
  - 1. Define a `NeuronConfig` class
  - 2. Define an `InferenceConfig` class
  - 3. Define a Neuron model

- 4. *Define an application/task head*
- 5. *Convert weights to a supported format*
- *Integrating Onboarded Model with vLLM*
- *Evaluating Models on Neuron*
  - *Logit Matching*
  - *Token Matching*
  - *Benchmarking*
  - *Profiling Models*
  - *Evaluating Models with the Inference Demo Script*
- *Debugging Models on Neuron*
- *Writing Tests on Neuron*

## Overview

This guide demonstrates how to adapt an existing PyTorch model to run on Neuron with the NeuronX Distributed (NxD) Inference library. At a high-level, you will do the following:

1. Define configuration classes. NxD Inference models include a `NeuronConfig`, which defines Neuron-specific configuration parameters, and an `InferenceConfig`, which defines model configuration parameters. When adapting a model that works with HuggingFace, `InferenceConfig` is synonymous to `PretrainedConfig`.
2. Define model classes. When you define model classes, you replace linear layers with parallel layers that are optimized for distributed inference on Neuron. NxD Inference also provides modules for attention, KV cache management, and more, which you can use to write model classes that work with Neuron. Model classes are compiled to run effectively on Neuron.
3. Define application heads. Application heads orchestrate passing inputs to the correct compiled model. Application heads also provide the interface to compile and load the model.
4. Convert weights to a supported format. NxD Inference supports safetensors and pickle formats.

### 1. Define a `NeuronConfig` class

Define a Neuron configuration class, which extends `NeuronConfig`. `NeuronConfig` includes Neuron-specific configuration parameters. In the config class for your model, you can define any additional Neuron-specific configuration parameters that your model requires.

- For MoE models, you can extend `MoENeuronConfig` instead of `NeuronConfig`. This class includes configuration parameters specific to MoE models.

```
from neuronx_distributed_inference.models.config import NeuronConfig

class NeuronLlamaConfig(NeuronConfig):
    def __init__(self, **kwargs):
        super().__init__(**kwargs)
        # Set any args/defaults
```

## 2. Define an InferenceConfig class

Define an inference configuration class, which extends `InferenceConfig`. `InferenceConfig` includes model parameters, such as those from a HuggingFace `PretrainedConfig` (like `LlamaConfig`). When users initialize your config, they can provide required attributes directly, or they can populate the config from a HuggingFace `PretrainedConfig`. You can also override `get_required_attributes` to enforce that certain attributes are present.

```
from neuronx_distributed_inference.models.config import InferenceConfig, NeuronConfig

class LlamaInferenceConfig(InferenceConfig):
    def get_required_attributes(self) -> List[str]:
        return [
            "hidden_size",
            "num_attention_heads",
            "num_hidden_layers",
            "num_key_value_heads",
            "pad_token_id",
            "vocab_size",
            "max_position_embeddings",
            "rope_theta",
            "rms_norm_eps",
            "hidden_act",
        ]

    @classmethod
    def get_neuron_config_cls(cls) -> Type[NeuronConfig]:
        return NeuronLlamaConfig
```

## 3. Define a Neuron model

Define a Neuron model. This class is a subclass of `NeuronBaseModel`, which is a PyTorch module.

1. In this class, you provide implementations for `setup_attr_for_model(self, config)` and `init_model(self, config)`.
  1. In `setup_attr_for_model`, set values for the following attributes. You can set these attributes from values in `config` and `config.neuron_config`.
    1. `self.on_device_sampling`
    2. `self.tp_degree`
    3. `self.hidden_size`
    4. `self.num_attention_heads`
    5. `self.num_key_value_heads`
    6. `self.max_batch_size`
    7. `self.buckets`
  2. In `init_model`, initialize the modules that make up the model.
    1. For attention modules, extend `NeuronAttentionBase`, which provides a group query attention (GQA) implementation adapted to Neuron.

2. Replace linear layers (such as in attention and MLP) with Neuron parallel layers (RowParallelLinear and ColumnParallelLinear).
  1. For more information about RowParallelLinear and ColumnParallelLinear layers, see [Tensor Parallelism Overview](#).
3. Replace embeddings with Neuron parallel embeddings (ParallelEmbedding).
4. Replace any other modules that require Neuron-specific implementations.

Note: This example demonstrates a simplified version of NeuronLlamaModel from the NxDI model hub.

```
from torch import nn
from transformers.activations import ACT2FN

from neuronx_distributed.parallel_layers import parallel_state
from neuronx_distributed.parallel_layers.layers import ColumnParallelLinear,
↳RowParallelLinear, ParallelEmbedding

from neuronx_distributed_inference.models.model_base import NeuronBaseModel
from neuronx_distributed_inference.modules.attention.attention_base import
↳NeuronAttentionBase
from neuronx_distributed_inference.modules.attention.utils import RotaryEmbedding
from neuronx_distributed_inference.modules.custom_calls import CustomRMSNorm

class NeuronLlamaMLP(nn.Module):
    """
    This class just replace the linear layers (gate_proj, up_proj and down_proj) with
    ↳column and row parallel layers
    """

    def __init__(self, config: InferenceConfig):
        super().__init__()
        self.config = config
        self.neuron_config = config.neuron_config
        self.tp_degree = config.neuron_config.tp_degree
        self.hidden_size = config.hidden_size
        self.intermediate_size = config.intermediate_size
        self.act_fn = ACT2FN[config.hidden_act]

        self.gate_proj = ColumnParallelLinear(
            self.hidden_size,
            self.intermediate_size,
            bias=False,
            gather_output=False,
            dtype=config.neuron_config.torch_dtype,
            pad=True,
        )
        self.up_proj = ColumnParallelLinear(
            self.hidden_size,
            self.intermediate_size,
            bias=False,
            gather_output=False,
            dtype=config.neuron_config.torch_dtype,
            pad=True,
```

(continues on next page)



(continued from previous page)

```

    )
    self.down_proj = RowParallelLinear(
        self.intermediate_size,
        self.hidden_size,
        bias=False,
        input_is_parallel=True,
        dtype=config.neuron_config.torch_dtype,
        pad=True,
    )

    def forward(self, x):
        return self.down_proj(self.act_fn(self.gate_proj(x)) * self.up_proj(x))

class NeuronLlamaAttention(NeuronAttentionBase):
    """
    Compared with LlamaAttention, this class just
    1. replaces the q_proj, k_proj, v_proj with column parallel layer
    2. replaces the o_proj with row parallel layer
    3. update self.num_head to be self.num_head / tp_degree
    4. update self.num_key_value_heads to be self.num_key_value_heads / tp_degree
    5. update forward() method to adjust to changes from self.num_head
    """

    def __init__(self, config: InferenceConfig):
        super().__init__()

        self.config = config
        self.neuron_config = config.neuron_config
        self.hidden_size = config.hidden_size
        self.num_attention_heads = config.num_attention_heads
        self.num_key_value_heads = config.num_key_value_heads
        self.head_dim = self.hidden_size // self.num_attention_heads
        self.max_position_embeddings = config.max_position_embeddings
        self.rope_theta = config.rope_theta
        self.padding_side = config.neuron_config.padding_side
        self.torch_dtype = config.neuron_config.torch_dtype

        self.tp_degree = parallel_state.get_tensor_model_parallel_size()

        self.fused_qkv = config.neuron_config.fused_qkv
        self.clip_qkv = None

        self.init_gqa_properties()
        self.init_rope()

    def init_rope(self):
        self.rotary_emb = RotaryEmbedding(
            self.head_dim,
            max_position_embeddings=self.max_position_embeddings,
            base=self.rope_theta,
        )

```

(continues on next page)

(continued from previous page)

```

class NeuronLlamaDecoderLayer(nn.Module):
    """
    Just replace the attention with the NXD version, and MLP with the NXD version
    """

    def __init__(self, config: InferenceConfig):
        super().__init__()
        self.hidden_size = config.hidden_size
        self.self_attn = NeuronLlamaAttention(config)
        self.mlp = NeuronLlamaMLP(config)
        self.input_layernorm = CustomRMSNorm(
            config.hidden_size,
            eps=config.rms_norm_eps,
        )
        self.post_attention_layernorm = CustomRMSNorm(
            config.hidden_size,
            eps=config.rms_norm_eps,
        )

    def forward(
        self,
        hidden_states: torch.Tensor,
        attention_mask: Optional[torch.Tensor] = None,
        position_ids: Optional[torch.LongTensor] = None,
        past_key_value: Optional[Tuple[torch.Tensor]] = None,
        **kwargs,
    ) -> Tuple[torch.FloatTensor, Optional[Tuple[torch.FloatTensor, torch.FloatTensor]]]:
        residual = hidden_states
        hidden_states = self.input_layernorm(hidden_states)

        # Self Attention
        attn_outs = self.self_attn(
            hidden_states=hidden_states,
            attention_mask=attention_mask,
            position_ids=position_ids,
            past_key_value=past_key_value,
            **kwargs,
        )

        hidden_states, present_key_value = attn_outs
        hidden_states = residual + hidden_states

        # Fully Connected
        residual = hidden_states
        hidden_states = self.post_attention_layernorm(hidden_states)
        hidden_states = self.mlp(hidden_states)
        hidden_states = residual + hidden_states

        return (hidden_states, present_key_value)

```

(continues on next page)

(continued from previous page)

```

class NeuronLlamaModel(NeuronBaseModel):
    """
    The neuron version of the LlamaModel
    """

    def setup_attr_for_model(self, config: InferenceConfig):
        # Needed for init_inference_optimization()
        self.on_device_sampling = config.neuron_config.on_device_sampling_config is not None

        self.tp_degree = config.neuron_config.tp_degree
        self.hidden_size = config.hidden_size
        self.num_attention_heads = config.num_attention_heads
        self.num_key_value_heads = config.num_key_value_heads
        self.max_batch_size = config.neuron_config.max_batch_size
        self.buckets = config.neuron_config.buckets

    def init_model(self, config: InferenceConfig):
        self.padding_idx = config.pad_token_id
        self.vocab_size = config.vocab_size

        self.embed_tokens = ParallelEmbedding(
            config.vocab_size,
            config.hidden_size,
            self.padding_idx,
            dtype=config.neuron_config.torch_dtype,
            shard_across_embedding=True,
            # We choose to shard across embedding dimension because this stops XLA from
            # rank specific constant parameters into the HLO. We could shard across
            # vocab, but that
            # would require us to use non SPMD parallel_model_trace.
            pad=True,
        )
        self.lm_head = ColumnParallelLinear(
            config.hidden_size,
            config.vocab_size,
            bias=False,
            pad=True,
        )

        self.layers = nn.ModuleList(
            [NeuronLlamaDecoderLayer(config) for _ in range(config.num_hidden_layers)]
        )
        self.norm = CustomRMSNorm(config.hidden_size, eps=config.rms_norm_eps)

```

#### 4. Define an application/task head

Define an application/task head. Applications includes causal LM, classification, and so on. This class extends a task-specific Neuron application head class (such as `NeuronBaseForCausalLM`), or the general `NeuronApplicationHead` class.

1. In this class, you provide an value for `_model_cls` which is the Neuron model class you defined.
2. You can also override any other functions as needed for your model, such as `get_compiler_args(self)` or `convert_hf_to_neuron_state_dict(model_state_dict, neuron_config)`.

Note: This example demonstrates a simplified version of `NeuronLlamaForCausalLM` from the NxD Inference model hub.

```
class NeuronLlamaForCausalLM(NeuronBaseForCausalLM):
    _model_cls = NeuronLlamaModel

    @classmethod
    def get_config_cls(cls):
        return LlamaInferenceConfig
```

NxD Inference offers *Asynchronous Runtime Support* as an alternative method to executing NEFFs in parallel with CPU Logic. To evaluate if your task can utilize `async_mode`, the following questions must be answered:

1. **Does your task repeatedly execute a model for a single user request? If not, then `async_mode` won't offer any benefits.**
  - Example: The Auto Regressive loops used in LLMs perform repeated execution of models for a given prompt, which can get some benefits from async mode.
2. **Does the output of one execution get passed onto the next execution without manipulation? If not, then `async_mode` is incompatible.**
  - NOTE: It might be possible to address this by moving some manipulation logic within the neff.
  - Example: For LLMs using on-device-sampling, we pass in the token generated as output as input to the next step in the auto regressive loop directly. Without on-device-sampling, the sampling logic will rely on logits as output, which is a data dependent compute pattern that is incompatible with async mode.
3. **Is there sufficient CPU logic that is independent of the previous outputs? If not, then `async_mode` likely won't offer major benefits.**
  - Example: In production workloads, these are typically server overheads (scheduling, logging, etc.), but this could also be some pre/post processing steps in the model execution pipeline.

Based on the answers above, `async_mode` will need to be set accordingly, and/or, be configured to work correctly with the application.

#### 5. Convert weights to a supported format

NxD Inference supports weights stored in the model path in the following formats:

| Format      | Sharded | File name                    |
|-------------|---------|------------------------------|
| Safetensors | No      | model.safetensors            |
| Safetensors | Yes     | model.safetensors.index.json |
| Pickle      | No      | pytorch_model.bin            |
| Pickle      | Yes     | pytorch_model.bin.index.json |

If your weights are in another format, you must convert them to one of these formats before you can compile and load the model to Neuron. See the following references for more information about these formats:

- Safetensors:
  - <https://github.com/huggingface/safetensors>
  - <https://huggingface.co/docs/safetensors/en/convert-weights>
- Pickle:
  - <https://docs.python.org/3/library/pickle.html>

## Integrating Onboarded Model with vLLM

After completing the model onboarding in NxDI using the steps outlined in this guide, you can follow these steps to run that model through vLLM.

1. Register your application/task head class in vLLM within `_NEURON_SUPPORTED_MODELS` (see `vllm/model_executor/model_loader/neuronx_distributed.py`).
2. Use the local directory as `model_name_or_path` within vLLM which contains the model weights and the `config.json` file that works with your model's `InferenceConfig` class.
3. Pass in any custom `NeuronConfig` attributes by using the `override_neuron_config` attribute while initializing the vLLM engine.
4. Run offline/online inference to get the model working with vLLM.

## Evaluating Models on Neuron

NxD Inference provides tools that you can use to evaluate the accuracy and performance of the models that you onboard to Neuron.

### Logit Matching

The logit matching evaluation tool verifies that output logits are within certain tolerances of expected logits. With this evaluation tool, NxD Inference runs generation on the Neuron device. Then, it compares the output logits against expected logits, which you can provide or generate with the HuggingFace model on CPU.

During logit validation, if the output tokens diverge, then this process runs generation on Neuron again, using the tokens up to the point where it diverged. This process is performed repeatedly each time the output diverges, until the entire output matches. This process uses greedy sampling to choose the most likely token at each index.

Once all tokens match, this process compares the logits generated on Neuron with the expected logits. If all logits are within expected tolerances, this accuracy check passes. Divergence difference tolerance is used to compare the logits at the token that diverges. Absolute and relative tolerance are used to compare the values of the logits for the top k highest scoring tokens. For best results, use a lower relative tolerance for smaller k values, and a higher relative tolerance for larger k values. A top k of `None` means to compare logits for all possible tokens at each index.

Logit matching uses the following tolerances by default, and you can customize these tolerances.

- Divergence difference tolerance: `0.001`
- Absolute tolerance:
  - Top k = 5: `1e-5`
  - Top k = 50: `1e-5`

- Top k = 1000: 1e-5
- Top k = None: 1e-5
- Relative tolerance:
  - Top k = 5: 0.01
  - Top k = 50: 0.02
  - Top k = 1000: 0.03
  - Top k = None: 0.05

If all logits are within expected thresholds, this accuracy check passes.

- Note: Logit matching cannot be used with on-device sampling.
- Note: Generating HuggingFace model outputs on CPU can take a significant amount of time for larger models or large sequence lengths.

### Example (check\_accuracy\_logits API)

```
from neuronx_distributed_inference.utils.accuracy import check_accuracy_logits

# Init Neuron model, HuggingFace tokenizer, and HuggingFace generation config.

check_accuracy_logits(
    model,
    tokenizer,
    generation_config,
)
```

### Token Matching

The token matching evaluation tool verifies that output tokens match expected tokens. With this evaluation tool, NeuronX Distributed Inference runs generation on the Neuron device. Then, it compares the output against expected tokens, which you can provide or generate with the HuggingFace model on CPU. If all tokens match, this accuracy check passes.

- Warning: Token mismatches are acceptable in many scenarios, especially with large models or large sequence lengths. This tool should only be used for small models and small sequence lengths.
- Note: Generating HuggingFace model outputs on CPU can take a significant amount of time for larger models or large sequence lengths.

### Example (check\_accuracy API)

```
from neuronx_distributed_inference.utils.accuracy import check_accuracy

# Init Neuron model, HuggingFace tokenizer, and HuggingFace generation config.

check_accuracy(
    model,
    tokenizer,
```

(continues on next page)

(continued from previous page)

```

    generation_config,
)

```

## Benchmarking

NxD Inference provides a benchmarking tool that evaluates the latency and throughput of a Neuron model and its sub-models (context encoding, token generation, etc.).

### Example (benchmark\_sampling API)

```

from neuronx_distributed_inference.utils.benchmark import benchmark_sampling

# Init Neuron model and HuggingFace generation config.

benchmark_sampling(model, generation_config)

```

### Example benchmarking result

```

{
  "e2e_model": {
    "latency_ms_p50": 28890.24031162262,
    "latency_ms_p90": 28977.734088897705,
    "latency_ms_p95": 28983.17071199417,
    "latency_ms_p99": 29032.21325159073,
    "latency_ms_p100": 29044.473886489868,
    "latency_ms_avg": 28879.499554634094,
    "throughput": 283.66142510545984
  },
  "context_encoding_model": {
    "latency_ms_p50": 705.0175666809082,
    "latency_ms_p90": 705.3698301315308,
    "latency_ms_p95": 705.6618571281433,
    "latency_ms_p99": 705.8443236351013,
    "latency_ms_p100": 705.8899402618408,
    "latency_ms_avg": 705.0377488136292,
    "throughput": 5809.618005408024
  },
  "token_generation_model": {
    "latency_ms_p50": 27.20165252685547,
    "latency_ms_p90": 27.295589447021484,
    "latency_ms_p95": 27.324914932250977,
    "latency_ms_p99": 27.655515670776367,
    "latency_ms_p100": 32.74345397949219,
    "latency_ms_avg": 27.19622969277793,
    "throughput": 147.22298324644066
  }
}

```

## Profiling Models

Neuron provides a profiling tool, `neuron-profile`, which you can use to analyze the performance of a compiled Neuron model. For more information, see *Neuron Profile User Guide*.

## Evaluating Models with the Inference Demo Script

NxD Inference provides an `inference_demo` console script, which you can run from the environment where you install `neuronx_distributed_inference`.

Note: Before you can use a custom model with the `inference_demo`, you must add it to the `MODEL_TYPES` dictionary in `inference_demo.py`.

This script provides the following arguments to configure evaluation tools:

- `--check-accuracy-mode` - Provide one of the following values:
  - `token-matching` - Perform a token matching accuracy check.
  - `logit-matching` - Perform a logit matching accuracy check.
  - `skip-accuracy-check` - Do not perform an accuracy check.
- `--num-tokens-to-check` - The number of tokens to check when performing token matching or logit matching.
- `--expected-outputs-path` - The path to a file that contains tokens or logits to compare against for the accuracy check. This file must contain an object saved with `torch.save()`.
- `--benchmark` - Run benchmarking.
- `--on-cpu` - Run inference on CPU. To simulate tensor parallelism, initialize `inference_demo.py` with `torchrun`.

## Debugging Models on Neuron

When you debug models on Neuron, you can enable debug logging to view information about inputs and outputs of the `NeuronBaseForCausalLM` forward function, which calls the `NeuronBaseModel`'s forward function.

```
import logging

logging.getLogger().setLevel(logging.DEBUG)
```

Because the forward function of `NeuronBaseModel` is compiled, you cannot use `log/print` statements to debug code that is called from this forward function (or any other compiled code).

Debugging Neuron modeling code on CPU isn't yet supported.



## Writing Tests on Neuron

NxD Inference provides tools to help you write unit and integration tests that validate your model works as expected. For more information, see *Testing modeling code with NxD Inference*.

## vLLM User Guide for NxD Inference

vLLM is a popular library for LLM inference and serving utilizing advanced inference features such as continuous batching. This guide describes how to utilize AWS Inferentia and AWS Trainium AI accelerators in vLLM by using NxD Inference (`neuronx-distributed-inference`).

### Table of contents

- *Overview*
- *Supported Models*
- *Setup*
  - *Prerequisite: Launch an instance and install drivers and tools*
  - *Installing the AWS Neuron fork of vLLM*
  - *Installing vLLM from vLLM main repository*
- *Usage*
  - *Neuron Framework Selection*
  - *Quickstart*
  - *Model Configuration*
  - *Scheduling and K/V Cache*
  - *Decoding*
  - *Quantization*
  - *Loading pre-compiled models / Serialization Support*
  - *Prefix Caching*
  - *Disaggregated Inference*
- *Examples*
  - *Offline Inference Example*
  - *Online Inference Example*

### Overview

NxD Inference integrates into vLLM by extending the model execution components responsible for loading and invoking models used in vLLM's LLMEngine (see [https://docs.vllm.ai/en/latest/design/arch\\_overview.html#llm-engine](https://docs.vllm.ai/en/latest/design/arch_overview.html#llm-engine) for more details on vLLM architecture). This means input processing, scheduling and output processing follow the default vLLM behavior.

Currently, we support continuous batching and streaming generation in the NxD Inference vLLM integration. We are working with the vLLM community to enable support for other vLLM features like PagedAttention and Chunked Prefill on Neuron instances through NxD Inference in upcoming releases.

### Supported Models

Refer to *Supported Model Architectures* for a list of models supported in vLLM through NxD Inference.

If you are adding your own model to NxD Inference, please see *Integrating Onboarded Model with vLLM* for instructions on how to setup vLLM integration for it.

### Setup

Before installing vLLM with the instructions below, you need to install the Neuron SDK.

#### Prerequisite: Launch an instance and install drivers and tools

Before installing vLLM with the instructions below, you will first need to launch an Inferentia or Trainium instance and install the necessary Neuron drivers and tools. Refer to *these setup instructions* for different ways to prepare your environment, including using Neuron DLAMIs and Neuron DLCs for quick setups.

### Installing the AWS Neuron fork of vLLM

We maintain a fork of vLLM that supports the latest features for NxD Inference.

#### Quickstart using Docker

Users can now use a preconfigured Deep Learning Container (DLC) with the AWS Neuron fork of vLLM pre-installed. Refer to the `vllm-inference-neuronx` container on <https://github.com/aws-neuron/deep-learning-containers> to get started.

#### Manually install from source

To manually install the AWS fork from source, use the following commands:

```
git clone -b neuron-2.24-vllm-v0.7.2 https://github.com/aws-neuron/upstreaming-to-vllm.  
↪ git  
cd upstreaming-to-vllm  
pip install -r requirements-neuron.txt  
VLLM_TARGET_DEVICE="neuron" pip install -e .
```

---

**Note:** The current AWS Neuron fork of vLLM (`neuron-2.24-vllm-v0.7.2`) is based on vLLM 0.7.2 as some of the new features in 2.24 release like prefix caching are only tested with vLLM 0.7.2. We intend to upgrade and support the latest vLLM version in future Neuron releases.

---



---

**Note:** Starting in release 2.24, the Neuron initialization in vLLM code no longer enables sequence parallel by default. This is to ensure better compatibility with models and configurations where sequence parallelism is not well supported. If you previously relied on the Neuron vLLM code to specify sequence parallel, you may now see increased TTFT times. To re-enable sequence parallelism, you can pass `--override-neuron-config "{\n\"sequence_parallel_enabled\":true}"`.

---

## Installing vLLM from vLLM main repository

A prior version of Neuron SDK 2.23 NxD Inference support was upstreamed onto vLLM v0.9.0. Additional details can be found in vLLM docs [here](#).

To install the official vLLM repository with Neuron support, use the following commands. Only Neuron SDK 2.23 and prior features are currently available in the official vLLM repository. See Neuron SDK 2.23 artifacts [here](#). It is recommended to re-install `neuronx-distributed` and `neuronx-distributed-inference` libraries after installing vLLM to avoid dependency version incompatibilities.

```
git clone -b releases/v0.9.0 https://github.com/vllm-project/vllm.git
cd vllm
pip install -U -r requirements/neuron.txt
VLLM_TARGET_DEVICE="neuron" pip install -e .

pip install neuronx-distributed==0.12.12111
pip install neuronx-distributed-inference==0.3.5591
```

## Usage

### Neuron Framework Selection

---

**Note:** The Neuron integration for vLLM supports both Transformers NeuronX and NxD Inference libraries. Set the `VLLM_NEURON_FRAMEWORK` environment variable to `neuronx-distributed-inference` to use the NxD Inference library. Set the `VLLM_NEURON_FRAMEWORK` environment variable to `transformers-neuronx` to use the Transformers NeuronX library. Make sure you have the corresponding library installed before running vLLM. If you have both libraries installed, and the `VLLM_NEURON_FRAMEWORK` environment variable is not set, the NxD Inference library will be used by default.

---

If you are migrating from Transformers NeuronX to NxD Inference, you can refer to this [Migration Guide](#) for additional support.

## Quickstart

Here is a quick and minimal example to get running.

```
import os
os.environ['VLLM_NEURON_FRAMEWORK'] = "neuronx-distributed-inference"

from vllm import LLM, SamplingParams
llm = LLM(
    model="TinyLlama/TinyLlama-1.1B-Chat-v1.0",
    max_num_seqs=8,
    max_model_len=128,
    device="neuron",
    tensor_parallel_size=2)

prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# note that top_k must be set to lower than the global_top_k defined in
# the neuronx_distributed_inference.models.config.OnDeviceSamplingConfig
sampling_params = SamplingParams(top_k=10, temperature=0.8, top_p=0.95)

outputs = llm.generate(prompts, sampling_params)

for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```

## Model Configuration

NxD Inference models provide many configuration options. When using NxD Inference through vLLM, we configure the model with a default configuration that sets the required fields from vLLM settings. It is recommended that you do not override these configuration settings unless you need it.

```
neuron_config = dict(
    tp_degree=parallel_config.tensor_parallel_size,
    ctx_batch_size=1,
    batch_size=scheduler_config.max_num_seqs,
    max_context_length=scheduler_config.max_model_len,
    seq_len=scheduler_config.max_model_len,
    enable_bucketing=True,
    is_continuous_batching=True,
    quantized=False,
    torch_dtype=TORCH_DTYPE_TO_NEURON_AMP[model_config.dtype],
    padding_side="right"
)
```

If you want to add or change any settings, you can use vLLM's `override_neuron_config` setting. You provide the settings you want to override as dictionary (or JSON object when starting vLLM from the CLI) containing basic types

e.g. to disable auto bucketing (for illustration), use

```
override_neuron_config={
    "enable_bucketing":False,
}
```

or when launching vLLM from the CLI

```
--override-neuron-config "{\"enable_bucketing\":false}"
```

For more information on NxD Inference features, see *NxD Inference Features Configuration Guide* and *NxD Inference API Reference*.

## Scheduling and K/V Cache

We currently use a contiguous memory layout for the K/V cache instead of PagedAttention support in NxD Inference. We integrated into vLLMs block manager by setting the block size to the maximum length supported by the model and allocating one block per maximum number of sequences configured. However, the vLLM scheduler currently does not introspect the blocks associated to each sequence when (re-)scheduling running sequences. It requires an additional free block regardless of space available in the current block resulting in preemption. This would lead to a large increase in latency for the preempted sequence because it would be rescheduled in the context encoding phase. Since we ensure each block is big enough to fit the maximum model length, preemption is never needed in our current integration. Therefore, we disabled the preemption checks done by the scheduler in our fork. This significantly improves E2E performance of the Neuron integration.

## Decoding

On-device sampling is enabled by default, which performs sampling logic on the Neuron devices rather than passing the generated logits back to CPU and sample through vLLM. This allows us to use Neuron hardware to accelerate sampling and reduce the amount of data transferred between devices leading to improved latency.

However, on-device sampling comes with some limitations. Currently, we only support the following sampling parameters: `temperature`, `top_k` and `top_p` parameters. Other sampling parameters ([https://docs.vllm.ai/en/latest/dev/sampling\\_params.html](https://docs.vllm.ai/en/latest/dev/sampling_params.html)) are currently not supported through on-device sampling.

When on-device sampling is enabled, we handle the following special cases:

- When `top_k` is set to -1, we limit `top_k` to 256 instead.
- When `temperature` is set to 0, we use greedy decoding to remain compatible with existing conventions. This is the same as setting `top_k` to 1.

By default, on-device sampling utilizes a greedy decoding strategy to select tokens with the highest probabilities. You can enable a different on-device sampling strategy by passing a `on_device_sampling_config` using the override neuron config feature (see Model Configuration). It is strongly recommended to make use of the `global_top_k` configuration limiting the maximum value of `top_k` a user can request for improved performance.

## Quantization

NxD Inference supports quantization but has not yet been integrated with vLLMs configuration for quantization. If you want to use quantization, **do not** set vLLM's `--quantization` setting to `neuron_quant`. Keep it unset and use the Neuron configuration of the model to configure quantization of the NxD Inference model directly. For more information on how to configure and use quantization with NxD Inference incl. requirements on checkpoints, refer to [Quantization](#) in the NxD Inference Feature Guide.

## Loading pre-compiled models / Serialization Support

Tracing and compiling the model can take a non-trivial amount of time depending on model size e.g. a small-ish model of 15GB might take around 15min to compile. Exact times depend on multiple factors. Doing this on each server start would lead to unacceptable application startup times. Therefore, we support storing and loading the traced and compiled models.

Both are controlled through the `NEURON_COMPILED_ARTIFACTS` variable. When pointed to a path that contains a pre-compiled model, we load the pre-compiled model directly, and any differing model configurations passed in to the vllm API will not trigger re-compilation. If loading from the `NEURON_COMPILED_ARTIFACTS` path fails, then we will recompile the model with the provided configurations and store the results in the provided location. If `NEURON_COMPILED_ARTIFACTS` is not set, we will compile the model and store it under a `neuron-compiled-artifacts` subdirectory in the directory of your model checkpoint.

## Prefix Caching

Starting in Neuron SDK 2.24, prefix caching is supported on the AWS Neuron fork of vLLM. Prefix caching allows developers to improve TTFT by re-using the KV Cache of the common shared prompts across inference requests. See [Prefix Caching](#) for more information on how to enable prefix caching with vLLM.

## Disaggregated Inference

Starting in Neuron SDK 2.24, disaggregated inference is supported on the AWS Neuron fork of vLLM. This feature allows different hardware resources to separately perform the compute intensive prefill phase and the memory bandwidth intensive decode phase of inference, thereby removing the prefill-decode interference and improving Goodput. See [Disaggregated Inference](#) for more information on how to use disaggregated inference with vLLM.

## Examples

For a list of examples for using vLLM with Neuron, refer to [upstreaming-to-vllm/examples /offline\\_inference/](#) folder. Look for example scripts with the `neuron_` prefix. We provide examples for use cases such as [automatic prefix caching](#), [disaggregated inference](#), [speculative decoding with a draft model](#), [speculative decoding using EAGLE](#), [multimodal models](#), [multi-LoRA](#), [quantization](#), and more.

For more in depth NxD Inference tutorials that include vLLM deployment steps, refer to [Tutorials](#).

The following examples use `meta-llama/Llama-3.1-8B-Instruct` on a `Trn1.32xlarge` instance.

If you have access to the model checkpoint locally, replace `meta-llama/Llama-3.1-8B-Instruct` with the path to your local copy. Otherwise, you need to request access through HuggingFace and login via `huggingface-cli login` using a [HuggingFace user access token](#) before running the examples.

If you use a different instance type, you need to adjust the `tp_degree` according to the number of Neuron Cores available on your instance type (for more information see: [Tensor-parallelism support](#)).

## Offline Inference Example

Here is an example for running offline inference. *Bucketing* is only disabled to demonstrate how to override Neuron configuration values. Keeping it enabled generally delivers better performance.

```
import os
os.environ['VLLM_NEURON_FRAMEWORK'] = "neuronx-distributed-inference"

from vllm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(top_k=1)

# Create an LLM.
llm = LLM(
    model="meta-llama/Llama-3.1-8B-Instruct",
    max_num_seqs=4,
    max_model_len=128,
    override_neuron_config={
        "enable_bucketing": False,
    },
    device="neuron",
    tensor_parallel_size=32)

outputs = llm.generate(prompts, sampling_params)

for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```

## Online Inference Example

You can start an OpenAI API compatible server with the same settings as the offline example by running the following command:

```
VLLM_NEURON_FRAMEWORK='neuronx-distributed-inference' python -m vllm.entrypoints.openai.
↪api_server \
    --model="meta-llama/Llama-3.1-8B-Instruct" \
    --max-num-seqs=4 \
    --max-model-len=128 \
    --tensor-parallel-size=8 \
    --port=8080 \
    --device "neuron" \
    --override-neuron-config '{"enable_bucketing":false}'
```

In addition to the sampling parameters supported by OpenAI, we also support `top_k`. You can change the sampling parameters and enable or disable streaming.

```
from openai import OpenAI

# Client Setup
openai_api_key = "EMPTY"
openai_api_base = "http://localhost:8000/v1"

client = OpenAI(
    api_key=openai_api_key,
    base_url=openai_api_base,
)

models = client.models.list()
model_name = models.data[0].id

# Sampling Parameters
max_tokens = 1024
temperature = 1.0
top_p = 1.0
top_k = 50
stream = False

# Chat Completion Request
prompt = "Hello, my name is Llama "
response = client.chat.completions.create(
    model=model_name,
    messages=[{"role": "user", "content": prompt}],
    max_tokens=int(max_tokens),
    temperature=float(temperature),
    top_p=float(top_p),
    stream=stream,
    extra_body={'top_k': top_k}
)

# Parse the response
generated_text = ""
if stream:
    for chunk in response:
        if chunk.choices[0].delta.content is not None:
            generated_text += chunk.choices[0].delta.content
else:
    generated_text = response.choices[0].message.content

print(generated_text)
```



## Testing modeling code with NxD Inference

To ensure that your model is accurate and performant, we recommend that you write tests for your modules, functions, and models. Run your tests each time you make a code change to check that your modeling code continues to work as expected.

### Table of contents

- [\*Testing models on Neuron\*](#)
- [\*Testing modules and functions on Neuron\*](#)
  - [\*Building modules to run on Neuron\*](#)
  - [\*Building functions to run on Neuron\*](#)
  - [\*Validating module and function accuracy on Neuron\*](#)
  - [\*Examples\*](#)

## Testing models on Neuron

NxD Inference provides utilities that you can use to test the performance and accuracy of a full model end-to-end. The [\*check\\_accuracy\\_logits\*](#) tool validates the accuracy of a Neuron model's output logits over the full sequence length. NxD Inference also includes a benchmarking tool, [\*benchmark\\_sampling\*](#), that you can use to evaluate the performance of your model and its submodels. You can use these utilities to define integration tests that validate your model. For more information, see [\*Evaluating Models on Neuron\*](#).

## Testing modules and functions on Neuron

NxD Inference provides common test utilities to help you validate that modules and functions run correctly on Neuron. The `build_module` and `build_function` utilities help you convert modules and functions into Neuron models. Then, you can use the `validate_accuracy` function to validate that the Neuron model is accurate for given inputs. You can use these utilities to write unit tests for your modeling code.

## Building modules to run on Neuron

```
neuronx_distributed_inference.utils.testing.build_module(
    module_cls,
    example_inputs: List[Tuple[torch.Tensor]],
    module_init_kwargs: Dict = {},
    tp_degree: int = 1,
    compiler_args: Optional[str] = None,
    compiler_workdir: Optional[str] = None,
    checkpoint_path: Optional[str] = None,
)
```

Builds a module into a Neuron model. This function traces the module using the example inputs, which is a list of tuples where each item is a tensor. Then, it compiles the traced module to produce a Neuron model. Arguments:

- `module_cls`: The module class to compile.

- `example_inputs`: The list of example inputs to use to trace the module. This list must contain exactly one tuple of tensors.
- `tp_degree`: The TP degree to use. Defaults to 1.
- `module_init_kwargs`: The kwargs to pass when initializing the module.
- `compiler_args`: The compiler args to use.
- `compiler_workdir`: Where to save compiler artifacts. Defaults to a tmp folder with a UUID for uniqueness.
- `checkpoint_path`: The path to the checkpoint to load. By default, this function saves the module state dict to use as the checkpoint.

## Building functions to run on Neuron

```
neuronx_distributed_inference.utils.testing.build_function(  
    func: Callable,  
    example_inputs: List[Tuple[torch.Tensor]],  
    tp_degree: int = 1,  
    compiler_args: Optional[str] = None,  
    compiler_workdir: Optional[str] = None,  
)
```

Builds a function into a Neuron model. You can use `build_function` to test an individual function, such as a `top_k` sampling function.

See `build_module` for more information about common arguments. If the function has non-tensor inputs, you must convert it to a function that only takes tensor inputs. You can use `partial` to do this, where you provide the non-tensor inputs as constants in the partial function. This step is necessary because all inputs must be tensors in a Neuron model.

```
import torch  
  
from neuronx_distributed_inference.utils.testing import build_module  
  
def top_k(input: torch.Tensor, k: int, dim: int):  
    return torch.topk(input, k, dim)  
  
top_k_partial = partial(top_k, 1, 0)  
model = build_fuction(top_k_partial, example_inputs=[(torch.rand(4)),])  
output = model(torch.rand(4))
```

## Validating module and function accuracy on Neuron

```
neuronx_distributed_inference.utils.testing.validate_accuracy(  
    neuron_model,  
    inputs: List[Tuple],  
    expected_outputs: Optional[List] = None,  
    cpu_callable: Optional[Callable] = None,  
    assert_close_kwargs: Dict = {},  
)
```

Validates the accuracy of a Neuron model. This function tests that the model produces expected outputs, which you can provide and/or produce on CPU. To compare outputs, this function uses `torch_neuronx.testing.assert_close`. If the output isn't similar, this function raises an `AssertionError`. Arguments:

- `neuron_model`: The Neuron model to validate.
- `inputs`: The list of inputs to use to run the model. Each input is passed to the model's forward function.
- `expected_outputs`: The list of expected outputs for each input. If not provided, this function compares against the CPU output for each input.
- `cpu_callable`: The callable to use to produce output on CPU.
- `assert_close_kwargs`: The kwargs to pass to `torch_neuronx.testing.assert_close`.

## Examples

### Example: Basic module test

This example demonstrates how to validate the accuracy of a basic module with a single linear layer. In this example, we initialize the module separately on Neuron and CPU (using the `distributed` arg in `ExampleModule`). This flag enables us run a parallel linear layer on Neuron and compare it to a standard linear layer on CPU.

```
import torch

from neuronx_distributed_inference.utils.testing import build_module, validate_accuracy

# Module to test.
class ExampleModule(torch.nn.Module):
    def __init__(self, distributed):
        super().__init__()
        if distributed:
            self.linear = ColumnParallelLinear(
                input_size=SAMPLE_SIZE,
                output_size=SAMPLE_SIZE,
                bias=False,
                dtype=torch.float32,
            )
        else:
            self.linear = torch.nn.Linear(
                in_features=SAMPLE_SIZE,
                out_features=SAMPLE_SIZE,
                bias=False,
                dtype=torch.float32,
            )

    def forward(self, x):
        return self.linear(x)

def test_validate_accuracy_basic_module():
    inputs = [(torch.arange(0, SAMPLE_SIZE, dtype=torch.float32),)]
    example_inputs = [(torch.zeros((SAMPLE_SIZE), dtype=torch.float32),)]
```

(continues on next page)

(continued from previous page)

```

module_cpu = ExampleModule(distributed=False)
neuron_model = build_module(ExampleModule, example_inputs, module_init_kwargs={
↪ "distributed": True})

validate_accuracy(neuron_model, inputs, cpu_callable=module_cpu)

```

### Example: Basic function test

This example demonstrates how to validate the accuracy of a basic function with tensor args.

```

import torch

from neuronx_distributed_inference.utils.testing import build_function, validate_accuracy

def example_sum(tensor):
    return torch.sum(tensor)

def test_validate_accuracy_basic_function():
    inputs = [(torch.tensor([1, 2, 3], dtype=torch.float32),)]
    example_inputs = [(torch.zeros((3), dtype=torch.float32),)]

    neuron_model = build_function(example_sum, example_inputs)
    validate_accuracy(neuron_model, inputs, cpu_callable=example_sum)

```

### Additional examples

For additional examples of `build_module`, `build_function`, and `validate_accuracy`, see the [testing.py unit tests](#).

### Migrating from NxD Core inference examples to NxD Inference

We have migrated the NeuronX Distributed Core [examples/inference](#) folder to a separate package, NeuronX Distributed (NxD) Inference (`neuronx-distributed-inference`), so you can import and use it as a proper library. This new library, NxD Inference, includes production ready models that you can deploy out of the box with model inference backends, such as vLLM. This library also provides modules that you can use to implement your own models to run with the Neuron SDK.

If you use the inference examples from NxD Core, follow this guide to migrate to NxD Inference. For more information about NxD Inference and to see examples of how to use it, see [NxD Inference Features Configuration Guide](#), [NxD Inference Tutorials](#), and the `generation_demo.py` script.

**Warning:** Previous inference examples (including Llama 2, Llama 3, Mixtral, and DBRX) in the NxD Core GitHub repository were removed in Neuron Release 2.23. The models and example code are implemented in the NxD Inference library, so you can easily integrate them with your inference scripts. If you use these examples in NxD Core, we recommend that you update your inference scripts to use the NxD Inference model hub instead. If your use case requires you to directly integrate with the NxD Core library (and not NxD Inference) then you can

continue to use the NxD Core library directly. For an example of how to integrate with NxD Core directly, see the newer [Llama3.2 1B sample](#) added in Neuron Release 2.23. For more information, see [announce-eos-nxd-examples](#).

## Table of contents

- *Changes*
  - 1. *New config interface*
  - 2. *New base application interface*
  - 3. *New generation inference*
  - 4. *New quantization interface*
  - 5. *Inference demo script (replaces runners)*

## Changes

### 1. New config interface

NxD Inference includes a new model config interface, `InferenceConfig`, where `NeuronConfig` is an attribute within the model config, and the model config no longer extends HuggingFace's `PretrainedConfig`. NxDI includes an adapter for loading an HuggingFace's config into this model config. The configurations are serialized into a file named `neuron_config.json`.

**This change means that the config structure is inverted compared to the NxD examples folder.**

- To access the model config (similar to HuggingFace's `PreTrainedConfig`), use `config` (or `model.config, self.config`).
- To access the `NeuronConfig`, use `config.neuron_config` (or `model.neuron_config, self.neuron_config`).

To onboard a custom model, you define config classes that extend `InferenceConfig` and `NeuronConfig`. The following example from DBRX shows how to define a DBRX-specific `NeuronConfig` (`NeuronDbrxConfig`) and `InferenceConfig` (`DbrxInferenceConfig`). `DbrxInferenceConfig` that defines required config attributes and specifies that `NeuronDbrxConfig` is the `NeuronConfig` class. The required attributes are typically set by loading a `PretrainedConfig` (in this case, HuggingFace's `DbrxConfig`) into the `InferenceConfig`. Alternatively, a user can manually provide these attributes to avoid depending on an HuggingFace config class.

```
class NeuronDbrxConfig(MoENeuronConfig):
    def __init__(self, **kwargs):
        super().__init__(**kwargs)
        self.fused_qkv = True

class DbrxInferenceConfig(InferenceConfig):
    def get_required_attributes(self) -> List[str]:
        return [
            "d_model",
            "n_heads",
            "max_seq_len",
            "emb_pdrop",
```

(continues on next page)

(continued from previous page)

```

        "resid_pdrop",
        "pad_token_id",
        "vocab_size",
        "attn_config",
        "ffn_config",
    ]

    @classmethod
    def get_neuron_config_cls(cls):
        return NeuronDbrxConfig

```

**Note:** NeuronDbrxConfig extends MoENeuronConfig, which is a subclass of NeuronConfig that includes attributes that are specific to mixture-of-experts (MoE) models.

To load the config from an HuggingFace checkpoint or a compiled checkpoint, pass `load_pretrained_config(path)` as the `load_config` hook when you create the InferenceConfig.

```

from neuronx_distributed_inference.utils.hf_adapter import load_pretrained_config

neuron_config = DbrxNeuronConfig() # Provide args
config = DbrxInferenceConfig(
    neuron_config,
    load_config=load_pretrained_config(model_path),
)

```

To serialize the config, call `save(path)`.

```
config.save(compiled_model_path)
```

To deserialize the config, call `load(path)`.

```
config = DbrxInferenceConfig.load(compiled_model_path)
```

NeuronConfig also supports nested configs now. For example, see the OnDeviceSamplingConfig class and its integration into NeuronConfig.

## 2. New base application interface

NeuronApplicationBase takes general purpose features from NeuronBaseForCausalLM, such as compile and load, and makes them available in a new abstract base class. You can extend this base class to define other types of application heads, such as for image classification.

## 3. New generation inference

The Neuron model classes no longer extend HuggingFace's PretrainedModel, so they no longer include a HuggingFace generate() function. Additionally, GenerationConfig arguments are no longer passed through the model config. To run HuggingFace generation in NxD Inference, wrap the Neuron model in a HuggingFaceGenerationAdapter, and pass a GenerationConfig when you call generate().

```
from transformers import GenerationConfig

from neuronx_distributed_inference.utils.hf_adapter import HuggingFaceGenerationAdapter

# Init config, model, and tokenizer.

generation_config = GenerationConfig.from_pretrained(model_path)
generation_config_kwargs = {
    "do_sample": True,
    "top_k": 1,
    "pad_token_id": generation_config.eos_token_id,
    "max_length": neuron_config.max_length,
}
generation_config.update(**generation_config_kwargs)

inputs = tokenizer(prompts, padding=True, return_tensors="pt")
generation_model = HuggingFaceGenerationAdapter(model)
outputs = generation_model.generate(
    inputs.input_ids,
    generation_config=generation_config,
    attention_mask=inputs.attention_mask,
)
```

## 4. New quantization interface

This new base class also includes an interface for quantization, which was previously part of the run\_llama\_quantized.py example in the old NxD examples folder. The following example saves a quantized checkpoint for a Llama model. In this example, the config includes a neuron\_config with quantization enabled.

```
NeuronLlamaForCausalLM.save_quantized_state_dict(model_path, config)
```

## 5. Inference demo script (replaces runners)

In place of `runner.py` and various `run_x.py` examples, NxD-I provides an `inference_demo` console script. When you run the script, you provide a model path and configuration parameters to use for inference. This script includes benchmarking and accuracy checking features that you can use to verify that your models and modules work correctly.

The following example demonstrates how to run Llama-3-8b with token matching and benchmarking enabled.

```
inference_demo \
--model-type llama \
--task-type causal-lm \
run \
--model-path /home/ubuntu/model_hf/Llama-3.1-8b/ \
--compiled-model-path /home/ubuntu/traced_model/Llama-3.1-8b/ \
--torch-dtype bfloat16 \
--tp-degree 32 \
--batch-size 2 \
--max-context-length 32 \
--seq-len 64 \
--on-device-sampling \
--enable-bucketing \
--top-k 1 \
--do-sample \
--pad-token-id 2 \
--prompt "I believe the meaning of life is" \
--prompt "The color of the sky is" \
--check-accuracy-mode token-matching \
--benchmark
```

For additional examples, see the `neuronx-distributed-inference` GitHub repository: <https://github.com/aws-neuron/neuronx-distributed-inference>.

## Migrating from Transformers NeuronX to NeuronX Distributed(NxD) Inference

### Table of contents

- *How is writing modeling code different in NxD Inference?*
- *How can I migrate from Transformers NeuronX to use NxD Inference with vLLM?*
  - *Update Environment Variable to force vLLM to use NxD Inference*
  - *Compiling and loading the model*
  - *Features currently not supported in NxD Inference through vLLM*
- *Serialization support*
  - *Transformers NeuronX*
  - *NeuronX Distributed Inference*
- *Models supported in Transformers NeuronX and NxD Inference model hubs*
- *Onboarding custom or private models with NxD Inference*
- *Neuron Config Migration*



For customers who are currently using Transformers NeuronX, this migration guide explains the steps involved in migrating from Transformers NeuronX to NxD Inference library.

### How is writing modeling code different in NxD Inference?

In Transformers NeuronX, you write modeling code in HLO format using a Python HLO interface. In NeuronX Distributed Inference, you write modeling code in native PyTorch and Python, and the library converts it to HLO for you. This change makes it easier to develop models to run on Neuron, because you can start from existing Pytorch or Python modeling code.

### How can I migrate from Transformers NeuronX to use NxD Inference with vLLM?

Transformers NeuronX library currently supports Llama and Mistral model architectures with vLLM integration. If you are using one of these models, like Llama 3.1, Llama 3, Llama 2, or Mistral-7b-V2, you can migrate to use NxD Inference library with vLLM using the following steps:

#### Update Environment Variable to force vLLM to use NxD Inference

As vLLM currently supports both Transformers NeuronX and NeuronX Distributed Inference libraries for the Llama and Mistral models, you need to update the following environment variable in the inference scripts to force vLLM to use NxD Inference.

```
# Force vLLM framework to use neuronx-distributed-inference
os.environ['VLLM_NEURON_FRAMEWORK'] = "neuronx-distributed-inference"
```

#### Compiling and loading the model

Transformers NeuronX uses Neuron Persistent Cache to load a pre-compiled model so that there is no additional delay in compilation when loading the model on vLLM. NxD Inference currently does not support Neuron Persistent Cache but provides the following way to load a pre-compiled model in NeuronX Distributed Inference.

For production use cases where customer wants to avoid compiling the model in NxD Inference for the first time, users can set the environment variable `NEURON_COMPILED_ARTIFACTS` which points to pre-compiled artifacts directory to avoid the compilation time. If the artifacts are not present within the specified directory, then compilation of the model would be triggered as a fallback mechanism and will store the artifacts by default in `neuron-compiled-artifacts/{unique_hash}/`

#### Features currently not supported in NxD Inference through vLLM

NxD Inference doesn't yet support the following features that TNx supports in vLLM integration.

- Multi-Node Inference
- Persistent Cache
- concurrency > 1 support for speculation

Users can use exactly the same set of parameters to test out vLLM with NxD Inference library as they specify with Transformers NeuronX with the exception of `override_neuron_config`. Both Transformers NeuronX and NxD

Inference allows overriding available NeuronConfig, but not all NeuronConfig parameters that are available with Transformers NeuronX are still valid/applicable in NxD Inference. Refer to the [Neuron Config Migration](#) to migrate your `override_neuron_config` params from Transformers NeuronX to NxD Inference.

## Serialization support

In both libraries, you serialize the compiled model, so you can use the model in subsequent runs without compiling it each time.

In Transformers NeuronX, the save function does not serialize sharded weights by default, and you can enable this functionality with the `sharded_weights` flag. In NeuronX Distributed Inference, the `compile` function serializes sharded weights by default, and you can disable this functionality with the `save_sharded_checkpoint` flag in NeuronConfig.

## Transformers NeuronX

```
# Create and compile the Neuron model
neuron_config = NeuronConfig()
model_neuron = LlamaForSampling.from_pretrained(
    'openlm-research/open_llama_3b',
    batch_size=1,
    tp_degree=8,
    n_positions=128,
    neuron_config=neuron_config
)

# Compile the model.
model_neuron.to_neuron()

# Save the presharded weights and compiled artifacts to a directory.
model_neuron.save('llama-artifacts', sharded_weights=True)
```

## NeuronX Distributed Inference

```
model_path = "/home/ubuntu/models/open_llama_3b"
compiled_model_path = "/home/ubuntu/compiled_models/open_llama_3b"

neuron_config = NeuronConfig(
    batch_size=1,
    tp_degree=8,
    seq_len=128
)

config = LlamaInferenceConfig(
    neuron_config,
    load_config=load_pretrained_config(model_path)
)

model = NeuronLlamaForCausalLM(model_path, config)
```

(continues on next page)

(continued from previous page)

```
# Compile the model, shard the weights, and save to the given path.
model.compile(compiled_model_path)
```

## Models supported in Transformers NeuronX and NxD Inference model hubs

The following table depicts the list of models currently supported by TNx and their status in the NxD Inference library. For a more detailed list of models currently supported in NeuronX Distributed Inference, please refer to [NxD Inference model hub guide](#)

| Model                   | Transformers NeuronX (TNx) |                    | NxD Inference (NxDI) |                              |
|-------------------------|----------------------------|--------------------|----------------------|------------------------------|
|                         | supported in TNx           | vLLM Support (TNx) | supported in NxDI    | vLLM Support (NxD Inference) |
| BLOOM                   | Yes                        | No                 | No                   | No                           |
| GPT2                    | Yes                        | No                 | No                   | No                           |
| GPT-J                   | Yes                        | No                 | No                   | No                           |
| GPT-Neox                | Yes                        | No                 | No                   | No                           |
| Llama 2                 | Yes                        | Yes                | Yes                  | Yes                          |
| Llama 3                 | Yes                        | Yes                | Yes                  | Yes                          |
| Llama 3.1               | Yes                        | Yes                | Yes                  | Yes                          |
| Llama 3.2 (1B and 3B)   | Yes                        | Yes                | Yes                  | Yes                          |
| Llama 3.2 (11B and 90B) | No                         | No                 | Yes                  | Yes                          |
| Mistral-V2              | Yes                        | Yes                | Yes                  | Yes                          |
| Mixtral                 | Yes                        | No                 | Yes                  | Yes                          |
| DBRX                    | No                         | No                 | Yes                  | Yes                          |

## Onboarding custom or private models with NxD Inference

If you need model support for one of the models not currently supported in NxD Inference or if you have a private model that you currently implemented support in Transformers Neuronx, you need to implement the model using NxD Inference library. You can use the [Onboarding models to run on NxD Inference](#) guide.

## Neuron Config Migration

There are differences in Neuron Config parameters in Transformers NeuronX and [NxD Inference](#) libraries. If you use TNx directly without vLLM, or if you use the `override_neuron_config` param in vLLM with TNx, then you must update config parameters according to the following table.

| Transformers NeuronX parameter | NxD Inference parameter         | Notes   |
|--------------------------------|---------------------------------|---|
| <code>sparse_attn</code>       | N/A                             |   |
| <code>quant.quant_dtype</code> | <code>quantization_dtype</code> | To use quantization, set <code>quantized</code> to <code>True</code> , and provide the <code>quantized_checkpoints_path</code> where the quantized model is stored (or will be stored). |

continues on next page

Table 3.3 – continued from previous page

| Transformers NeuronX parameter                | NxD Inference parameter             | Notes   |
|---|-------------------------------------|---|
| quant.dequant_dtype                           | torch_dtype                         | NxD Inference uses the inference dtype as the dequant dtype.  |
| quant.quantize_method                         | quantization_type                   |   |
| quant.quantize_attn                           | N/A                                 |   |
| quant.no_quantize_list                        | N/A                                 |   |
| kv_cache_quant.quant_dtype                    | N/A                                 | NxD Inference uses FP8 (torch.float8_e4m3fn) for KV cache quantization. To use KV cache quantization, set kv_cache_quant to True. |
| kv_cache_quant.dequant_dtype                  | torch_dtype                         | NxD Inference uses the inference dtype as the dequant dtype.  |
| kv_cache_quant.quantize_method                | N/A                                 | NxD Inference uses direct cast.   |
| continuous_batching.max_num_seqs              | max_batch_size                      | To use continuous batching, set is_continuous_batching to True, and set tkg_batch_size to the max batch size.                     |
| continuous_batching.max_model_len             | seq_len                             |   |
| continuous_batching.optimized_paged_attention | N/A                                 |   |
| continuous_batching.block_size                | N/A                                 |   |
| continuous_batching.num_blocks                | N/A                                 |   |
| attention_layout                              | N/A                                 | NxD Inference uses BHSD layout.   |
| collectives_layout                            | N/A                                 | NxD Inference uses BHSD layout.   |
| cache_layout                                  | N/A                                 | NxD Inference uses BHSD layout.   |
| padding_side                                  | padding_side                        | NxD Inference defaults to padding on the right side.  |
| group_query_attention                         | N/A                                 |   |
| sequence_parallel_norm                        | sequence_parallel_enabled           |   |
| sequence_parallel_norm_threshold              | N/A                                 |   |
| bf16_rms_norm                                 | N/A                                 | NxD Inference upcasts RMS norm inputs to fp32.  |
| on_device_embedding                           | N/A                                 |   |
| on_device_generation                          | on_device_sampling_config           |   |
| on_device_generation.max_length               | seq_len                             | NxD Inference uses the model's sequence length.   |
| on_device_generation.do_sample                | on_device_sampling_config.do_sample |   |
| on_device_generation.top_k                    | on_device_sampling_config.top_k     | NxD Inference supports top_k through dynamic sampling. Pass the top_k values to the model inputs.                                 |
| on_device_generation.top_p                    | N/A                                 | NxD Inference supports top_p through dynamic sampling. Pass the top_p values to the model inputs.                                 |
| on_device_generation.temperature              | N/A                                 | NxD Inference supports temperature through dynamic sampling. Pass the temperature values to the model inputs.                     |

continues on next page

Table 3.3 – continued from previous page

| Transformers NeuronX parameter        | NxD Inference parameter                 | Notes  |
|---------------------------------------|---|--|
| on_device_generation.top_p_min_tokens | tokens                                  | NxD Inference defaults to a minimum of 1 token.                                      |
| on_device_generation.global_top_k     | on_device_sampling_config.global_topk   |  |
| on_device_generation.eos_token_id     | N/A                                     | NxD Inference sampling treats EOS like any other token.                              |
| on_device_generation.dynamic          | on_device_sampling_config.dynamic       |  |
| on_device_generation.deterministic    | on_device_sampling_config.deterministic |  |
| on_device_generation.per_batch_line   | N/A                                     |  |
| all_reduce_dtype                      | rpl_reduce_dtype                        | NxD Inference applies this dtype to only the all_reduce in attention's o_proj layer. |
| cast_logits_dtype                     | N/A                                     |  |
| fuse_qkv                              | fused_qkv                               |  |
| qkv_tiling                            | N/A                                     |  |
| weight_tiling                         | N/A                                     |  |
| mlp_in_weight_tiling_permute_order    | N/A                                     |  |
| mlp_out_weight_tiling_permute_order   | N/A                                     |  |
| mlp_out_weight_transpose              | N/A                                     |  |
| log_softmax_scores                    | N/A                                     |  |
| shard_over_sequence                   | flash_decoding_enabled                  |  |
| duplicate_q_weight_sos                | N/A                                     |  |
| output_all_logits                     | N/A                                     |  |
| fused_rmsnorm_qkv                     | qkv_kernel_enabled                      |  |
| fused_rmsnorm_mlp                     | mlp_kernel_enabled                      |  |
| attn_output_transposed                | N/A                                     |  |
| compilation_worker_count              | N/A                                     |  |

## LLM Inference Benchmarking guide

This guide gives an overview of the metrics that are tracked for LLM Inference and guidelines in using LLMPerf library to benchmark for LLM Inference.

### Table of contents

- *LLM Inference metrics*
- *Using LLMPerf to benchmark LLM Inference performance*
  - *Using the relevant HF tokenizer*
  - *Excluding TTFT in TPOT calculation*
  - *Benchmarking Data parallel inference with multiple model copies*

## LLM Inference metrics

Following are the essential metrics for monitoring LLM inference server performance.

| Metric                       | Description   |
|------------------------------|---|
| Time To First Token (TTFT)   | Average time taken for the LLM to process the prompt and output the first output token to the user. This is typically measured in milli seconds.  |
| Time per Output Token (TPOT) | Average time taken for LLM to generate an output token for an inference request. This is typically measured in milli seconds. This metric is also referred as Inter Token Latency (ITL) or Per Token Latency(PTL) |
| End-to-End Response Latency  | Time taken for the LLM to generate the entire response, including all output tokens. This metric is computed as end-to-end latency = (TTFT) + (TPOT) * (Number of output tokens).                                 |
| Output Token Throughput      | Number of output tokens generated per second by the inference server across all concurrent users and requests.  |

## Using LLMPerf to benchmark LLM Inference performance

LLMPerf is an open source library to benchmark LLM Inference performance. However, there are few changes that need to be applied to LLMPerf to accurately benchmark and reproduce the metrics that are published by Neuron.

All the changes outlined below are provided as a patch file that you can easily download and apply. We will work in upstreaming these changes to public LLMPerf in the future.

### Using the relevant HF tokenizer

In public LLMPerf, `hf-internal-testing` tokenizer is used by default for all the models that can impact accuracy of performance. Instead, there is a change to pass the tokenizer config of the model from Hugging Face which is being benchmarked for Neuron.

### Excluding TTFT in TPOT calculation

LLMPerf includes TTFT in Time per Output Token(or Inter Token Latency) calculation. As TPOT and TTFT are two different metrics, a change is done to LLMPerf to exclude TTFT from TPOT calculation to keep it consistent with how other industry standard performance benchmarks are done.

Following are the instructions to apply the patch to the LLMPerf library.

- Step 1: Get the Neuron git patch file  
Download the `neuron_perf.patch` file into the `llmperf` directory.
- Step 2: Apply the git patch  
Run `git apply neuron_perf.patch`. Confirm changes with `git diff`.

## Benchmarking Data parallel inference with multiple model copies

To measure performance with data parallel inference by using multiple model copies, we need to make additional changes to LLMPerf by applying the following patch:

- Step 1: Get the Neuron git patch file for data parallel inference  
Download the `llmperf_dp.patch` file into the `llmperf` directory.
- Step 2: Apply the git patch  
Run `git apply llmperf_dp.patch`. Confirm changes with `git diff`.

This patch enables data parallelism by allowing requests to be distributed across multiple model server endpoints. When multiple addresses are specified in `OPENAI_API_BASE` (e.g. “<http://server1>;<http://server2>;<http://server3>”), each request will be routed to a different server either randomly or in round-robin fashion, allowing concurrent processing across multiple model servers.

## Accuracy Evaluation of Models on Neuron Using Open Source Datasets

This guide demonstrates how to evaluate accuracy of models on Trainium and Inferentia instances using open source datasets. This approach expands on the accuracy evaluation using logits and enables you to evaluate accuracy using open source datasets like MMLU and GSM8K for tasks such as instruction following and mathematical reasoning.

Under the hood, this accuracy suite uses vLLM server to serve the model and can use benchmarking clients such as [lm-eval](#) and [LongBench](#) to evaluate on their supported datasets. In future we will add support for other benchmarking clients.

The code used in this guide is located at <https://github.com/aws-neuron/aws-neuron-samples/tree/master/inference-benchmarking/>

For a tutorial that you can follow and run on a trainium or inferentia instance please look at [Evaluating Accuracy of Llama-3.1-70B on Neuron using open source datasets](#).

## Configuration Setup

### Creating the Configuration File

Create a `test_config.yaml` file that defines your server settings and accuracy test configurations:

```
server:
  name: "test-model-server"
  model_path: "/path/to/model"
  model_s3_path: "s3://bucket/path/to/model"
  max_seq_len: 2048
  context_encoding_len: 1024
  tp_degree: 2
  n_vllm_threads: 16
  server_port: 8000
  continuous_batch_size: 2

test:
  accuracy:
    mmlu_test:
      client: "lm_eval"
```

(continues on next page)

(continued from previous page)

```

datasets: ["mmlu"]
max_concurrent_requests: 1
timeout: 3600
client_params:
    limit: 100

longbench_test:
    client: "longbench"
    datasets: ["qasper", "multifieldqa"]
    max_concurrent_requests: 1
    timeout: 7200
    client_params:
        max_length: 4096

```

## Configuration Parameters

### Server Configuration

| Parameter             | Description                      |
|-----------------------|----------------------------------|
| name                  | Identifier for your model server |
| model_path            | Local path to model files        |
| model_s3_path         | S3 location of model files       |
| max_seq_len           | Maximum sequence length          |
| context_encoding_len  | Length of context encoding       |
| tp_degree             | Tensor parallelism degree        |
| n_vllm_threads        | Number of vLLM threads           |
| server_port           | Server port number               |
| continuous_batch_size | Size of continuous batches       |

if `model_s3_path` is specified, the model will be downloaded into `model_path`, otherwise model should already exist in `model_path`.

### Accuracy Test Configuration

| Parameter               | Description  |
|-------------------------|--|
| client                  | Evaluation framework (e.g., "lm_eval", "longbench")                  |
| datasets                | List of datasets for evaluation from the supported set by the client |
| max_concurrent_requests | Maximum parallel requests  |
| timeout                 | Maximum execution time (seconds)                                     |
| client_params           | Client-specific parameters   |



## Running Evaluations

Execute accuracy tests using the CLI command:

```
python accuracy.py --config test_config.yaml
```

For more detailed information and advanced configurations, please refer to: - [lm-eval Documentation](#) - [LongBench Documentation](#)

These resources provide comprehensive guides on client-specific parameters and advanced evaluation scenarios.

## Custom Quantization

### Overview

This document gives an overview of customizable quantization feature in the NxD Inference. Users can specify which modules should not be converted during quantization, allowing custom quantized model inference. Users can take an un-quantized model and apply selective quantization to specific layers while keeping others in full precision.

The document also explains how to use external libraries like [llmcompressor](#), including quantization config setup and applying necessary patches. It also covers running inference with quantized models and specifying unconverted modules through either command-line arguments or NeuronConfig kwargs.

### Quantization

Custom quantization allows users to have fine-grained control over which layers of the model are quantized. This can be particularly useful for maintaining model accuracy while still benefiting from the reduced memory footprint of quantization. For more detailed information on quantization techniques and implementation, please refer to the [quantization feature guide](#).

### Quantize Using NxD

Quantization can significantly reduce the model size and inference time, making it more suitable for deployment of large models that typically cannot fit on a single instance. However, not all layers of the model benefit equally from quantization.

- Some layers, especially those involved in critical computations like normalizations or certain types of activations, may see a significant drop in accuracy if quantized. Leaving these layers in full precision helps maintain the overall performance of the model.
- Quantization can also introduce small errors in each layer's computation. When these errors accumulate through the network, they can lead to a noticeable degradation in performance. Keeping certain layers in full precision can mitigate this accumulation.

To leverage the customizable quantization feature in NxD, follow the steps below. This process involves importing necessary libraries, defining the model and output paths, specifying modules to not convert, and utilizing a quantization function to create a quantized model.

```
import torch
from typing import Optional, List
from transformers import AutoModelForCausalLM, AutoTokenizer
from neuronx_distributed_inference.modules.checkpoint import prune_state_dict, save_state_
    dict_safetensors
```

(continues on next page)

(continued from previous page)

```

from neuronx_distributed.quantization.quantization_utils import quantize_pytorch_model_
↳per_channel_symmetric, convert_qint8_to_int8_state_dict

model_path = "<model_path/llama-3.1-405b-instruct-4layers/>"
output_path = "<save_quantized_checkpoints>"

modules_to_not_convert = [
    "lm_head",
    "layers.0.self_attn",
    "layers.1.self_attn",
    "layers.2.self_attn",
    "layers.1.mlp"
]

def quantize(model: torch.nn.Module, dtype=torch.qint8, modules_to_not_convert:
↳Optional[List[str]] = None) -> torch.nn.Module:
    quant_model = quantize_pytorch_model_per_channel_symmetric(model, dtype=dtype,
↳modules_to_not_convert=modules_to_not_convert)
    model_quant_sd = quant_model.state_dict()
    convert_qint8_to_int8_state_dict(model_quant_sd)
    quantized_state_dict = prune_state_dict(model_quant_sd)
    return quantized_state_dict

model = AutoModelForCausalLM.from_pretrained(model_path)
tokenizer = AutoTokenizer.from_pretrained(model_path)

state_dict = quantize(model, torch.float8_e4m3fn, modules_to_not_convert)

save_state_dict_safetensors(state_dict=state_dict, state_dict_dir=output_path)
tokenizer.save_pretrained(output_path)

```

## Quantize using external libraries

In addition to the built-in quantization features of NxD, users can also leverage external libraries for more flexible and advanced quantization options. One such library is `llmcompressor`, which offers a robust set of tools for quantizing models. To use the `llmcompressor` library for quantization, follow the steps below.

This process involves importing necessary libraries, specifying modules to not convert, setting up a quantization recipe, and applying the quantization to create a quantized model. `llmcompressor` gives us a range from `-/+448`, so it is important to ensure the scale range is set from `-/+240` if you need to run inference on the quantized model later using NxD Inference. Values outside the range of `-/+240` on Neuron devices result in NaNs.

The LLaMA model is an example where not all layers are quantized.

- By keeping the attention layers, first and last MLP layers, and the LM head in full precision, the model maintains high accuracy in tasks like language generation and comprehension.
- Quantizing the remaining layers (e.g., intermediate MLP layers) reduces the model size and inference time without significantly compromising performance.
- This strategy allows for a balanced trade-off between model efficiency and accuracy, making the model suitable for high performance deployment.

```

import torch
from llmcompressor.transformers import oneshot, SparseAutoModelForCausalLM
from transformers import AutoTokenizer
from compressed_tensors.quantization.utils.helpers import calculate_range
from compressed_tensors.quantization.quant_args import QuantizationType
import compressed_tensors.quantization.utils.helpers as helpers

model_path = "<model_path>/llama-3.1-405b-instruct-4layers/"
output_path = "<save_quantized_checkpoints>"

modules_to_not_convert = ['lm_head',
    "model.layers.0.mlp.down_proj",
    "model.layers.0.mlp.gate_proj",
    "model.layers.0.mlp.up_proj",
    "model.layers.3.mlp.down_proj",
    "model.layers.3.mlp.gate_proj",
    "model.layers.3.mlp.up_proj",
    "model.layers.0.self_attn.k_proj",
    "model.layers.0.self_attn.o_proj",
    "model.layers.0.self_attn.q_proj",
    "model.layers.0.self_attn.v_proj",
    "model.layers.1.self_attn.k_proj",
    "model.layers.1.self_attn.o_proj",
    "model.layers.1.self_attn.q_proj",
    "model.layers.1.self_attn.v_proj",
    "model.layers.2.self_attn.k_proj",
    "model.layers.2.self_attn.o_proj",
    "model.layers.2.self_attn.q_proj",
    "model.layers.2.self_attn.v_proj",
    "model.layers.3.self_attn.k_proj",
    "model.layers.3.self_attn.o_proj",
    "model.layers.3.self_attn.q_proj",
    "model.layers.3.self_attn.v_proj"]

recipe = f"""
quant_stage:
  quant_modifiers:
    QuantizationModifier:
      ignore: {modules_to_not_convert}
  config_groups:
    group_0:
      weights:
        num_bits: 8
        type: float
        strategy: channel
        dynamic: false
        symmetric: true
      input_activations:
        num_bits: 8
        type: float
        strategy: token
        dynamic: true
        symmetric: true

```

(continues on next page)

(continued from previous page)

```

        targets: ["Linear"]
"""

model = SparseAutoModelForCausalLM.from_pretrained(
    model_path, torch_dtype="auto"
)

# Monkey patch to rescale weights from -/+448 to -/+240
original_calculate_range = helpers.calculate_range
def calculate_range(*args, **kwargs):
    q_min, q_max = original_calculate_range(*args, **kwargs)
    if args[0].type == QuantizationType.FLOAT and args[0].num_bits == 8:
        return torch.tensor(-240.0, device=args[1]), torch.tensor(240.0, device=args[1])
    return q_min, q_max

# Patch it
helpers.calculate_range = calculate_range
oneshot(model=model, recipe=recipe)

for name, module in model.named_modules():
    if hasattr(module, 'weight_scale'):
        module.weight_scale.data = module.weight_scale.data.to(torch.float32)

tokenizer = AutoTokenizer.from_pretrained(model_path)

model.save_pretrained(output_path)
tokenizer.save_pretrained(output_path)

```

## Quantization Commands

To utilize the quantization commands in NxD Inference, users can follow the instructions below. These commands cover the required flags to enable running inference with quantized models.

### First Quantize then Inference

If you have a model in full precision and need to quantize it on the CPU first before using it for inference, you can set the following flags to enable quantization during inference:

```

inference_demo --model-type llama --task-type causal-lm run \
--model-path /your_model_path/ \
--compiled-model-path /save_to_path/ \
--torch-dtype bfloat16 \
--tp-degree 32 \
--batch-size 1 \
--max-context-length 1024 \
--quantized \
--quantization-dtype f8e4m3 \
--quantization-type per_channel_symmetric \
--quantized-checkpoints-path /save_to_path/ \
--seq-len 2048 \

```

(continues on next page)

(continued from previous page)

```
--fused-qkv \
--pad-token-id 2 \
--on-device-sampling \
--sequence-parallel-enabled \
--attn-kernel-enabled \
--prompt "I believe the meaning of life is" \
--is-continuous-batching \
--enable-fused-speculation \
--enable-eagle-speculation \
--speculation-length 4 \
--draft-model-path /your_draft_model_path \
--modules-to-not-convert-file /path/modules_to_not_convert.json
```

### Inference Using Already quantized checkpoint

To utilize the quantization commands in NxD, users can follow the instructions below. These commands cover the required flags to enable running inference with quantized models. The `modules-to-not-convert-file` allows you to specify the list of modules to not quantize, useful for quantizing models that explicitly require having some modules left in their original precision.

### How to Use

- Pass `modules_to_not_convert` using Inference Demo

```
inference_demo --model-type llama --task-type causal-lm run \
  --model-path <path> \
  --compiled-model-path <path> \
  --torch-dtype bfloat16 \
  --tp-degree <value> \
  --batch-size <value> \
  --max-context-length <value> \
  --seq-len <value> \
  --on-device-sampling \
  --mlp-kernel-enabled \
  --quantized-mlp-kernel-enabled \
  --quantization-dtype <dtype> \
  --quantization-type <type> \
  --prompt "I believe the meaning of life is" \
  --modules-to-not-convert-file /<your_path>/modules_to_not_convert.json
```

- Pass `modules_to_not_convert` using NeuronConfig kwargs

```
neuron_config = NeuronConfig(
    tp_degree=32,
    batch_size=2,
    max_context_length=32,
    seq_len=64,
    on_device_sampling_config=OnDeviceSamplingConfig(top_k=1),
    enable_bucketing=True,
    flash_decoding_enabled=False,
```

(continues on next page)

(continued from previous page)

```
modules_to_not_convert=["lm_head", "layers.0.self_attn", "layers.1.mlp", ...],
draft_model_modules_to_not_convert=["lm_head", "layers.0.self_attn", "layers.1.mlp",
→ ..., "fc"]
)
```

*Note: If you are creating different NeuronConfig for draft and target models, you only need to pass the modules\_to\_not\_convert list for both.*

## JSON File Structure

The JSON structure is a crucial component for specifying which modules should not be converted during the quantization if you are using inference demo. This section provides detailed examples of how to format the JSON file. The JSON structure depends on whether fused speculation is used.

### 1. Basic Structure

For simple cases:

```
{
  "modules_to_not_convert": [
    "lm_head",
    "layers.0.self_attn",
    "layers.1.self_attn",
    "layers.2.self_attn",
    "layers.3.self_attn",
    "layers.0.mlp",
    "layers.3.mlp"
  ]
}
```

OR

```
{
  "model": {
    "modules_to_not_convert": [
      "lm_head",
      "layers.0.self_attn",
      "layers.1.self_attn",
      "layers.2.self_attn",
      "layers.3.self_attn",
      "layers.0.mlp",
      "layers.3.mlp"
    ]
  }
}
```

### 1. With Fused Speculation

```
{
  "model": {
    "modules_to_not_convert": [
      "lm_head",
      "layers.0.self_attn",

```

(continues on next page)

(continued from previous page)

```

        "layers.1.self_attn",
        "layers.2.self_attn",
        "layers.3.self_attn",
        "layers.0.mlp",
        "layers.3.mlp"
    ]
},
"draft_model": {
    "modules_to_not_convert": [
        "lm_head",
        "layers.0.self_attn",
        "layers.0.mlp",
        "fc"
    ]
}
}

```

## Important Notes

- Make sure to assign partial names in modules to avoid conversion, as shown in the examples above. This is necessary due to different naming schemes between the model layers being read from the source and the model we create for inference. The above examples include the partial parts of the names which are common between the two naming schemes.
  - For example: Original model names are like `model.layers.0.self_attn.q_proj`, whereas the names we give are like `layers.0.self_attn.qkv_proj.q_proj`
- Quantization with Fused Speculation
  - We currently do not quantize the draft model, Include these in the `draft_model.modules_to_not_convert` section of your JSON file

## Backward Incompatible Changes:

- Now running the quantization workflow will need the `modules-to-not-convert-file` flag while running with `inference demo` because we no longer hard-code the layers to incorporate quantized layers.

## NxD Inference Weights Sharding Guide

NxD Inference provides two approaches to shard model weights and load them onto Neuron Devices, enabling parallel processing (e.g. Tensor Parallelism) on each device. This guide demonstrates the usage of both approaches using *Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances*, and provides insights into selecting the appropriate method based on the usage pattern and performance requirements.

**Note:** Sharding speed on different storage volumes can vary. We recommend to use NVMe solid state drive (SSD) storage to achieve the best sharding performance. This guide shows sharding results on NVMe SSD. For more information about NVMe storage on EC2 instances, see the following: \* [Instance store volumes](#) in the Amazon EC2 User Guide. Instance store volumes are drives attached to EC2 instances that you can use for temporary storage. Neuron instances such as Trn1 and Trn2 include NVMe drives that you can use as instance store volumes. \* [EBS volumes and](#)

NVMe in the Amazon EBS User Guide. For persistent storage on NVMe, you can use EBS volumes built on AWS Nitro.

### Table of contents

- *Shard on compile (Pre-shard)*
- *Shard on load*

## Shard on compile (Pre-shard)

The shard on compile (pre-shard) approach loads the supported pretrained *checkpoints*, converts to Neuron compatible format, shards for each parallel rank and serializes sharded weights to disk as safetensors files. The entire sharding and serialization process can take a few minutes to hours depending on the model size and throughput of the storage volume. This approach is optimized to minimize the future model loading time.

The following example demonstrates how to run shard on compile with Llama3.1-405b.

First, complete the *prerequisites* for running Llama3.1-405b on a Trn2.48xlarge instance.

Next, enable shard on compile by adding `--save-sharded-checkpoint` to the command. The sharded checkpoints will be saved to the `/weights` folder under the specified `COMPILED_MODEL_PATH`.

Full command to run shard on compile for Llama3.1-405b:

```
# Replace this with the path where model files are downloaded.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct/"
# This is where the compiled model will be saved.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-405B-Instruct/"

NUM_CORES=128
TP_DEGREE=64
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
    --model-path $MODEL_PATH \
    --compiled-model-path $COMPILED_MODEL_PATH \
    --torch-dtype bfloat16 \
    --start_rank_id 0 \
    --local_ranks_size $TP_DEGREE \
    --tp-degree $TP_DEGREE \
    --batch-size 1 \
    --max-context-length 2048 \
    --seq-len 2048 \
    --on-device-sampling \
    --top-k 1 \
```

(continues on next page)



(continued from previous page)

```
--fused-qkv \
--sequence-parallel-enabled \
--qkv-kernel-enabled \
--attn-kernel-enabled \
--mlp-kernel-enabled \
--cc-pipeline-tiling-factor 1 \
--pad-token-id 2 \
--save-sharded-checkpoint \
--prompt "What is annapurna labs?" 2>&1 | tee log
```

You should see the outputs below in your logs. The duration can slightly vary between runs. Note that model loading started only after sharding is completed.

```
INFO:Neuron:Sharding Weights for ranks: 0...63
INFO:Neuron:Done sharding weights in 1856.5586961259833 seconds
Loading model to Neuron...
Total model loading time: 107.76132441597292 seconds
```

Now that sharded checkpoints have been serialized to disk, you may save sharding time in your next run by adding `--skip-sharding` to the command. Sharded weights will be directly loaded from the disk for inference, which saves you 30+ minutes of sharding for each subsequent run in this example.

The total model loading time in each subsequent run is expected to be comparable with the first run.

### Shard on load

The shard on load approach significantly reduces sharding overheads by parallelizing tensor movement in sharding/loading and skipping sharded checkpoints serialization. This approach is preferred when you are working with weights that are frequently retrained/fine-tuned so re-sharding becomes a bottleneck when serving with new weights. Since Neuron 2.23 release, Shard on load is enabled by default in NxD Inference.

Full command to run shard on load for Llama3.1-405b is shown below. Note that `--save-sharded-checkpoint` is excluded from the command.

```
# Replace this with the path where model files are downloaded.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct/"
# This is where the compiled model will be saved.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-405B-Instruct/"

NUM_CORES=128
TP_DEGREE=64
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
  --model-path $MODEL_PATH \
```

(continues on next page)

(continued from previous page)

```

--compiled-model-path $COMPILED_MODEL_PATH \
--torch-dtype bfloat16 \
--start_rank_id 0 \
--local_ranks_size $TP_DEGREE \
--tp-degree $TP_DEGREE \
--batch-size 1 \
--max-context-length 2048 \
--seq-len 2048 \
--on-device-sampling \
--top-k 1 \
--fused-qkv \
--sequence-parallel-enabled \
--qkv-kernel-enabled \
--attn-kernel-enabled \
--mlp-kernel-enabled \
--cc-pipeline-tiling-factor 1 \
--pad-token-id 2 \
--prompt "What is annapurna labs?" 2>&1 | tee log

```

You should see the outputs below in your logs. The duration can slightly vary between runs. Note that sharding happened while model was being loaded (i.e. shard on load).

```

Loading model to Neuron...
INFO:Neuron:Done Sharding weights in 49.31190276599955 seconds
Total model loading time: 187.3972628650372 seconds

```

As you can see, weights sharding of shard on load is much faster than that of shard on compile.

When the current run finishes, no sharded checkpoints will be saved. Therefore, you cannot use `--skip-sharding` for your next run. In each subsequent run, NxN Inference will do the exact same amount of sharding work, so the total model loading time is expected to be comparable with the first run. It's also expected that the total model loading time is longer than that of shard on compile, due to the extra sharding work it has to do during loading time.

## Disaggregated Inference [BETA]

### Overview

Disaggregated Inference (DI), also known as disaggregated serving, disaggregated prefill, P/D disaggregation, is an LLM serving architecture that separates the prefill and decode phases of inference onto different hardware resources. To achieve this, prefill worker needs to transfer the computed KV cache to the decode worker to resume decoding. Separating the compute intensive prefill phase from the memory bandwidth intensive decode phase can improve the LLM serving experience by

1. Removing prefill interruptions to decode from continuous batching to reduce inter token latency (ITL). These gains can be used to achieve higher throughput by running with a higher decode batch size while staying under Service Level Objectives (SLO).
2. Adapt to changing traffic patterns while still remaining under application SLOs.
3. Enable independent scaling of resources and parallelism strategies for prefill (compute bound) and decode (memory bound).

**Note:** This feature is still in beta. Currently only a single decode server and a single prefill server is supported (1P1D).

Automatic Prefix Caching is not supported with DI.

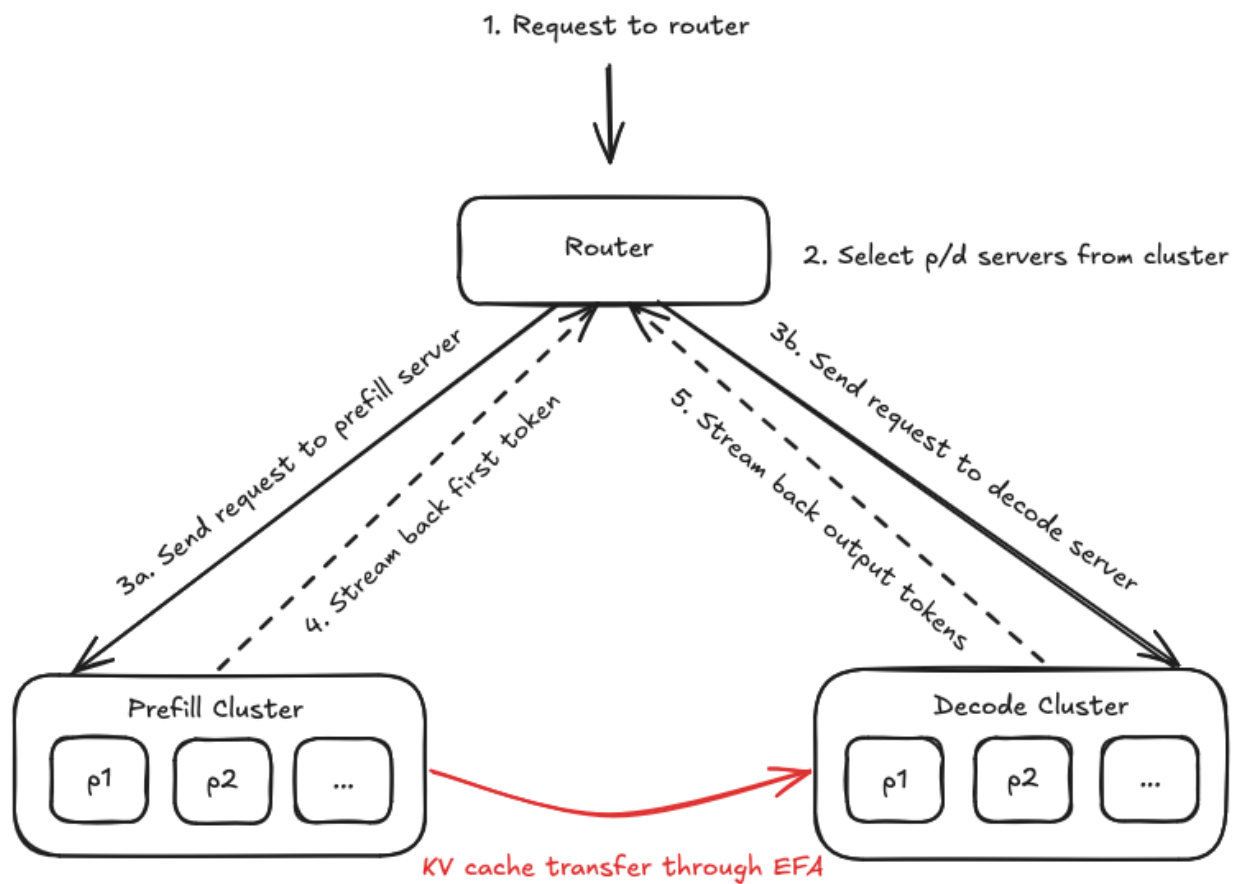
## Neuron Implementation Details

Disaggregated Inference is mainly implemented through Neuron's vLLM fork <https://github.com/aws-neuron/upstreaming-to-vllm/tree/neuron-2.24-vllm-v0.7.2> and the Neuron Runtime.

There are three main components to a DI workflow.

1. The router. Its job is to orchestrate requests to servers inside the prefill and decode clusters.
2. The prefill cluster. This represents all of the prefill servers ready to run a DI workload.
3. The decode cluster. This represents all of the decode servers ready to run a DI workload.

Below is an example lifespan of a single request through the DI flow.



1. A request is sent to the router (1), a component responsible for orchestrating (2) the requests to both the prefill and decode servers. It receives responses from the prefill and decode servers and streams the results back to the user.

2. The prefill server receives the request from the router (3a) and starts prefilling. After the prefill completes (4), it updates the status of the request for the decode server by sending information through another ZMQ server. Then, it listens for a “pull request” from the decode server to initiate the KV cache transfer. We use Neuron runtime APIs to transfer the KV cache through EFA from Neuron device to Neuron device. This is a zero copy transfer, meaning that we do not copy the KV cache from a Neuron device to CPU to transfer, but rather directly transfer KV cache from Neuron device to Neuron device. The transfer is also asynchronous. This means that the prefill server can immediately start

prefilling the next request while the KV cache of the previous request is being transferred. This ensures that TTFT is not impacted for other requests while the KV cache for older request is being transferred to decode.

3. The decode server also receives a request from the router at the same time as the prefill server (3b). It waits until it receives a signal that its corresponding prefill is done from the prefill server by listening on the ZMQ server. Then, if there is a free spot in the decode batch, the scheduler will schedule the request and send a “pull request” to the prefill server. This initiates the asynchronous KV cache transfer (red arrow) through EFA by calling the Neuron Runtime API. The receive also needs to be asynchronous to ensure smooth ITL. While the receive is happening other decode requests will still run. As soon as the receive is finished the scheduler will add the request to the next decode batch (5).

## Prefill Decode Interference When Colocating Prefill and Decode

In traditional continuous batching, prefill requests are prioritized over decode requests. Prefills are run as batch size 1 because they are compute intensive whereas decodes can be run at a higher batch size because it is constrained on memory bandwidth not compute. To ensure the highest throughput, continuous batching schedulers prioritize new prefills if the decode batch is not at max capacity. As soon as a decode request finishes, another prefill is scheduled to fill the finished request’s place. However, all other ongoing decodes pause while the new prefill is running because that prefill uses the same compute resources. This effect is known as prefill stall or prefill contention.

Disaggregated Inference avoids prefill stall because the decode workflow is never interrupted by a prefill as it receives KV caches asynchronously while decoding. The overall ITL on DI is affected by the transfer time of the KV cache but this does not scale with batch size. For example, in a continuous batching workload of batch size 8 each request will on average be interrupted 7 times whereas in DI each request is only affected by a single transfer since it happens asynchronously.

Another advantage of DI is its ability to adapt to traffic patterns while maintaining a consistent ITL. For example, if prefill requests double in length the application can double the amount of available prefill servers in the prefill cluster to match the new traffic pattern. Continuous batching workloads would suffer because longer prefill requests increase tail ITL whereas DI workloads continue to deliver a low variance and a predictable customer experience.

Additionally, DI also allows users to tailor their parallelism strategies differently for prefill and decode. For example, a model with 32 attention heads may prefer to run two decode servers Data Parallel=2 (DP). each with Tensor Parallel=32 (TP) in order to reduce KV replication instead of TP=64. Such replication will get worse if using Group Query Attention (GQA).

DI does not necessarily improve throughput directly but it can help depending on the workload. Continuous batching is a technique optimized for throughput at the cost of ITL. An application may have an SLO to ensure that ITL is under a certain threshold. Because increasing the batch size increases the amount of prefill stall, and therefore increases ITL, many applications run on smaller than ideal batch sizes when using continuous batching. DI can allow an application to run at a higher batch size while still keeping ITL under the application defined SLO.

## Trade-Offs

Because DI runs prefill and decode separately, each part of the inference process needs to operate at an equal level of efficiency to maximize throughput and hardware resources. For example, if you can process 4 prefill requests per second and two decode requests per second the application will be stuck processing two requests per second. It is also important to note that the prefill and decode efficiency can vary based on the prompt length and the number of tokens for a response respectively. Continuous batching and chunked prefill do not have this issue as these techniques run prefill and decode on the same hardware.

One technique to remediate this is to run with a dynamic amount of prefill and decode servers. We call this dynamic xPyD. In the above example, we could run with 1 prefill and 2 decode servers so that our prefill and decode efficiency will be balanced. This is being actively worked on and currently only static configurations of one prefill to one decode (1PID) are supported.

## Example Usage

Refer to the [offline inference DI example](#) for a quick example to get started.

Refer to the [Disaggregated Inference Tutorial](#) for a detailed usage guide.

Use the NxD Inference (**neuronx-distributed-inference**) Developer Guides to learn how to use NxD Inference.

- [NxD Inference Features Configuration Guide](#)
- [NxD Inference - Production Ready Models](#)
- [Onboarding models to run on NxD Inference](#)
- [vLLM User Guide for NxD Inference](#)
- [Testing modeling code with NxD Inference](#)
- [Migrating from NxD Core inference examples to NxD Inference](#)
- [Migrating from Transformers NeuronX to NeuronX Distributed\(NxD\) Inference](#)
- [LLM Inference Benchmarking guide](#)
- [Accuracy Evaluation of Models on Neuron Using Open Source Datasets](#)
- [Custom Quantization](#)
- [NxD Inference Weights Sharding Guide](#)
- [Disaggregated Inference \[BETA\]](#)

## 3.2.5 Tutorials

### Tutorial: Deploying Llama3.1 405B (Trn2)

NeuronX Distributed (NxD) Inference enables you to deploy Llama3.1 405B on a single Trn2 instance.

You can run Llama3.1 405B with default configuration options. NxD Inference also provides several features and configuration options that you can use to optimize and tune the performance of Llama3.1 405B on Trn2. This guide walks through how to run Llama3.1 405B on Trn2 with vLLM, and how to enable these optimizations for optimal performance. In addition, we also have a separate tutorial for running Llama3.1 405B with vanilla fused speculative decoding [Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances](#).

#### Table of contents

- [Background, Concepts, and Optimizations](#)
  - [Logical NeuronCore Configuration \(LNC\)](#)
  - [Tensor parallelism \(TP\) on Trn2](#)
  - [Optimizing Performance](#)
- [Tutorial: Run Llama3.1 405B on Trn2](#)
  - [Step 1: Connect to the Trn2 instance](#)
  - [Step 2: Install the vLLM version that supports NxD Inference](#)
  - [Step 3: Deploy Llama 3.1 405B sample code](#)

## Background, Concepts, and Optimizations

### Logical NeuronCore Configuration (LNC)

On Trn2, the Neuron SDK supports *Logical NeuronCore Configuration (LNC)*, which determines the number of NeuronCores visible to the Neuron SDK. When running on Trn2, the Neuron SDK is optimized for LNC=2, which means each NeuronCore visible to the Neuron SDK is two physical NeuronCores. The LNC configuration also affects what TP degree options you can use.

NxD Inference automatically chooses the correct LNC configuration based on the target platform.

For more information about LNC, see [Logical NeuronCore configuration](#).

### Tensor parallelism (TP) on Trn2

Each Trn2 instance has 128 Neuron cores. With LNC=2, you can set a TP degree up to 64. We recommend that you use LNC=2 for all models on Trn2.

For more information about tensor parallelism in NxD Inference, see [nxd-tensor-parallelism](#).

## Optimizing Performance

### EAGLE Speculative Decoding

Speculative decoding is a performance optimization technique where a smaller *draft* LLM model predicts the next tokens, and the larger *target* LLM model verifies those predictions.

NxD Inference supports EAGLE v1 speculative decoding with a flat draft structure. To use EAGLE v1, you must use an EAGLE checkpoint for a draft model that is not tree-based and is specifically fine-tuned for EAGLE speculation. For more information about EAGLE, see the official implementation on GitHub: [SafeAILab/EAGLE](#).

To optimize performance for EAGLE speculative decoding, NxD Inference uses a feature called *fused speculation*, where the draft model and target model are fused into a single compiled model artifact to improve performance. Fused speculation uses a different config called `FusedSpecNeuronConfig`, which specifies the model class, draft config, and draft model path to fuse with the target model.

For more information about speculative decoding in NxD Inference, including other types of speculative decoding supported, see [Speculative Decoding](#).

### FP8 Quantization

NxD Inference supports FP8 quantization, where model weights and data are converted to a smaller data type to reduce memory bandwidth usage. FP8 quantization enables optimal usage of memory bandwidth to improve model performance. For more information, see [Model Weight Quantization](#).

NxD Inference also supports KV cache quantization, where the KV cache is quantized to FP8. For more information, see [KV Cache Quantization](#).

## Optimized Kernels

NxD Inference supports kernels that optimize parts of the modeling code for best performance.

- **Flash attention.** This kernel uses a sharded flash attention implementation to improve performance during the context encoding pass. This kernel is enabled automatically at supported sequence lengths. For LNC2, NxD Inference automatically enables flash attention for sequence lengths of 256 and larger that are divisible by 256. For LNC1, NxD Inference automatically enables flash attention for sequence lengths of 4096 and larger. You can also enable it with `attn_kernel_enabled=True` in `NeuronConfig`. NxD Inference automatically enables the flash attention kernel at supported sequence lengths even if `attn_kernel_enabled` is `false`.
- **QKV.** This kernel fuses the QKV layers to improve performance during the attention forward pass. To enable this kernel, set `qkv_kernel_enabled=True` in `NeuronConfig`.
- **MLP.** This kernel implements the MLP module used in decoder layers. To enable this kernel, set `mlp_kernel_enabled=True` in `NeuronConfig`.
- **Quantized MLP.** This kernel implements a quantized version of the MLP kernel. This kernel uses FP8 compute to improve performance. To enable this kernel, set `quantized_mlp_kernel_enabled=True`. This kernel requires `mlp_kernel_enabled=True`.

---

**Note:** To use the QKV and MLP kernels, you must set `torch_dtype` to `torch.bfloat16` in `NeuronConfig`.

---

## Tutorial: Run Llama3.1 405B on Trn2

As a prerequisite, this tutorial requires that you have a Trn2 instance created from a Deep Learning AMI that has the Neuron SDK pre-installed.

To set up a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK, see [NxD Inference Setup Guide](#).

### Step 1: Connect to the Trn2 instance

Use SSH to connect to the Trn2 instance using the key pair that you chose when you launched the instance.

After you are connected, activate the Python virtual environment that includes the Neuron SDK.

```
source ~/aws_neuronx_venv_pytorch_2_5_nxd_inference/bin/activate
```

Run `pip list` to verify that the Neuron SDK is installed.

```
python -m pip list
```

You should see Neuron packages including `neuronx-distributed-inference` and `neuronx-cc`.

## Step 2: Install the vLLM version that supports NxD Inference

NxD Inference supports running models with vLLM. This functionality is available in the AWS Neuron fork of the vLLM GitHub repository. Install the latest release branch of vLLM from the AWS Neuron fork following instructions in the [vLLM User Guide for NxD Inference](#).

## Step 3: Deploy Llama 3.1 405B sample code

Choose one of the following examples to run on the Trn2 instance:

1. Deploy Llama3.1 405B with vLLM offline inference. This example demonstrates how to deploy on Trn2 with vLLM and topK sampling.
2. Deploy Llama3.1 405B with EAGLE speculative decoding. This example demonstrates how to use EAGLE to optimize Llama3.1 405B on Trn2.

### Example 1: Deploy Llama3.1 405B on Trn2 with vLLM offline inference

This example demonstrates how to deploy Llama3.1 405B on Trn2 with vLLM offline inference and the following configuration options:

- Sequence length: 2048 tokens
- Max context length: 1024 tokens
- Speculation length: 6 tokens
- Flash attention, QKV, and MLP kernels
- On-device sampling with topK sampling

To use this sample, you must first download a 405B model checkpoint from Hugging Face to a local path on the Trn2 instance. For more information, see [Downloading models](#) in the Hugging Face documentation. You can download and use [meta-llama/Llama-3.1-405B-Instruct](#) for this tutorial.

```
import os
import torch

from vllm import LLM, SamplingParams

# Force vLLM framework to use neuronx-distributed-inference
os.environ['VLLM_NEURON_FRAMEWORK'] = "neuronx-distributed-inference"

model_path = "/home/ubuntu/models/Llama-3.1-405B-Instruct/"

def run_llama_generate():
    # Initialize vLLM.
    llm = LLM(
        model=model_path,
        tensor_parallel_size=64,
        max_num_seqs=1,
        max_model_len=2048,
        block_size=2048,
        dtype=torch.bfloat16,
```

(continues on next page)



(continued from previous page)

```

# Configure NeuronConfig.
override_neuron_config={
    "max_context_length": 1024,
    "skip_warmup": True,
},
device="neuron"
)

# Run vLLM to generate outputs.
prompts = ["I believe the meaning of life is"]
sampling_params = SamplingParams(top_k=50)
outputs = llm.generate(prompts, sampling_params)
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

if __name__ == "__main__":
    run_llama_generate()

```

### Example 2: Deploy Llama3.1 405B on Trn2 with EAGLE speculative decoding

This example demonstrates how to deploy Llama3.1 405B on Trn2 with EAGLE speculative decoding.

**Note:** To use this example, you must provide an EAGLE-trained Llama3.1 405B checkpoint to use for EAGLE speculative decoding. For more information about EAGLE checkpoint compatibility with NxD Inference, see [EAGLE Speculative Decoding](#).

This example uses the following configuration options:

- Sequence length: 2048 tokens
- Max context length: 1024 tokens
- Speculation length: 6 tokens
- Flash attention, QKV, and MLP kernels
- On-device sampling with greedy sampling
- Sequence parallelism enabled
- Auto-bucketing enabled, which automatically selects buckets to use. For more information about bucketing and how to customize the buckets used, see [Bucketing](#).

```

import copy
import os
import torch

from transformers import AutoTokenizer, GenerationConfig

from neuronx_distributed_inference.models.config import FusedSpecNeuronConfig,

```

(continues on next page)

(continued from previous page)

```

↳NeuronConfig, OnDeviceSamplingConfig
from neuronx_distributed_inference.models.llama.modeling_llama import _
↳LlamaInferenceConfig, NeuronLlamaForCausalLM
from neuronx_distributed_inference.utils.hf_adapter import HuggingFaceGenerationAdapter, _
↳load_pretrained_config

model_path = "/home/ubuntu/models/llama-3.1-405b-Instruct/"
draft_model_path = "/home/ubuntu/models/EAGLE-llama-3-405b/"
compiled_model_path = "/home/ubuntu/neuron_models/llama-3-405b-instruct-EAGLE/"

# Set environment variables for Trn2.
os.environ["XLA_DENSE_GATHER_FACTOR"] = "0"
os.environ["NEURON_RT_EXEC_TIMEOUT"] = "600"

def run_llama_generate():
    top_k = 1
    do_sample = False

    # Initialize tokenizer.
    tokenizer = AutoTokenizer.from_pretrained(model_path, padding_side="right")
    tokenizer.pad_token = tokenizer.eos_token

    # Initialize target model config.
    neuron_config = NeuronConfig(
        torch_dtype=torch.bfloat16,
        tp_degree=64,
        batch_size=1,
        max_context_length=1024,
        seq_len=2048,
        on_device_sampling_config=OnDeviceSamplingConfig(
            dynamic=False,
            do_sample=do_sample,
            top_k=top_k
        ),
        enable_eagle_speculation=True,
        enable_fused_speculation=True,
        speculation_length=6,
        trace_tokengen_model=False,
        enable_bucketing=True,
        fused_qkv=True,
        sequence_parallel_enabled=True,
        attn_kernel_enabled=True,
        qkv_kernel_enabled=True,
        mlp_kernel_enabled=True,
        cc_pipeline_tiling_factor=1,
    )
    config = LlamaInferenceConfig(
        neuron_config,
        load_config=load_pretrained_config(model_path),
    )

    # Initialize draft model config.

```

(continues on next page)

(continued from previous page)

```

draft_neuron_config = copy.deepcopy(neuron_config)
draft_neuron_config.trace_tokengen_model = True
draft_neuron_config.enable_fused_speculation = False
draft_neuron_config.is_eagle_draft = True
draft_neuron_config.sequence_parallel_enabled = False
draft_config = LlamaInferenceConfig(
    draft_neuron_config,
    load_config=load_pretrained_config(draft_model_path)
)

# Initialize fused speculation config.
fused_spec_config = FusedSpecNeuronConfig(
    NeuronLlamaForCausalLM._model_cls,
    draft_config=draft_config,
    draft_model_path=draft_model_path,
)
config.fused_spec_config = fused_spec_config

# Compile and save model.
print("\nCompiling and saving model...")
model = NeuronLlamaForCausalLM(model_path, config)
model.compile(compiled_model_path)
tokenizer.save_pretrained(compiled_model_path)

# Load from compiled checkpoint.
print("\nLoading model from compiled checkpoint...")
model = NeuronLlamaForCausalLM(compiled_model_path)
model.load(compiled_model_path)
tokenizer = AutoTokenizer.from_pretrained(compiled_model_path)

# Initialize generation config.
generation_config = GenerationConfig.from_pretrained(model_path)
generation_config_kwargs = {
    "do_sample": do_sample,
    "top_k": top_k,
    "pad_token_id": 0,
    "prompt_lookup_num_tokens": neuron_config.speculation_length,
}
generation_config.update(**generation_config_kwargs)

# Generate outputs.
print("\nGenerating outputs...")
prompts = ["I believe the meaning of life is"]
print(f"Prompts: {prompts}")
inputs = tokenizer(prompts, padding=True, return_tensors="pt")
generation_model = HuggingFaceGenerationAdapter(model)
outputs = generation_model.generate(
    inputs.input_ids,
    generation_config=generation_config,
    attention_mask=inputs.attention_mask,
    max_length=model.config.neuron_config.max_length,
)

```

(continues on next page)

(continued from previous page)

```

output_tokens = tokenizer.batch_decode(outputs, skip_special_tokens=True, clean_up_
↪tokenization_spaces=False)
print("Generated outputs:")
for i, output_token in enumerate(output_tokens):
    print(f"Output {i}: {output_token}")

if __name__ == "__main__":
    run_llama_generate()

```

## Tutorial: Deploying Llama3.2 Multimodal Models

NeuronX Distributed Inference (NxDI) enables you to deploy Llama-3.2-11B-Vision-Instruct and Llama-3.2-90B-Vision-Instruct models on Neuron Trainium and Inferentia instances.

You can run Llama3.2 Multimodal with default configuration options. NxD Inference also provides several features and configuration options that you can use to optimize and tune the performance for inference. This guide walks through how to run Llama3.2 Multimodal with vLLM, and how to enable optimizations for inference performance on Trn1/Inf2 instances. It takes about 20-60 minutes to complete.

### Table of contents

- *Step 1: Set up Development Environment*
- *Step 2: Download and Convert Checkpoints*
- *Step 3: Deploy with vLLM Inference*
  - *Configurations*
  - *Model Inputs*
  - *Offline Example*
  - *Online Example*

## Step 1: Set up Development Environment

1. Launch a trn1.32xlarge or inf2.48xlarge instance on Ubuntu 22 with Neuron Multi-Framework DLAMI. Please refer to the setup guide if you don't have one yet. If you are looking to install NxD Inference library without using pre-existing DLAMI, please refer to the [NxD Inference Setup Guide](#).
2. Use default virtual environment pre-installed with the Neuron SDK.

```
source /opt/aws_neuronx_venv_pytorch_2_6_nxd_inference/bin/activate
```

3. Install the latest release branch of vLLM from the AWS Neuron fork following the [vLLM User Guide for NxD Inference](#).
4. You should now have the Neuron SDK and other necessary packages installed, including neuronx-distributed-inference, neuronx-cc, torch, torchvision, and vllm-neuronx.

## Step 2: Download and Convert Checkpoints

Download Llama3.2 Multimodal models from either [Meta's official website](#) or HuggingFace(HF) (11B, 90B).

NxDI supports HF checkpoint out-of-the-box. To use the Meta checkpoint, you need to run the following script to convert the downloaded Meta checkpoint into NxDI supported format.

```
python -m neuronx_distributed_inference.models.mllama.convert_mllama_weights_to_neuron \
  --input-dir <path_to_meta_pytorch_checkpoint> \
  --output-dir <path_to_neuron_checkpoint> \
  --instruct \
  --num-shards 8 #(1 for 11B and 8 for 90B)
```

After the script is finished running, you should have the following configuration and checkpoint files. Verify by `ls <path_to_neuron_checkpoint>`:

```
chat_template.json      model.safetensors      tokenizer_config.json
config.json            special_tokens_map.json
generation_config.json tokenizer.json
```

**Note:** The following code examples use HF checkpoint, as it is supported by default, no script needs to be run.

## Step 3: Deploy with vLLM Inference

We provide two examples to deploy Llama3.2 Multimodal with vLLM:

1. Offline inference: you can provide prompts in a python script and execute it.
2. Online inference: you will serve the model in an online server and send requests.

If you already have a compiled model artifact in `MODEL_PATH` with the same specified configuration, or if you have set an environment variable `NEURON_COMPILED_ARTIFACTS`, the vLLM engine will load the compiled model and run inference directly. Otherwise, it will automatically compile and save a new model artifact. See [vLLM User Guide for NxDI Inference](#) for more. We provide example configurations here, continue reading on how to tune them per your use case.

## Configurations

You should specifically tune these configurations when optimizing performance for Llama3.2 Multimodal models. Please refer to [NxDI Inference Features Configuration Guide](#) for detailed explanation of each configuration.

- `MODEL_PATH` - The directory containing NxDI-supported configs, checkpoints, and neuron compiled artifacts.
- `BATCH_SIZE` - Batch size and sequence length together are bounded by device memory. For sequence shorter than 16k, we support up to batch size 4. For longer sequence, we support batch size 1.
- `SEQ_LEN` - The entire sequence length combining input and output sequence. We support sequence length up to 128k for 11B model, and 16k for 90B model.
- `TENSOR_PARALLEL_SIZE` - For best performance, choose the maximum supported value by your instance, that is divisible by the model's hidden sizes and number of attention heads: 32 for `trn1.32xlarge` and 16 for `inf2.48xlarge`.

- `CONTEXT_ENCODING_BUCKETS` - Set based on your distribution of input/context length. For example, suppose 90% of the input traffic is shorter than 1k sequence, and all are less than 2k, then we should set the context encoding buckets to be `[1024, 2048]`.
- `TOKEN_GENERATION_BUCKETS` - Set based on your distribution of entire sequence length. Use similar principle as above.

---

**Note:** Longer sequence takes up more memory, so we should use less buckets. For example, to compile the 90B model on `trn1.32xlarge` with `SEQ_LEN=16384`, `BATCH_SIZE=4`, we can use buckets `[1024, 2048, 16384]` to cover the longest possible sequence as well as shorter sequence where the majority of traffic comes from. We also set an environment variable by `export NEURON_SCRATCHPAD_PAGE_SIZE=1024` to increase the scratchpad size in our direct memory access engine to fit the large tensors.

---

- `SEQUENCE_PARALLEL_ENABLED` - Set to `True` to enable sequence parallel. In principle, sequence parallel helps scaling to long sequence length by splitting tensors along the sequence dimension. However, for short sequence length less than 2k, it is not worth to pay for the collectives overhead when compute workload is manageable. So in this example, as we configured sequence length to be no more than 2k, we disabled the sequence parallel.
- `IS_CONTINUOUS_BATCHING` - Set based on your input traffic. For example, suppose end-to-end latency to generate an entire output sequence (batch size 1) is 1 second in average. However, you receive a request every 0.5 second. Then it is beneficial to enable continuous batching so that new request can get generation started before prior request is finished.
- `ON_DEVICE_SAMPLING_CONFIG` - We enable on-device sampling to perform sampling logic on the Neuron device (rather than on the CPU) to achieve better performance.

## Model Inputs

- `PROMPTS`: `List[str]` - Batch of text prompts.
- `IMAGES`: `List[Union[PIL.Image.Image, torch.Tensor]]` - Batch of image prompts. We currently support one image per prompt as recommended by [Meta](#). If the prompt has no image, use an empty tensor.
- `SAMPLING_PARAMS`: `List[Dict]` - Batch of sampling parameters. With dynamic sampling, you can pass different `top_k`, `top_p`, and `temperature` values for each input in a batch.

## Offline Example

```
import torch
import requests
from PIL import Image

from vllm import LLM, SamplingParams
from vllm import TextPrompt

from neuronx_distributed_inference.models.mllama.utils import add_instruct

def get_image(image_url):
    image = Image.open(requests.get(image_url, stream=True).raw)
    return image
```

(continues on next page)

(continued from previous page)

```

# Configurations
MODEL_PATH = "/home/ubuntu/model_hf/Llama-3.2-90B-Vision-Instruct-hf"
BATCH_SIZE = 4
SEQ_LEN = 2048
TENSOR_PARALLEL_SIZE = 32
CONTEXT_ENCODING_BUCKETS = [1024, 2048]
TOKEN_GENERATION_BUCKETS = [1024, 2048]
SEQUENCE_PARALLEL_ENABLED = False
IS_CONTINUOUS_BATCHING = True
ON_DEVICE_SAMPLING_CONFIG = {"global_topk":64, "dynamic": True, "deterministic": False}

# Model Inputs
PROMPTS = ["What is in this image? Tell me a story",
           "What is the recipe of mayonnaise in two sentences?" ,
           "Describe this image",
           "What is the capital of Italy famous for?",
           ]
IMAGES = [get_image("https://github.com/meta-llama/llama-models/blob/main/models/scripts/
↳resources/dog.jpg?raw=true"),
           torch.empty((0,0)),
           get_image("https://awsdocs-neuron.readthedocs-hosted.com/en/latest/_images/nxd-
↳inference-block-diagram.jpg"),
           torch.empty((0,0)),
           ]
SAMPLING_PARAMS = [dict(top_k=1, temperature=1.0, top_p=1.0, max_tokens=256),
                    dict(top_k=1, temperature=0.9, top_p=1.0, max_tokens=256),
                    dict(top_k=10, temperature=0.9, top_p=0.5, max_tokens=512),
                    dict(top_k=10, temperature=0.75, top_p=0.5, max_tokens=1024),
                    ]

def get_VLLM_mllama_model_inputs(prompt, single_image, sampling_params):
    # Prepare all inputs for mllama generation, including:
    # 1. put text prompt into instruct chat template
    # 2. compose single text and single image prompt into Vllm's prompt class
    # 3. prepare sampling parameters
    input_image = single_image
    has_image = torch.tensor([1])
    if isinstance(single_image, torch.Tensor) and single_image.numel() == 0:
        has_image = torch.tensor([0])

    instruct_prompt = add_instruct(prompt, has_image)
    inputs = TextPrompt(prompt=instruct_prompt)
    inputs["multi_modal_data"] = {"image": input_image}
    # Create a sampling params object.
    sampling_params = SamplingParams(**sampling_params)
    return inputs, sampling_params

def print_outputs(outputs):
    # Print the outputs.
    for output in outputs:
        prompt = output.prompt

```

(continues on next page)

(continued from previous page)

```

generated_text = output.outputs[0].text
print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

if __name__ == '__main__':
    assert len(PROMPTS) == len(IMAGES) == len(SAMPLING_PARAMS), \
        f"""Text, image prompts and sampling parameters should have the same batch size,
        got {len(PROMPTS)}, {len(IMAGES)}, and {len(SAMPLING_PARAMS)}"""

    # Create an LLM.
    llm = LLM(
        model=MODEL_PATH,
        max_num_seqs=BATCH_SIZE,
        max_model_len=SEQ_LEN,
        block_size=SEQ_LEN,
        device="neuron",
        tensor_parallel_size=TENSOR_PARALLEL_SIZE,
        override_neuron_config={
            "context_encoding_buckets": CONTEXT_ENCODING_BUCKETS,
            "token_generation_buckets": TOKEN_GENERATION_BUCKETS,
            "sequence_parallel_enabled": SEQUENCE_PARALLEL_ENABLED,
            "is_continuous_batching": IS_CONTINUOUS_BATCHING,
            "on_device_sampling_config": ON_DEVICE_SAMPLING_CONFIG,
        }
    )

    batched_inputs = []
    batched_sample_params = []
    for pmpt, img, params in zip(PROMPTS, IMAGES, SAMPLING_PARAMS):
        inputs, sampling_params = get_VLLM_mllama_model_inputs(pmpt, img, params)
        # test batch-size = 1
        outputs = llm.generate(inputs, sampling_params)
        print_outputs(outputs)
        batched_inputs.append(inputs)
        batched_sample_params.append(sampling_params)

    # test batch-size = 4
    outputs = llm.generate(batched_inputs, batched_sample_params)
    print_outputs(outputs)

```

This script will print the outputs. Below is an example output from image-text prompt:

```

Prompt: '<|begin_of_text|><|start_header_id|>user<|end_header_id|>\n\n<|image|>What is
in this image? Tell me a story<|eot_id|><|start_header_id|>assistant<|end_header_id|>\n\n
→',
Generated text: 'The image shows a dog riding a skateboard. The dog is standing on the
skateboard, which is in the middle of the road. The dog is looking at the camera with its
mouth open, as if it is smiling. The dog has floppy ears and a long tail. It is wearing a
collar around its neck. The skateboard is black with red wheels. The background is_
→blurry,
but it appears to be a city street with buildings and cars in the distance.'

```



## Online Example

First, open a terminal and spin up a server of the model. If you specify a new set of configurations, a new neuron model artifact will be compiled now.

```
MODEL_PATH="/home/ubuntu/model_hf/Llama-3.2-90B-Vision-Instruct-hf"
python3 -m vllm.entrypoints.openai.api_server \
  --model $MODEL_PATH \
  --tensor-parallel-size 32 \
  --max-model-len 2048 \
  --max-num-seqs 4 \
  --device neuron \
  --override-neuron-config '{
    "context_encoding_buckets": [1024, 2048],
    "token_generation_buckets": [1024, 2048],
    "sequence_parallel_enabled": false,
    "is_continuous_batching": true,
    "on_device_sampling_config": {
      "global_topk": 64,
      "dynamic": true,
      "deterministic": false
    }
  }'
```

If you see the below logs, that means your server is up and running:

```
INFO: Started server process [284309]
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: Uvicorn running on http://0.0.0.0:8000 (Press CTRL+C to quit)
```

Then open a new terminal as the client where you can send requests to the server. We've enabled continuous batching by default, so you can open up to --max-num-seqs client terminals to send requests. To send a text-only request:

```
MODEL_PATH="/home/ubuntu/model_hf/Llama-3.2-90B-Vision-Instruct-hf"
curl http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "'$MODEL_PATH'",
  "messages": [
    {
      "role": "user",
      "content": "What is the capital of Italy?"
    }
  ]
}'
```

You should receive outputs shown in the client terminal shortly:

```
{
  "id": "chat-2df3e876738b470ab27b090e0a09736e",
  "object": "chat.completion",
  "created": 1734401826,
  "model": "/home/ubuntu/model_hf/Llama-3.2-90B-Vision-Instruct-hf/",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "The capital of Italy is Rome."
      },
      "logprobs": null,
      "finish_reason": "stop",
      "stop_reason": null
    }
  ]
}
```

(continues on next page)

(continued from previous page)

```
"usage":{"prompt_tokens":42,"total_tokens":50,"completion_tokens":8},"prompt_logprobs":
↪null}
```

If the request fails, increase the value of the `VLLM_RPC_TIMEOUT` environment variable using `export VLLM_RPC_TIMEOUT=180000`, then restart the server. The timeout value depends on the model and deployment configuration used.

To send a request with both text and image prompts:

```
curl http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "'$MODEL_PATH'",
  "messages": [
    {
      "role": "user",
      "content": [
        {
          "type": "text",
          "text": "Describe this image"
        },
        {
          "type": "image_url",
          "image_url": {
            "url": "https://awsdocs-neuron.readthedocs-hosted.com/en/latest/_images/
↪nxd-inference-block-diagram.jpg"
          }
        }
      ]
    }
  ]
}'
```

You can expect results appear in the client terminal shortly:

```
{"id":"chat-fd1319865bd44d6aa60a4739cce61c9d","object":"chat.completion",
"created":1734401984,"model":"/home/ubuntu/model_hf/Llama-3.2-90B-Vision-Instruct-hf/",
"choices":[{"index":0,"message":{"role":"assistant","content":"The image presents a
diagram illustrating the components of NxD Inference, with a focus on inference modules
and additional modules. The diagram is divided into two main sections: \"Inference
Modules\" and \"Additional Modules.\" \n\n**Inference Modules:**\n\n  Attention
Techniques\n*   KV Caching\n*   Continuous Batching\n\n**Additional Modules:**\n\n*
Speculative Decoding (Draft model and Draft heads (Medusa / Eagle))\n\nThe diagram also
includes a section titled \"NxD Core (Distributed Strategies, Distributed Model Tracing)\n
↪\"
and a logo for PyTorch at the bottom."},"tool_calls":[],"logprobs":null,
"finish_reason":"stop","stop_reason":null}], "usage":{"prompt_tokens":14,"total_tokens":
↪137,
"completion_tokens":123},"prompt_logprobs":null}
```

## Tutorial: Using Speculative Decoding to improve Llama-3.3-70B inference performance on Trn2 instances

NeuronX Distributed (NxD) Inference allows you to deploy Llama3.3 70B on a single Trn2 or Trn1 instance. This tutorial provides a step-by-step guide to deploy Llama3.3 70B on a Trn2 instance using two different configurations, one without speculative decoding and the other with draft model based speculative decoding enabled (with Llama-3.2 1B as the draft model). We will also measure performance by running a load test using LLMPPerf and compare key metrics between the two configurations. While this tutorial uses batch size 1 for demonstration purposes, the model configuration provides support for batch sizes up to 4.

### Table of contents

- *Prerequisites:*
  - *Set up and connect to a Trn2.48xlarge instance*
  - *Install packages*
- *Scenario 1: Run Llama3.3 70B on Trn2*
  - *Step 1: Compile the model*
  - *Step 2: Run the model using vLLM*
  - *Step 3: Measure performance using LLMPPerf*
- *Scenario 2: Run Llama3.3 70B on Trn2 with Speculative Decoding*
  - *Step 1: Compile the model*
  - *Step 2: Run the model using vLLM*
  - *Step 3: Measure performance using LLMPPerf*
- *Conclusion*

### Prerequisites:

#### Set up and connect to a Trn2.48xlarge instance

As a prerequisite, this tutorial requires that you have a Trn2 instance created from a Deep Learning AMI that has the Neuron SDK pre-installed.

To set up a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK, see [NxD Inference Setup Guide](#).

After setting up an instance, use SSH to connect to the Trn2 instance using the key pair that you chose when you launched the instance.

After you are connected, activate the Python virtual environment that includes the Neuron SDK.

```
source ~/aws_neuronx_venv_pytorch_2_5_nxd_inference/bin/activate
```

Run `pip list` to verify that the Neuron SDK is installed.

```
pip list | grep neuron
```

You should see Neuron packages including `neuronx-distributed-inference` and `neuronx-cc`.

## Install packages

NxD Inference supports running models with vLLM. This functionality is available in the AWS Neuron fork of the vLLM GitHub repository. Install the latest release branch of vLLM from the AWS Neuron fork following instructions in the [vLLM User Guide for NxD Inference](#).

In this tutorial, you will use [llmperf](#) to measure the performance. We will use the [load test](#) feature of LLMPerf and measure the performance for accepting 10,000 tokens as input and generating 1500 tokens as output. Install llmperf into the virtual environment.

```
git clone https://github.com/ray-project/llmperf.git
cd llmperf
pip install -e .
```

## Download models

To use this sample, you must first download a 70B model checkpoint from Hugging Face to a local path on the Trn2 instance. For more information, see [Downloading models](#) in the Hugging Face documentation. You can download and use [meta-llama/Llama-3.3-70B-Instruct](#) for this tutorial.

Since we will be using Speculative Decoding in the second configuration, you will also need a draft model checkpoint. You can download and use [meta-llama/Llama-3.2-1B-Instruct](#).

---

**Note:** NxD Inference supports batch sizes up to 4 for this model configuration. To determine the optimal batch size for your specific use case, we recommend incrementally testing batch sizes from 1 to 4 while monitoring your application's performance metrics.

---

## Scenario 1: Run Llama3.3 70B on Trn2

In this scenario, you will run Llama3.3 70B on Trn2 without Speculative Decoding using bfloat16 precision.

### Step 1: Compile the model

We will first compile and run generation on a sample prompt using a command installed by `neuronx-distributed-inference`. Save the contents of the below script to your favorite shell script file, for example, `compile_model.sh` and then run it.

Note that we are using the following features as described in the tutorial for running 405B model [Tutorial: Deploying Llama3.1 405B \(Trn2\)](#)

- Logical NeuronCore Configuration (LNC)
- Tensor parallelism (TP) on Trn2
- Optimized Kernels

The script compiles the model and runs generation on the given input prompt. Note the path we used to save the compiled model. This path should be used when launching vLLM server for inference so that the compiled model can be loaded without recompilation. Please refer to [NxD Inference API Reference](#) for more information on these `inference_demo` flags.

---

**Note:** Known issue: Using kernels with bucket length of 1024 or less may lead to `Numerical Error` in inference.

---

```
RuntimeError: Failed to execute the model status=1003 message=Numerical Error
```

```
# Replace this with the path where you downloaded and saved the model files.
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
# This is where the compiled model will be saved. The same path
# should be used when launching vLLM server for inference.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

NUM_CORES=128
TP_DEGREE=64
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600
export XLA_DENSE_GATHER_FACTOR=0
export NEURON_RT_INSPECT_ENABLE=0

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
    --model-path $MODEL_PATH \
    --compiled-model-path $COMPILED_MODEL_PATH \
    --torch-dtype bfloat16 \
    --start_rank_id 0 \
    --local_ranks_size $TP_DEGREE \
    --tp-degree $TP_DEGREE \
    --batch-size 1 \
    --max-context-length 12288 \
    --seq-len 12800 \
    --on-device-sampling \
    --top-k 1 \
    --do-sample \
    --fused-qkv \
    --sequence-parallel-enabled \
    --qkv-kernel-enabled \
    --attn-kernel-enabled \
    --mlp-kernel-enabled \
    --cc-pipeline-tiling-factor 1 \
    --pad-token-id 2 \
    --enable-bucketing \
    --context-encoding-buckets 2048 4096 8192 12288 \
      --token-generation-buckets 2048 4096 8192 12800 \
    --prompt "What is annapurna labs?" 2>&1 | tee log
```

## Step 2: Run the model using vLLM

After compiling the model, you can run the model using vLLM. Save the contents of the below script to another shell script file, for example, `start_vllm.sh` and then run it.

```
export NEURON_RT_INSPECT_ENABLE=0
export NEURON_RT_VIRTUAL_CORE_SIZE=2

# These should be the same paths used when compiling the model.
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs 1 \
    --max-model-len 128000 \
    --tensor-parallel-size 64 \
    --device neuron \
    --use-v2-block-manager \
    --override-neuron-config '{"on_device_sampling_config": {"do_sample": true}, \
    ↪ "skip_warmup": true}' \
    --port 8000 &
PID=$!
echo "vLLM server started with PID $PID"
```

## Step 3: Measure performance using LLMPerf

After the above steps, the vllm server should be running. You can now measure the performance using LLMPerf. Before we can use the `llmperf` package, we need to make a few changes to its code. Follow [benchmarking with LLMPerf guide](#) to apply the code changes.

Below is a sample shell script to run LLMPerf. To provide the model with 10000 tokens as input and generate 1500 tokens as output on average, we use the following parameters from LLMPerf:

```
--mean-input-tokens 10000 \
--mean-output-tokens 1500 \
```

More information about several arguments used in the script can be found in the [llmperf open source code](#).

```
# This should be the same path to which the model was downloaded (also used in the above ↪
↪ steps).
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
# This is the name of directory where the test results will be saved.
OUTPUT_PATH=llmperf-results-sonnets

export OPENAI_API_BASE="http://localhost:8000/v1"
export OPENAI_API_KEY="mock_key"

python token_benchmark_ray.py \
    --model $MODEL_PATH \
```

(continues on next page)

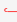
(continued from previous page)

```

--mean-input-tokens 10000 \
--stddev-input-tokens 0 \
--mean-output-tokens 1500 \
--stddev-output-tokens 0 \
--num-concurrent-requests 1\
--timeout 3600 \
--max-num-completed-requests 50 \
--tokenizer $MODEL_PATH \
--additional-sampling-params '{}' \
--results-dir $OUTPUT_PATH \
--llm-api "openai"

```

A sample output from the above script is shown below:

Results **for** token benchmark **for** /home/ubuntu/models/Llama-3.3-70B-Instruct/ queried **with**  the openai api.

```

inter_token_latency_s
p25 = 0.01964743386193489
p50 = 0.01965969146322459
p75 = 0.019672998415771872
p90 = 0.01969826815724373
p95 = 0.019810569172135244
p99 = 0.020350346909947692
mean = 0.01969182239660784
min = 0.0196275211258056
max = 0.020702997242410977
stddev = 0.00015700734112322808
ttft_s
p25 = 0.8109508841298521
p50 = 0.8142827898263931
p75 = 30.46490489714779
p90 = 30.513100237119943
p95 = 30.521608413150535
p99 = 48.876512633068415
mean = 11.503728219866753
min = 0.8080519903451204
max = 66.4881955627352
stddev = 15.692731777293613
end_to_end_latency_s
p25 = 30.296781020238996
p50 = 30.326033774763346
p75 = 59.9560666854959
p90 = 60.001504834741354
p95 = 60.028880204679446
p99 = 79.1842334462329
mean = 41.04328096391633
min = 30.265212223865092
max = 97.54387667682022
stddev = 15.796048923358924
request_output_throughput_token_per_s
p25 = 25.044969421803977

```

(continues on next page)

(continued from previous page)

```
p50 = 49.49542857484997
p75 = 49.543217224244
p90 = 49.583184869985566
p95 = 49.58588728343319
p99 = 49.592597790896676
mean = 40.91042833304163
min = 15.387946954098137
max = 49.59489426003143
stddev = 11.825984480587056
number_input_tokens
p25 = 10000.0
p50 = 10000.0
p75 = 10000.0
p90 = 10000.0
p95 = 10000.0
p99 = 10000.0
mean = 10000.0
min = 10000
max = 10000
stddev = 0.0
number_output_tokens
p25 = 1501.0
p50 = 1501.0
p75 = 1501.0
p90 = 1501.0
p95 = 1501.0
p99 = 1502.02
mean = 1501.04
min = 1501
max = 1503
stddev = 0.282842712474619
Number Of Errored Requests: 0
Overall Output Throughput: 36.55567822866449
Number Of Completed Requests: 50
Completed Requests Per Minute: 1.4612140207588533
```

## Scenario 2: Run Llama3.3 70B on Trn2 with Speculative Decoding

In this scenario, you will run Llama3.3 70B on Trn2 with Speculative Decoding. Specifically, we will use the below variations from the supported variants as described in *Speculative Decoding*

- Speculative Decoding with Llama-3.2-1B as the draft model *Speculative Decoding with a Draft model*
- Fused Speculation for improved performance *Fused Speculation*



## Step 1: Compile the model

When compiling the model to use speculative decoding, you need to provide a draft model checkpoint and a few additional parameters to the `inference_demo` command.

For a quick review, here are the additional arguments provided:

```
--draft-model-path $DRAFT_MODEL_PATH \
--enable-fused-speculation \
--speculation-length 7 \
```

Please refer to [NxN Inference API Reference](#) for more information on these `inference_demo` flags. The complete script to compile the model for this configuration is shown below:

**Note:** Known issue: Using kernels with bucket length of 1024 or less may lead to Numerical Error in inference.

```
RuntimeError: Failed to execute the model status=1003 message=Numerical Error
```

```
# This is the same path as in the previous scenario.
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
# This is the path where the draft model is downloaded and saved.
DRAFT_MODEL_PATH="/home/ubuntu/models/Llama-3.2-1B-Instruct/"
# As in the previous scenario, this is where the compiled model will be saved.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

NUM_CORES=128
TP_DEGREE=64
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600
export XLA_DENSE_GATHER_FACTOR=0
export NEURON_RT_INSPECT_ENABLE=0

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
  --model-path $MODEL_PATH \
  --compiled-model-path $COMPILED_MODEL_PATH \
  --torch-dtype bfloat16 \
  --start_rank_id 0 \
  --local_ranks_size $TP_DEGREE \
  --tp-degree $TP_DEGREE \
  --batch-size 1 \
  --max-context-length 12288 \
  --seq-len 12800 \
  --on-device-sampling \
  --top-k 1 \
  --fused-qkv \
  --sequence-parallel-enabled \
```

(continues on next page)

(continued from previous page)

```

--qkv-kernel-enabled \
--attn-kernel-enabled \
--mlp-kernel-enabled \
--cc-pipeline-tiling-factor 1 \
--draft-model-path $DRAFT_MODEL_PATH \
--enable-fused-speculation \
--speculation-length 7 \
--pad-token-id 2 \
--enable-bucketing \
--context-encoding-buckets 2048 4096 8192 12288 \
    --token-generation-buckets 2048 4096 8192 12800 \
--prompt "What is annapurna labs?" 2>&1 | tee log

```

## Step 2: Run the model using vLLM

Similar to compiling the model, we need to specify parameters specific to speculative decoding when running the model using vLLM.

For a quick glance, these are the parameters that are different for running vLLM server with model compiled using speculative decoding:

```

--speculative-max-model-len 12800 \
--speculative-model $DRAFT_MODEL_PATH \
--num-speculative-tokens 7 \
--override-neuron-config '{"enable_fused_speculation":true}' \

```

Here is the complete script to run the model using vLLM with speculative decoding:

```

export NEURON_RT_INSPECT_ENABLE=0
export NEURON_RT_VIRTUAL_CORE_SIZE=2

# These should be the same paths used when compiling the model.
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
DRAFT_MODEL_PATH="/home/ubuntu/models/Llama-3.2-1B-Instruct/"
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs 1 \
    --max-model-len 12800 \
    --tensor-parallel-size 64 \
    --device neuron \
    --speculative-max-model-len 12800 \
    --speculative-model $DRAFT_MODEL_PATH \
    --num-speculative-tokens 7 \
    --use-v2-block-manager \
    --override-neuron-config '{"enable_fused_speculation":true}' \
    --port 8000 &
PID=$!

```

(continues on next page)

(continued from previous page)

```
echo PID=$PID
echo "vLLM server started with PID $PID"
```

### Step 3: Measure performance using LLMPerf

The script to measure the performance using LLMPerf is same as the one used in the first scenario. Before we can use the llmperf package, we need to make a few changes to its code. Follow [benchmarking with LLMPerf guide](#) to apply the code changes.

For convenience, here's the script once again:

```
# This should be the same path to which the model was downloaded (also used in the above
↳ steps).
MODEL_PATH="/home/ubuntu/models/Llama-3.3-70B-Instruct/"
# This is the name of directory where the test results will be saved. Use a different
↳ name for this scenario.
OUTPUT_PATH=llmperf-results-sonnets-speculative

export OPENAI_API_BASE="http://localhost:8000/v1"
export OPENAI_API_KEY="mock_key"

python token_benchmark_ray.py \
    --model $MODEL_PATH \
    --mean-input-tokens 10000 \
    --stddev-input-tokens 0 \
    --mean-output-tokens 1500 \
    --stddev-output-tokens 0 \
    --num-concurrent-requests 1 \
    --timeout 3600 \
    --max-num-completed-requests 50 \
    --tokenizer $MODEL_PATH \
    --additional-sampling-params '{}' \
    --results-dir $OUTPUT_PATH \
    --llm-api "openai"
```

A sample output from the above script is shown below:

Results **for** token benchmark **for** /home/ubuntu/models/Llama-3.3-70B-Instruct/ queried **with**
↳ the openai api.

```
inter_token_latency_s
p25 = 0.0053349758717231455
p50 = 0.005386366705410183
p75 = 0.005441084293027719
p90 = 0.005499971026182175
p95 = 0.005520176071580499
p99 = 0.005911254031351169
mean = 0.00540780140378178
min = 0.005264532127728065
max = 0.006265544256816307
stddev = 0.00013951778334019935
```

(continues on next page)

(continued from previous page)

```
ttft_s
  p25 = 0.8693495176266879
  p50 = 0.870149074587971
  p75 = 0.8710820493288338
  p90 = 0.8725412225350737
  p95 = 0.8742059985175729
  p99 = 36.83790613239617
  mean = 2.280795605443418
  min = 0.8676468348130584
  max = 71.38881027325988
  stddev = 9.97280539681726
end_to_end_latency_s
  p25 = 8.873123338911682
  p50 = 8.950916013680398
  p75 = 9.030085149221122
  p90 = 9.120021602977067
  p95 = 9.150626054406166
  p99 = 45.70815015356973
  mean = 10.393093119114637
  min = 8.766328778117895
  max = 80.78758085798472
  stddev = 10.158917239418473
request_output_throughput_token_per_s
  p25 = 166.22213179149702
  p50 = 167.69243252025473
  p75 = 169.16253286110174
  p90 = 169.52692450439133
  p95 = 169.81518762962915
  p99 = 170.85438941846397
  mean = 164.631719334475
  min = 18.579588397857652
  max = 171.2233293995004
  stddev = 21.152953887186314
number_input_tokens
  p25 = 10000.0
  p50 = 10000.0
  p75 = 10000.0
  p90 = 10000.0
  p95 = 10000.0
  p99 = 10000.0
  mean = 10000.0
  min = 10000
  max = 10000
  stddev = 0.0
number_output_tokens
  p25 = 1501.0
  p50 = 1501.0
  p75 = 1501.0
  p90 = 1501.0
  p95 = 1501.0
  p99 = 1502.02
  mean = 1501.04
```

(continues on next page)

(continued from previous page)

```

min = 1501
max = 1503
stddev = 0.282842712474619
Number Of Errored Requests: 0
Overall Output Throughput: 144.17136914316023
Number Of Completed Requests: 50
Completed Requests Per Minute: 5.76285918335928

```

## Conclusion

As seen in the table below, TPOT reduced by 3.6x and output token throughput increased by 4x when using speculative decoding with draft model combined with fused speculative decoding, compared to baseline without speculative decoding. Please note that batch size of 1 is used in this tutorial for computing the below metrics.

| Scenario (all using BF16)                       | TTFT (P50 in ms) | TPOT (P50 in ms) | Output token Throughput (per second) |
|---|------------------|------------------|--------------------------------------|
| No speculative decoding                         | 814.2            | 19.6             | 36                                   |
| Fused speculative decoding (Llama 3.2 1B Draft) | 870.1            | 5.3              | 144                                  |

## Tutorial: Scaling LLM Inference with Data Parallelism on Trn2

### Introduction

This tutorial demonstrates how to implement data parallelism (DP) LLM inference with multiple model copies on Neuron. The following sections provides a sequence of steps to stand up multiple Llama 3.3 70B model endpoints on a single *trn2.48xlarge* instance with NxD Inference and vLLM and run data parallel inference.

### Data Parallel Inference

We can achieve Data Parallelism by using multiple copies of the same model hosted on the instance to process multiple requests simultaneously. Using NxD Inference and vLLM, you can deploy multiple model endpoints by adjusting the tensor parallel degree (Tensor Parallelism (TP) refers to sharding model weight matrices onto multiple NeuronCores within each model copy) and allocating appropriate NeuronCore ranges for each model endpoint. While increasing the batch size with a single copy of the model increases throughput, introducing data parallelism with multiple model endpoints combined with tensor parallelism allows further increase in instance throughput with some impact to latency. Use this technique when you can relax the latency constraint of your application to further maximize the throughput of the instance.

In this tutorial we use Llama 3.3 70B with DP=2 and TP=32, However, you can follow the same sequence of steps to deploy additional model copies by appropriately changing the tensor parallel degree. You can also use this guide to deploy multiple copies of any other models on Trn1 or Inf2 instances as long as the model fits and the DP x TP degree does not exceed the number of model cores.

### Prerequisites:

#### Setup and Connect to an Amazon EC2 Trn2 Instance

An Amazon EC2 `trn2.48xlarge` instance with AWS Neuron SDK version 2.23.0 or later ([Neuron 2.24.1 \(06/30/2025\)](#)) is required.

To launch a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK and NxD Inference dependencies, see [NxD Inference Setup Guide](#).

Make sure to activate the Neuron virtual environment

```
source /opt/aws_neuronx_venv_pytorch_2_6_nxd_inference/bin/activate
```

To verify that NxD Inference has installed successfully, check that you can run the `inference_demo` console script.

```
inference_demo --help
```

#### Download Model Weights

To use this tutorial, you must first download a Llama 3.3 70B Instruct model checkpoint from Hugging Face to a local path on the Trn2 instance. For more information, see [Downloading Models](#) in the Hugging Face documentation. You can download and use `meta-llama/Llama-3.3-70B-Instruct` for this tutorial.

#### Install Neuron vLLM Fork

NxD Inference supports running models with vLLM. This functionality is available in the AWS Neuron fork of the vLLM GitHub repository. Install the latest release branch of vLLM from the AWS Neuron fork following instructions in the [vLLM User Guide for NxD Inference](#).

#### Install LLMPerf

In this tutorial, you will use `LLMPerf` to measure the performance.

Install `llmperf` into the virtual environment.

```
git clone --branch v2.0 https://github.com/ray-project/llmperf.git
cd llmperf
pip install -e .
```

Once you have installed LLMPerf, please apply relevant patches as described in [LLM Inference Benchmarking guide](#). Ensure that you apply all the patches described there including the data parallelism support patch.

## Step-by-Step Tutorial Instructions

### Step 1: Compile the model

Before we launch the model endpoint with vLLM, we'll use the NxD Inference library to compile the model with an appropriate configuration. Refer to [NxD Inference Features Configuration Guide](#) for more information. To compile a model for data parallelism inference, set the `NUM_CORES`, `TP_DEGREE`, `BATCH_SIZE` to allow for strategic workflow distribution. For `DP=2` with `BATCH_SIZE>=1`, `TP_DEGREE` should be set to  $64/2=32$  to maximize NeuronCore utilization across all model copies. Simply create and run a shell script as illustrated below:

*compile\_model.sh*

```
#!/bin/bash
# Replace with path to your downloaded Hugging Face model checkpoints
MODEL_PATH="/ubuntu/model_hf/Llama-3.3-70B-Instruct/"

# This is where the compiled model will be saved. The same path
# should be used when launching vLLM server for inference.
export COMPILED_MODEL_PATH="/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

NUM_CORES=128
TP_DEGREE=32
LNC=2
BATCH_SIZE=4

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600
export XLA_DENSE_GATHER_FACTOR=0
export NEURON_RT_INSPECT_ENABLE=0

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
  --model-path $MODEL_PATH \
  --compiled-model-path $COMPILED_MODEL_PATH \
  --torch-dtype bfloat16 \
  --start_rank_id 0 \
  --local_ranks_size $TP_DEGREE \
  --tp-degree $TP_DEGREE \
  --batch-size $BATCH_SIZE \
  --max-context-length 8192 \
  --seq-len 8192 \
  --on-device-sampling \
  --top-k 1 \
  --do-sample \
  --fused-qkv \
  --qkv-kernel-enabled \
  --attn-kernel-enabled \
  --mlp-kernel-enabled \
  --pad-token-id 2 \
  --compile-only \
```

(continues on next page)

(continued from previous page)

```
--prompt "What is annapurna labs?" 2>&1 | tee log
```

To compile the model, run this script with command: `./compile_model.sh`

It's important to specify the path to which the compiled model is saved, as this same path must be used when you later launch the vLLM server for inference, allowing you to use the pre-compiled model without having to compile it again.

**Note:** To run this script on trn1, set LNC=1. For more information about LNC, see [Logical NeuronCore configuration](#). Also appropriately change NUM\_CORES & TP\_DEGREE (eg. 16 for DP=2)

For detailed information about the inference\_demo flags, you can consult the [Nx D Inference API Reference](#).

## Step 2: Launch model endpoints

Create a deployment script (`deploy_vllm_endpoint.sh`) containing below code snippet that configures and launches a model endpoint. The script is parameterized so that you can pass a specific port number, range of neuron cores, tensor parallel degree and batch size.

Key Parameters Explained:

- **MODEL\_PATH:** The Hugging Face model identifier or local model\_hf path containing Meta-Llama-3.3-70B-Instruct hugging face checkpoints. Eg. `/home/ubuntu/model_hf/Llama-3.3-70B-Instruct/`
- **port:** Network port for the endpoint Eg. 8000. The port number should be unique for each model endpoint.
- **cores:** Range of NeuronCores allocated to this endpoint. This should be a non overlapping range of cores when deploying multiple model endpoints on the same instance. For example, when allocated 32 NeuronCores to a model endpoint specify 0-31 or 32-63.
- **tp\_degree:** Degree of tensor parallelism for model sharding. To maximize NeuronCores utilization, reduce tp\_degree while increasing dp\_degree.
- **bs :** Batch size specified for model endpoint.

These parameters should match the values used during compilation step above.

*deploy\_vllm\_endpoint.sh*

```
# Model deployment script with detailed configuration

# Default values for arguments
DEFAULT_PORT=8000
DEFAULT_CORES="0-31"
DEFAULT_TP_DEGREE=32
DEFAULT_BS=4

# Help function
show_help() {
    echo "Usage: $0 [options]"
    echo "Options:"
    echo "  -p port          Port number for vLLM endpoint (default: $DEFAULT_PORT)"
    echo "  -c cores         Range of neuron cores (default: $DEFAULT_CORES)"
    echo "  -t tp_degree     Tensor parallel degree (default: $DEFAULT_TP_DEGREE)"
    echo "  -b bs           Batch size (default: $DEFAULT_BS)"
    echo "  -h              Show this help message"
```

(continues on next page)



(continued from previous page)

```

}

# Parse single-letter arguments
while getopts "p:c:t:b:h" opt; do
    case $opt in
        p) port="$OPTARG" ;;
        c) cores="$OPTARG" ;;
        t) tp_degree="$OPTARG" ;;
        b) bs="$OPTARG" ;;
        h) show_help; exit 0 ;;
        ?) show_help; exit 1 ;;
    esac
done

# Set defaults if not provided
port=${port:-$DEFAULT_PORT}
cores=${cores:-$DEFAULT_CORES}
tp_degree=${tp_degree:-$DEFAULT_TP_DEGREE}
bs=${bs:-$DEFAULT_BS}

# Environment configurations
export NEURON_RT_INSPECT_ENABLE=0
export NEURON_RT_VIRTUAL_CORE_SIZE=2

# These should be the same paths used when compiling the model.
MODEL_PATH="/ubuntu/model_hf/Llama-3.3-70B-Instruct/"
COMPILED_MODEL_PATH="/ubuntu/traced_model/Llama-3.3-70B-Instruct/"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
export NEURON_RT_VISIBLE_CORES=${cores}

VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs ${bs} \
    --max-model-len 12800 \
    --tensor-parallel-size ${tp_degree} \
    --device neuron \
    --use-v2-block-manager \
    --override-neuron-config '{"on_device_sampling_config": {"do_sample": true, \
    ↪ "global_topk": 64}}' \
    --port ${port} &
PID=$!
echo "vLLM server started with PID $PID"

```

Run this script to launch 2 vLLM servers. You can run these commands as background processes in the same terminal or run two separate terminals for each command. We launch two servers, each with a tensor parallel degree of 32 and batch size of 4. Note that the first vLLM server uses neuron cores 0-31 and the second one 32-63. You can pick any ports that are available.

```
./deploy_vllm_endpoint.sh -p 8000 -c 0-31 -t 32 -b 4 &
```

and

```
./deploy_vllm_endpoint.sh -p 8001 -c 32-63 -t 32 -b 4 &
```

The server start up time can take a few minutes since the model weights are getting loaded. Once the vLLM servers have been launched, you should see the following log output. This implies that the model server has been deployed.

```
INFO:      Started server process [221607]
INFO:      Waiting for application startup.
INFO:      Application startup complete.
INFO:      Uvicorn running on http://0.0.0.0:8000 (Press CTRL+C to quit)
```

### Step 3: Benchmark the deployed model endpoints

After the above steps, the vLLM server should be running. You can now measure the performance using LLMPerf. Ensure you have made the required changes to use LLMPerf with DP>1 by following [Install LLMPerf](#)

Below is a sample shell script to run LLMPerf. The script allows the user to specify tensor parallelism degree, data parallelism degree, and batch size through command-line arguments, with default values provided. It calculates the concurrency based on batch size and data parallelism, sets up the environment for benchmarking with input tokens N(7936, 30) and output tokens N(256,30), and then runs LlmPerf's `token_benchmark_ray.py` with various parameters to measure the model endpoints' performance. The benchmark simulates requests with specific input and output token distributions, and collects results for analysis.

More information about several arguments used in the script can be found in the [llmperf open source code](#)

*benchmark\_model.sh*

```
#!/bin/bash

# Default values for arguments
DEFAULT_TP_DEGREE=32
DEFAULT_DP_DEGREE=2
DEFAULT_BS=1

# Help function
show_help() {
    echo "Usage: $0 [options]"
    echo "Options:"
    echo "  -t tp_degree      Tensor parallel degree (default: $DEFAULT_TP_DEGREE)"
    echo "  -d dp_degree      Data parallel degree (default: $DEFAULT_DP_DEGREE)"
    echo "  -b bs             Batch size (default: $DEFAULT_BS)"
    echo "  -h               Show this help message"
}

# Parse single-letter arguments
while getopts "t:d:b:h" opt; do
    case $opt in
        t) tp_degree="$OPTARG" ;;
        d) dp_degree="$OPTARG" ;;
        b) bs="$OPTARG" ;;
        h) show_help; exit 0 ;;
        ?) show_help; exit 1 ;;
    esac
done
```

(continues on next page)

(continued from previous page)

```

# Set defaults if not provided
tp_degree=${tp_degree:-$DEFAULT_TP_DEGREE}
dp_degree=${dp_degree:-$DEFAULT_DP_DEGREE}
bs=${bs:-$DEFAULT_BS}

# Calculate total concurrent requests (batch_size * data_parallelism)
# If result is less than 1, default to batch_size
concurrency=$(awk -v batch="$bs" -v dp_degree="$dp_degree" 'BEGIN {
    concurrency = int(batch * dp_degree)
    print (concurrency >= 1 ? concurrency : batch)
}')
```

echo "concurrency: \$concurrency"

MODEL\_PATH="/ubuntu/model\_hf/Llama-3.3-70B-Instruct/"

# Modify OpenAI's API key and API base to use vLLM's API server.

export OPENAI\_API\_KEY=EMPTY

#if you have more vLLM servers, append the required number of ports like so:

http://localhost:8001/v1;http://localhost:8002/v1"

export OPENAI\_API\_BASE="http://0.0.0.0:8000/v1;http://0.0.0.0:8001/v1"

```

python /root/llmperf/token_benchmark_ray.py \
--model ${MODEL_PATH} \
--mean-input-tokens 7936 \
--stddev-input-tokens 30 \
--mean-output-tokens 256 \
--stddev-output-tokens 30 \
--num-concurrent-requests ${concurrency} \
--results-dir "/ubuntu/results/" \
--timeout 21600 \
--max-num-completed-requests 1000 \
--additional-sampling-params '{"temperature": 0.7, "top_k": 50}' \
--llm-api "openai"
```

Run this script with `./benchmark_model.sh -t 32 -d 2 -b 4`. These args match the args set while launching vLLM servers above.

Once the script starts executing, you will see output like:

```

INFO worker.py:1852 -- Started a local Ray instance.
4%|          | 39/1000 [01:29<30:14, 1.89s/it]
```

Once benchmarking is complete, results can be found in the directory specified with the `--results-dir` flag in the `benchmark_vllm.sh` script

## Conclusion

This tutorial demonstrates how

data parallelism using multiple model copies can help increase the throughput. While standard batching (DP=1, BS>1) processes multiple requests through a single model copy, data parallelism deploys multiple independent model copies that can process different requests simultaneously. Our experiments with batch sizes 1 & 4 show that as we decrease Tensor Parallelism (TP) from 64 to 16 and increase Data Parallelism (DP) from 1 to 4, we see up to 2x throughput improvement with non optimized configurations. However, this comes with an increase in Time To First Token (TTFT) latency. This illustrates a key consideration: while DP can improve overall system throughput by processing more concurrent requests, it can lead to higher latency

When to choose Data parallel with multiple model copies over using single model copy in an instance:

- Use DP when your workload is collective-bound rather than memory or compute-bound. At high batch sizes, TP64 / TP128 collectives can become slow due to the number of hops and increasing throughput requirements. At high enough batch size, it can be better to pay the cost of duplicated weight loads and use DP with multiple model copies in order to reduce collective latencies.
- Consider DP when you need to handle many concurrent requests and can tolerate moderate latency increases

Implementation requires careful consideration of your total memory budget, as each additional model copy increases memory consumption. You'll need to balance the number of model copies against the resources allocated to each model copy based on your specific throughput and latency requirements. By understanding these trade-offs and following the implementation guidelines in this tutorial, users can select the most appropriate approach for their specific use case and optimize their inference setup accordingly.

## Tutorial: Multi-LoRA serving for Llama-3.1-8B on Trn2 instances

NeuronX Distributed (NxD) Inference supports multi-LoRA serving. This tutorial provides a step-by-step guide for multi-LoRA serving with Llama-3.1-8B as the base model on a Trn2 instance. It describes two different ways of running multi-LoRA serving with NxD Inference directly and through vLLM (with NxD Inference). We will use LoRA adapters downloaded from HuggingFace as examples for serving.

### Table of contents

- *Prerequisites:*
  - *Set up and connect to a Trn2.48xlarge instance*
  - *Install packages*
  - *Download base model and LoRA adapters*
- *Run multi-LoRA serving on Trn2 from NxD Inference*
- *Using vLLM for multi-LoRA serving on Trn2*
  - *Multi-LoRA Configurations*
  - *Offline inference example*
  - *Online server example*

## Prerequisites:

### Set up and connect to a Trn2.48xlarge instance

As a prerequisite, this tutorial requires that you have a Trn2 instance created from a Deep Learning AMI that has the Neuron SDK pre-installed.

To set up a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK, see [Nx\*D\* Inference Setup Guide](#).

After setting up an instance, use SSH to connect to the Trn2 instance using the key pair that you chose when you launched the instance.

After you are connected, activate the Python virtual environment that includes the Neuron SDK.

```
source ~/aws_neuronx_venv_pytorch_2_5_nxd_inference/bin/activate
```

Run `pip list` to verify that the Neuron SDK is installed.

```
pip list | grep neuron
```

You should see Neuron packages including `neuronx-distributed-inference` and `neuronx-cc`.

## Install packages

NxD Inference supports running models with vLLM. This functionality is available in the AWS Neuron fork of the vLLM GitHub repository. Install the latest release branch of vLLM from the AWS Neuron fork following instructions in the [vLLM User Guide for Nx\*D\* Inference](#).

## Download base model and LoRA adapters

To use this sample, you must first download a [Llama-3.1-8B-Instruct](#) model checkpoint from Hugging Face to a local path on the Trn2 instance. Note that you may need access from Meta for model download. For more information, see [Downloading models](#) in the Hugging Face documentation.

You also need to download LoRA adapters from Hugging Face for multi-LoRA serving. As examples, you can download [nvidia/llama-3.1-nemoguard-8b-topic-control](#) and [reissbaker/llama-3.1-8b-abliterated-lora](#).

## Run multi-LoRA serving on Trn2 from Nx*D* Inference

We will run multi-LoRA serving from Nx*D* inference with `inference_demo` on Trn2 using Llama-3.1-8B and two LoRA adapters. The data type is bfloat16 precision.

You should specifically set the following configurations when enabling multi-LoRA serving with `inference_demo`.

- `enable_lora` - The flag to enable multi-LoRA serving in Nx*D* Inference. Defaults to *False*.
- `max_loras` - The maximum number of concurrent LoRA adapters in device memory. Defaults to 1.
- `max_lora_rank` - The highest LoRA rank that needs to be supported. Defaults to 16. If it is not specified, the maximum LoRA rank of the LoRA adapter checkpoints will be used.
- `lora_ckpt_path` - The checkpoint path for LoRA adapter in the format of `adapter_id : path`. Please set this flag multiple times if multiple LoRA adapters are needed.
- `adapter_id` - The adapter ID for prompt. Each prompt comes with an adapter ID.

Save the contents of the below script to your favorite shell script file, for example, `multi_lora_model.sh` and then run it. The script compiles the model and runs generation on the given input prompt.

```
# Replace this with the path where you downloaded and saved the model files.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-8B-Instruct/"
# Replace the following with the paths where you downloaded and saved the LoRA adapters.
LORA_PATH_1="/home/ubuntu/models/loras/llama-3.1-nemoguard-8b-topic-control"
LORA_PATH_2="/home/ubuntu/models/loras/llama-3.1-8b-abliterated-lora"
# This is where the compiled model will be saved.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-8B-Lora/"

NUM_CORES=128
TP_DEGREE=32
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$TP_DEGREE
export NEURON_RT_EXEC_TIMEOUT=600
export XLA_DENSE_GATHER_FACTOR=0
export NEURON_RT_INSPECT_ENABLE=0

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
  --model-path $MODEL_PATH \
  --compiled-model-path $COMPILED_MODEL_PATH \
  --torch-dtype bfloat16 \
  --start_rank_id 0 \
  --local_ranks_size $TP_DEGREE \
  --tp-degree $TP_DEGREE \
  --batch-size 2 \
  --max-context-length 12288 \
  --seq-len 64 \
  --on-device-sampling \
  --top-k 1 \
  --do-sample \
  --pad-token-id 2 \
  --enable-bucketing \
  --enable-lora \
  --max-loras 2 \
  --lora-ckpt-path "lora_id_1 : ${LORA_PATH_1}" \
  --lora-ckpt-path "lora_id_2 : ${LORA_PATH_2}" \
  --prompt "I believe the meaning of life is" \
  --adapter-id lora_id_1 \
  --prompt "I believe the meaning of life is" \
  --adapter-id lora_id_2 \
  | tee log
```

NxD Inference expects the same number of prompts and adapter IDs in the script. A prompt is mapped to the adapter ID with the same order. For example, the first prompt in the script associates with `lora_id_1` and the second one associates with `lora_id_2`. Although the two prompts are the same, NxD Inference will generate different outputs due to different adapter IDs.

## Using vLLM for multi-LoRA serving on Trn2

We can run multi-LoRA serving on Trn2 with vLLM for Llama models. Please refer to *vLLM User Guide for NxD Inference* for more details on how to run model inference on TRN2 with vLLM.

### Multi-LoRA Configurations

You should specifically set the following configurations when enabling multi-LoRA serving with vLLM.

- `enable_lora` - The flag to enable multi-LoRA serving in NxD Inference. Defaults to `False`.
- `max_loras` - The maximum number of concurrent LoRA adapters in device memory. Defaults to 1.
- `max_lora_rank` - The highest LoRA rank that needs to be supported. Defaults to 16. If it is not specified, the maximum LoRA rank of the LoRA adapter checkpoints will be used.
- `lora_modules` - Set the LoRA checkpoint paths and their adapter IDs in the format of `adapter_id_1=path1 adapter_id_2=path2 ....`

### Offline inference example

You can also run multi-LoRA serving offline on TRN2 with vLLM.

```
import os
os.environ['VLLM_NEURON_FRAMEWORK'] = "neuronx-distributed-inference"
from vllm import LLM, SamplingParams
from vllm.entrypoints.openai.serving_models import LoRAModulePath
from vllm.lora.request import LoRARequest

MODEL_PATH="/home/ubuntu/models/Llama-3.1-8B-Instruct/"
# LoRA checkpoint paths.
LORA_PATH_1="/home/ubuntu/models/loras/llama-3.1-nemoguard-8b-topic-control"
LORA_PATH_2="/home/ubuntu/models/loras/llama-3.1-8b-abliterated-lora"

# Sample prompts.
prompts = [
    "The president of the United States is",
    "The capital of France is",
]

# Create a sampling params object.
sampling_params = SamplingParams(top_k=1)

# Create an LLM with multi-LoRA serving.
llm = LLM(
    model=MODEL_PATH,
    max_num_seqs=2,
    max_model_len=64,
    tensor_parallel_size=32,
    device="neuron",
    override_neuron_config={
        "sequence_parallel_enabled": False,
    },
)
```

(continues on next page)

(continued from previous page)

```

lora_modules=[
    LoRAModulePath(name="lora_id_1", path=LORA_PATH_1),
    LoRAModulePath(name="lora_id_2", path=LORA_PATH_2),
],
enable_lora=True,
max_loras=2,
)
"""
NxN Inference enables static loading of LoRA adapters: https://docs.vllm.ai/en/v0.9.0/
↳ features/lora.html on vLLM server start and does
not optionally support dynamic serving of LoRA adapters: https://docs.vllm.ai/en/v0.9.0/
↳ features/lora.html#dynamically-serving-lora-adapters
Only the lora_name needs to be specified.
The lora_id and lora_path are supplied at the LLM class/server initialization, after
↳ which the paths are
handled by NxN Inference.
"""
lora_req_1 = LoRARequest("lora_id_1", 0, " ")
lora_req_2 = LoRARequest("lora_id_2", 1, " ")
outputs = llm.generate(prompts, sampling_params, lora_request=[lora_req_1, lora_req_2])

for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

```

### Online server example

Save the contents of the below script to another shell script file, for example, `start_vllm.sh` and then run it.

```

export NEURON_RT_INSPECT_ENABLE=0
export NEURON_RT_VIRTUAL_CORE_SIZE=2

# These should be the same paths used when compiling the model.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-8B-Instruct/"
# Replace the following with the paths where you downloaded and saved the LoRA adapters.
LORA_PATH_1="/home/ubuntu/models/loras/llama-3.1-nemoguard-8b-topic-control"
LORA_PATH_2="/home/ubuntu/models/loras/llama-3.1-8b-abliterated-lora"
# This is where the compiled model will be saved.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-8B-Lora/"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs 2 \
    --max-model-len 64 \
    --tensor-parallel-size 32 \
    --device neuron \
    --use-v2-block-manager \
    --enable-lora \

```

(continues on next page)



(continued from previous page)

```
--max-loras 2 \
--override-neuron-config "{\"sequence_parallel_enabled\": false}" \
--lora-modules lora_id_1=${LORA_PATH_1} lora_id_2=${LORA_PATH_2} \
--port 8000 &
PID=$!
echo "vLLM server started with PID $PID"
```

After the vLLM server is launched, we can send requests to the server for serving. A sample request is:

```
curl http://localhost:8000/v1/chat/completions -H "Content-Type: application/json" -
-d '{
  "model": "lora_id_1",
  "messages": [
    {
      "role": "user",
      "content": "The president of the United States is"
    }
  ]
}'
```

## Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances

NeuronX Distributed (NxD) Inference allows you to deploy Llama3.1 405B on a single Trn2 instance. This tutorial will show you how to optimize inference performance for Llama3.1 405B on a Trn2 instance with speculative decoding and quantization. We will compile and load the model into a VLLM server and measure performance using LLMPPerf. This tutorial consists of two parts. In the first part, we will collect performance metrics for our base configuration with bf16 model weights. In the second part, we will optimize inference performance with fp8 quantized weights and speculative decoding. The performance is then compared with the results from part 1.

### Table of contents

- *Prerequisites*
  - *Set up and connect to a Trn2.48xlarge instance*
  - *Install packages*
  - *Download models*
- *Scenario 1: Run Llama-3.1-405b inference with base configuration using bf16 weights*
  - *Step 1: Compile the model and run generate*
  - *Step 2: Start the Vllm server with the compiled Neuron model*
  - *Step 3: Measure performance using LLMPPerf*
- *Scenario 2: Run Llama-3.1-405b inference with fp8 weights and fused speculation (with draft model)*
  - *Step 1: Rescale the model weights to use Neuron FP8 format and save the modules to not convert file in model path*
  - *Step 2: Compile the model and run generate*
  - *Step 3: Start the Vllm server with the compiled Neuron model*

– *Step 4: Measure performance using LLMPerf*

- *Conclusion*

## Prerequisites

### Set up and connect to a Trn2.48xlarge instance

As a prerequisite, this tutorial requires that you have a Trn2 instance created from a Deep Learning AMI that has the Neuron SDK pre-installed.

To set up a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK, see [Nx\*D\* Inference Setup Guide](#).

After setting up an instance, use SSH to connect to the Trn2 instance using the key pair that you chose when you launched the instance.

After you are connected, activate the Python virtual environment that includes the Neuron SDK.

```
source ~/aws_neuronx_venv_pytorch_2_5_nxd_inference/bin/activate
```

Run `pip list` to verify that the Neuron SDK is installed.

```
pip list | grep neuron
```

You should see Neuron packages including `neuronx-distributed-inference` and `neuronx-cc`.

## Install packages

NxD Inference supports running models with vLLM. This functionality is available in the AWS Neuron fork of the vLLM GitHub repository. Install the latest release branch of vLLM from the AWS Neuron fork following instructions in the [vLLM User Guide for Nx\*D\* Inference](#).

In this tutorial, you will use `llmperf` to measure the inference performance of the base Llama-3.1-405b-Instruct configuration and the more optimized configuration. We will use the `load test` feature of LLMPerf and measure the performance for accepting 10,000 tokens as input and generating 1500 tokens as output. Install `llmperf` into the virtual environment.

```
git clone https://github.com/ray-project/llmperf.git
cd llmperf
pip install -e .
```

## Download models

To run inference in the first part of the tutorial, you need to download the Llama-3.1-405b-Instruct model checkpoint with bf16 weights from Hugging Face ([meta-llama/Llama-3.1-405B-Instruct](#)). For the second part of the tutorial, you will run a more optimized inference configuration. For this part, you need to download an fp8-quantized Llama3.1-405B-FP8 model checkpoint ([meta-llama/Llama-3.1-405B-Instruct-FP8](#)). With Speculative Decoding, you will also need to specify a draft model. You can download and use the model checkpoint from [meta-llama/Llama-3.2-1B-Instruct](#). For more information, see [Downloading models](#) in the Hugging Face documentation.

## Scenario 1: Run Llama-3.1-405b inference with base configuration using bf16 weights

### Step 1: Compile the model and run generate

We will first compile and run generation on a sample prompt using a command installed by neuronx-distributed-inference. Save the contents of the below script to your favorite shell script file, for example, `compile_model.sh` and then run it.

Note that we are using the following features as described in the tutorial for running 405B model [Tutorial: Deploying Llama3.1 405B \(Trn2\)](#)

- Logical NeuronCore Configuration (LNC)
- Tensor parallelism (TP) on Trn2
- Optimized Kernels

The script compiles the model and runs generation on the given input prompt. Please refer to [Nx D Inference API Reference](#) for more information on these `inference_demo` flags. Note the path we used to save the compiled model. This path should be used when launching vLLM server for inference so that the compiled model can be loaded without recompilation.

**Note:** Known issue: Using kernels with bucket length of 1024 or less may lead to Numerical Error in inference.

```
RuntimeError: Failed to execute the model status=1003 message=Numerical Error
```

```
# Replace this with the path where you downloaded and saved the model files.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct/"
# This is where the compiled model will be saved. The same path
# should be used when launching vLLM server for inference.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-405B-Instruct/"

NUM_CORES=128
TP_DEGREE=64
LNC=2

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
  --model-path $MODEL_PATH \
  --compiled-model-path $COMPILED_MODEL_PATH \
  --torch-dtype bfloat16 \
  --start_rank_id 0 \
  --local_ranks_size $TP_DEGREE \
  --tp-degree $TP_DEGREE \
  --batch-size 1 \
  --max-context-length 12288 \
```

(continues on next page)

(continued from previous page)

```

--seq-len 12800 \
--on-device-sampling \
--top-k 1 \
--fused-qkv \
--sequence-parallel-enabled \
--qkv-kernel-enabled \
--attn-kernel-enabled \
--mlp-kernel-enabled \
--cc-pipeline-tiling-factor 1 \
--pad-token-id 2 \
--enable-bucketing \
--context-encoding-buckets 2048 4096 10240 12288 \
--token-generation-buckets 12800 \
--prompt "What is annapurna labs?" 2>&1 | tee log

```

The above script will compile a Neuron model for this base-case configuration, and also run generate on the example prompt specified with the `-prompt` flag. You can change this prompt to your prompt of choice. The script's output will be written into `log`, a log file in the working directory.

In addition, in the subsequent runs of this script, you can add a `--skip-compile` flag to skip the compiling step since the model is already compiled in the first run of the script. This will allow you to test the model with different prompts.

## Step 2: Start the VLLM server with the compiled Neuron model

After compiling the model, you can run the model using vLLM. Save the contents of the below script to another shell script file, for example, `start_vllm.sh` and then run it.

```

export NEURON_RT_VIRTUAL_CORE_SIZE=2

MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct"
COMPILED_MODEL_PATH="/home/ubuntu/traced_models/Llama-3.1-405B-Instruct"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs 1 \
    --max-model-len 12800 \
    --tensor-parallel-size 64 \
    --device neuron \
    --use-v2-block-manager \
    --override-neuron-config "{}" \
    --port 8000 & PID=$!
echo "vLLM server started with PID $PID"

```

### Step 3: Measure performance using LLMPerf

After the above steps, the vllm server should be running. Before we can use the llmperf package, we need to make a few changes to its code. Follow [benchmarking with LLMPerf guide](#) to apply the code changes.

We can now measure the performance using llmperf. Below is a sample shell script to run llmperf. More information about several arguments used in the script can be found in the [llmperf open source code](#).

```
# This should be the same path to which the model was downloaded (also used in the above
steps).
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct"
# This is the name of directory where the test results will be saved.
OUTPUT_PATH=llmperf-results-sonnets

export OPENAI_API_BASE="http://localhost:8000/v1"
export OPENAI_API_KEY="mock_key"

python token_benchmark_ray.py \
    --model $MODEL_PATH \
    --mean-input-tokens 10000 \
    --stddev-input-tokens 0 \
    --mean-output-tokens 1500 \
    --stddev-output-tokens 0 \
    --num-concurrent-requests 1 \
    --timeout 3600 \
    --max-num-completed-requests 50 \
    --additional-sampling-params '{}' \
    --results-dir $OUTPUT_PATH \
    --llm-api "openai"
```

The output for this llama-3.1-405B model run for the base case is shown below. Please note that the numbers can slightly vary between runs but should be in the same order of magnitude.

Results for token benchmark for /home/ubuntu/models/llama-3.1-405b queried with the openai api.

```
inter_token_latency_s
p25 = 0.03783673520494379
p50 = 0.037929154633788834
p75 = 0.03799374728198055
p90 = 0.03806084386428147
p95 = 0.03818095359194858
p99 = 0.03862880035825585
mean = 0.03790912092492011
min = 0.03711292916794487
max = 0.03867580939426865
stddev = 0.0002364662521116205
ttft_s
p25 = 2.437347081664484
p50 = 2.441959390998818
p75 = 2.4439403364085592
p90 = 2.444729209714569
p95 = 2.445114637189545
p99 = 79.22927707570342
```

(continues on next page)

(continued from previous page)

```

    mean = 5.451600373298861
    min = 2.427013176959008
    max = 153.00210832804441
    stddev = 21.29264628138615
end_to_end_latency_s
    p25 = 70.06310007086722
    p50 = 70.09642704750877
    p75 = 70.1557097924524
    p90 = 70.28295350184199
    p95 = 70.56055794338462
    p99 = 148.28325726192182
    mean = 73.19207735829521
    min = 70.00512732309289
    max = 222.50397142698057
    stddev = 21.54750467688136
request_output_throughput_token_per_s
    p25 = 25.417755028050912
    p50 = 25.463487985775544
    p75 = 25.522234144656743
    p90 = 25.6487981126861
    p95 = 25.729858763245502
    p99 = 25.90146713883131
    mean = 25.13808905954906
    min = 8.080754642125802
    max = 26.021214285642255
    stddev = 2.465472136291901
number_input_tokens
    p25 = 10000.0
    p50 = 10000.0
    p75 = 10000.0
    p90 = 10000.0
    p95 = 10000.0
    p99 = 10000.0
    mean = 10000.0
    min = 10000
    max = 10000
    stddev = 0.0
number_output_tokens
    p25 = 1783.0
    p50 = 1785.0
    p75 = 1789.75
    p90 = 1798.1
    p95 = 1803.55
    p99 = 1816.67
    mean = 1787.92
    min = 1779
    max = 1825
    stddev = 8.54720386310933
Number Of Errored Requests: 0
Overall Output Throughput: 24.421011092151268
Number Of Completed Requests: 50
Completed Requests Per Minute: 0.8195336846889548

```

## Scenario 2: Run Llama-3.1-405b inference with fp8 weights and fused speculation (with draft model)

### Step 1: Rescale the model weights to use Neuron FP8 format and save the modules to not convert file in model path

Since Neuron device only supports the FP8\_EXP4 (IEEE-754) data type, and the HuggingFace FP8 checkpoint for Llama-405b is in a different FP8 format (OCP FP8 E4M3/e4m3fn) which has a different range, we need to rescale the public model weights. Follow this guide to rescale the FP8 model weights from HuggingFace: [link](#).

Running a quantized model requires us to create modules to not convert json file to explicitly mention the layers which are not quantized in the model. For this tutorial we can use the following file.

Download: `modules_to_not_convert.json`

Next we will compile and run the model and record performance metrics.

### Step 2: Compile the model and run generate

We will first compile and run generation on a sample prompt using a command installed by `neuronx-distributed-inference`. Save the contents of the below script to your favorite shell script file, for example, `compile_model.sh` and then run it.

Note that we are using the following features as described in the tutorial for running 405B model *Tutorial: Deploying Llama3.1 405B (Trn2)*

- Logical NeuronCore Configuration (LNC)
- Tensor parallelism (TP) on Trn2
- Optimized Kernels

The compiling script is similar to the one in part 1. Note that we have added the path for the draft model.

**Note:** Known issue: Using kernels with bucket length of 1024 or less may lead to Numerical Error in inference.

**RuntimeError:** Failed to execute the model status=1003 message=Numerical Error

```
# Replace this with the path where you downloaded and saved the model files.
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct-FP8-rescaled/"
# Replace this with the path where you downloaded and saved the draft model files.
DRAFT_MODEL_PATH="/home/ubuntu/models/Llama-3.2-1b-instruct/"
# This is where the compiled model (.pt file) and sharded checkpoints will be saved. The
↪ same path
# should be used when launching vLLM server for inference.
COMPILED_MODEL_PATH="/home/ubuntu/traced_model/Llama-3.1-405B-Instruct/"
# Add a modules to not convert json file to the model path to specify non quantized
↪ modules.
MTNC_FILE_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct-FP8-rescaled/modules_to_not_
↪ convert.json"

NUM_CORES=128
TP_DEGREE=64
LNC=2
```

(continues on next page)

(continued from previous page)

```

export NEURON_RT_VIRTUAL_CORE_SIZE=$LNC
export NEURON_RT_NUM_CORES=$((NUM_CORES/NEURON_RT_VIRTUAL_CORE_SIZE))
export NEURON_RT_EXEC_TIMEOUT=600
export XLA_HANDLE_SPECIAL_SCALAR=1
export UNSAFE_FP8FNCAST=1

inference_demo \
  --model-type llama \
  --task-type causal-lm \
  run \
    --model-path $MODEL_PATH \
    --compiled-model-path $COMPILED_MODEL_PATH \
    --torch-dtype bfloat16 \
    --start_rank_id 0 \
    --local_ranks_size $TP_DEGREE \
    --tp-degree $TP_DEGREE \
    --batch-size 1 \
    --max-context-length 12288 \
    --seq-len 12800 \
    --on-device-sampling \
    --top-k 1 \
    --fused-qkv \
    --sequence-parallel-enabled \
    --qkv-kernel-enabled \
    --attn-kernel-enabled \
    --mlp-kernel-enabled \
    --cc-pipeline-tiling-factor 1 \
    --draft-model-path $DRAFT_MODEL_PATH \
    --enable-fused-speculation \
    --speculation-length 7 \
    --pad-token-id 2 \
    --quantized-mlp-kernel-enabled \
    --quantization-type per_channel_symmetric \
    --rmsnorm-quantize-kernel-enabled \
    --enable-bucketing \
    --prompt "What is annapurna labs?" \
    --modules-to-not-convert-file $MTNC_FILE_PATH \
    --context-encoding-buckets 2048 4096 10240 12288 \
    --token-generation-buckets 12800 2>&1 | tee compile_and_generate_log

```

The above script will compile a Neuron model with fused speculation, and also run generate on the example prompt specified with the `-prompt` flag. Please refer to [Nx\*D\* Inference API Reference](#) for more information on these `inference_demo` flags.

You can change this prompt to your prompt of choice. The script's output will be written into `compile_and_generate_log`, a log file in the working directory.

In this script, we also turn on some additional environment variables: `XLA_HANDLE_SPECIAL_SCALAR` and `UNSAFE_FP8FNCAST` to enable Neuron compiler to treat rescaled FP8FN weights as FP8\_EXP4 weights.

In addition, in the subsequent runs of this script, you can add a `--skip-compile` flag to skip the compiling step since the model is already compiled in the first run of the script. This will allow you to test the model with different prompts.



### Step 3: Start the Vllm server with the compiled Neuron model

After compiling the model, you can run the model using vLLM. Save the contents of the below script to another shell script file, for example, `start_vllm.sh` and then run it.

```
export NEURON_RT_INSPECT_ENABLE=0
export NEURON_RT_VIRTUAL_CORE_SIZE=2
export XLA_HANDLE_SPECIAL_SCALAR=1
export UNSAFE_FP8FNCAST=1

MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct-FP8-rescaled"
DRAFT_MODEL_PATH="/home/ubuntu/models/Llama-3.2-1b-instruct"
COMPILED_MODEL_PATH="/home/ubuntu/traced_models/Llama-3.1-405B-Instruct_fp8"

export VLLM_NEURON_FRAMEWORK="neuronx-distributed-inference"
export NEURON_COMPILED_ARTIFACTS=$COMPILED_MODEL_PATH
VLLM_RPC_TIMEOUT=1000000 python -m vllm.entrypoints.openai.api_server \
    --model $MODEL_PATH \
    --max-num-seqs 1 \
    --max-model-len 12800 \
    --tensor-parallel-size 64 \
    --device neuron \
    --speculative-max-model-len 12800 \
    --speculative-model $DRAFT_MODEL_PATH \
    --num-speculative-tokens 7 \
    --use-v2-block-manager \
    --override-neuron-config '{"enable_fused_speculation":true, "quantized-mlp-kernel-
    enabled":true, "quantization-type":"per_channel_symmetric", "skip_warmup": true}
    ' \
    --port 8000 & PID=$!
echo "vLLM server started with PID $PID"
```

### Step 4: Measure performance using LLMPerf

After the above steps, the vllm server should be running. Before we can use the `llmperf` package, we need to make a few changes to its code. Follow [benchmarking with LLMPerf guide](#) to apply the code changes.

We can now measure the performance using `llmperf`. Run the following script with the modified `llmperf` package.

```
# This should be the same path to which the model was downloaded (also used in the above
steps).
MODEL_PATH="/home/ubuntu/models/Llama-3.1-405B-Instruct-FP8-rescaled"
# This is the name of directory where the test results will be saved.
OUTPUT_PATH=llmperf-results-sonnets

export OPENAI_API_BASE="http://localhost:8000/v1"
export OPENAI_API_KEY="mock_key"

python token_benchmark_ray.py \
    --model $MODEL_PATH \
```

(continues on next page)

(continued from previous page)

```

--mean-input-tokens 10000 \
--stddev-input-tokens 0 \
--mean-output-tokens 1500 \
--stddev-output-tokens 0 \
--num-concurrent-requests 1\
--timeout 3600 \
--max-num-completed-requests 50 \
--additional-sampling-params '{}' \
--results-dir $OUTPUT_PATH \
--llm-api "openai"

```

The output for this llama-3.1-405B model run with fused speculation with fused spec is shown below. Please note that the numbers can slightly vary between runs but should be in the same order of magnitude.

Results **for** token benchmark **for** /home/ubuntu/models/Llama-3.1-405B-Instruct-FP8-rescaled\_ queried **with** the openai api.

```

inter_token_latency_s
p25 = 0.008220573497974934
p50 = 0.008265312568750231
p75 = 0.008438719224417583
p90 = 0.00848199803312309
p95 = 0.008495625438929224
p99 = 0.011143428944987235
mean = 0.008419798457414533
min = 0.008173695931987216
max = 0.01364151847269386
stddev = 0.0007612118573477839
ttft_s
p25 = 2.2543624382815324
p50 = 2.254961202503182
p75 = 2.2576071268413216
p90 = 2.2596270388457924
p95 = 2.260639927221928
p99 = 2.2628143909573555
mean = 2.256157155628316
min = 2.2534945809748024
max = 2.2629711360204965
stddev = 0.0023667267664955545
end_to_end_latency_s
p25 = 14.586015026085079
p50 = 14.65608573507052
p75 = 14.91364526405232
p90 = 14.977840351965279
p95 = 15.000083449739032
p99 = 18.969864878777866
mean = 14.886235136194154
min = 14.520539953839034
max = 22.716861865017563
stddev = 1.1415236552464672
request_output_throughput_token_per_s
p25 = 100.64608830743339

```

(continues on next page)

(continued from previous page)

```

p50 = 102.4148205461138
p75 = 102.90679421801005
p90 = 103.02201242683091
p95 = 103.26614794565539
p99 = 103.36118277211666
mean = 101.22055373532301
min = 66.0742671641385
max = 103.37081160698546
stddev = 5.19249551094185
number_input_tokens
p25 = 10000.0
p50 = 10000.0
p75 = 10000.0
p90 = 10000.0
p95 = 10000.0
p99 = 10000.0
mean = 10000.0
min = 10000
max = 10000
stddev = 0.0
number_output_tokens
p25 = 1501.0
p50 = 1501.0
p75 = 1501.0
p90 = 1501.0
p95 = 1501.0
p99 = 1501.0
mean = 1501.0
min = 1501
max = 1501
stddev = 0.0
Number Of Errored Requests: 0
Overall Output Throughput: 100.69986490153724
Number Of Completed Requests: 50
Completed Requests Per Minute: 4.025311055357918

```

## Conclusion

As seen from the table below, draft model based fused speculative decoding and quantization significantly improved inference performance: TPOT reduced by 4x and output token throughput increased by 4x, while TTFT decreased from 2442 ms to 2255 ms compared to baseline without speculative decoding. Please note that batch size of 1 is used in this tutorial for computing the below metrics.

| Scenario (all using BF16)   | TTFT (P50<br>in ms) | TPOT (P50<br>in ms) | Output token Throughput<br>(per second) |
|---|---------------------|---------------------|---|
| No speculative decoding   | 2442                | 37.9                | 25.46                                   |
| Fused speculative decoding + rescaled weights<br>(Llama 3.2 1B Draft) | 2255                | 8.27                | 102.41                                  |

## Tutorial: Evaluating Accuracy of Llama-3.1-70B on Neuron using open source datasets

### Introduction

This tutorial provides a step-by-step guide to measure the accuracy of Llama3.1 70B on Trn1 with evaluation on two distinct tasks: mathematical reasoning and logical analysis.

For this tutorial we use two datasets available in lm-eval, namely `gsm8k_cot` (high school math questions) and `mmlu_flan_n_shot_generative_logical_fallacies` (multiple choice questions on the subject) to demonstrate accuracy evaluation on Trn1. The metrics in these task are two variants of `ExactMatch` metrics called `StrictMatch` and `FlexibleExtract` which differ in how strict they are in extracting the final answer from the generated output from the model. To see the exact task definition used in lm-eval please look at [gsm8k-cot](#) and [mmlu template](#).

We also need the instruction-tuned version of llama-3.1 70b [meta-llama/Llama-3.1-70B-Instruct](#) available hugging face.

### Task Overview

#### 1. GSM8K with Chain-of-Thought (`gsm8k_cot`)

The GSM8K dataset focuses on grade school math word problems, testing LLMs' mathematical reasoning capabilities. Using Chain-of-Thought (CoT) prompting, we evaluate models' ability to:

- Solve complex math word problems
- Show step-by-step reasoning
- Arrive at accurate numerical answers

#### 2. MMLU Logical Fallacies (`mmlu_flan_n_shot_generative_logical_fallacies`)

This evaluation focuses on the model's ability to identify and explain logical fallacies, a subset of the MMLU benchmark. The task tests:

- Understanding of common logical fallacies
- Ability to analyze arguments
- Explanation of reasoning flaws

### Environment Setup Guide

#### Prerequisites

This tutorial requires that you have a Trn1 instance created from a Deep Learning AMI that has the Neuron SDK pre-installed. Also we depend on our fork of vLLM as described in the [vLLM User Guide for NxN Inference](#).

Before running evaluations, ensure your environment is properly configured by following these essential setup guides:

1. [NxN Inference Setup Guide](#)

- Configure AWS Neuron environment
- Set up required dependencies
- Verify system requirements

## 2. vLLM User Guide for NxD Inference

- Setup vLLM according to the guide

### Installing dependencies

Copy the `inference-benchmarking` directory to some location on your instance. Change directory to the your copy of `inference-benchmarking`. Install other required dependencies in the same python env (e.g `aws_neuron_venv_pytorch` if you followed [manual install NxD Inference](#) ) by:

```
!pip install -r requirements.txt
```

### Download llama-3.1 70B

To use this sample, you must first download `meta-llama/Llama-3.1-70B-Instruct` model checkpoint from Hugging Face `/home/ubuntu/models/Llama-3.1-70B-Instruct/` on the Trn1 instance. For more information, see [Downloading models](#) in the Hugging Face documentation.

### Running Evaluations

There are two methods that you can use [the evaluation scripts](#) to run your evaluation.

1. Using a yaml configuration file and `accuracy.py` script
2. writing your own python script that uses several components provided in `accuracy.py` and `server_config.py`

We demonstrate each use case separately here.

#### 1. Running eval with yaml config file

In this method all you need is to create a yaml config file that specifies the server configuration and testing scenario you want to run. Create `config.yaml` with the following content.

```
server:
  name: "Llama-3.1-70B-Instruct"
  model_path: "/home/ubuntu/models/Llama-3.1-70B-Instruct/"
  model_s3_path: null
  compiled_model_path: "/home/ubuntu/traced_models/Llama-3.1-70B-Instruct"
  max_seq_len: 16384
  context_encoding_len: 16384
  tp_degree: 32
  n_vllm_threads: 32
  server_port: 8000
  continuous_batch_size: 1

test:
  accuracy:
    mytest:
      client: "lm_eval"
      datasets: ["gsm8k_cot", "mmlu_flan_n_shot_generative_logical_fallacies"]
      max_concurrent_requests: 1
```

(continues on next page)

(continued from previous page)

```

timeout: 3600
client_params:
  limit: 200
  use_chat: True

```

For tasks that require higher sequence length you need to adjust `max_seq_len`. For the tasks in this tutorial 16384 would suffice.

Run `python accuracy.py --config config.yaml`

## 2. Running eval through your own python code

You might be interested in running the evaluation in you python code. For instance if you want to change the configuration programatically or post-process the results. This is possible using 3 main components provided in `accuracy.py` and `server_config.py`.

1. Server Configuration: Using `ServerConfig` to define the vLLM server settings
2. Accuracy Scenario: Using `AccuracyScenario` to specify evaluation parameters
3. Test Execution: Running the evaluation with the configured settings

## Step-by-Step Implementation

First, import the necessary components:

```

from accuracy import AccuracyScenario, run_accuracy_test
from server_config import ServerConfig

```

### 1. Configure the Server

Set up your server configuration with `ServerConfig`. This example uses Llama 3.1-8b Instruct:

```

name = "Llama-3.1-70B-Instruct"
server_config = ServerConfig(
    name=name,
    model_path=f"/home/ubuntu/models/{name}", # Local model path
    model_s3_path=None, # S3 model path
    max_seq_len=16384, # Maximum sequence length
    context_encoding_len=16384, # Context window size
    tp_degree=32, # Tensor parallel degree
    n_vllm_threads=32, # Number of vLLM threads
    server_port=8000, # Server port
    continuous_batch_size=1, # Batch size for continuous batching
)

```

## 2. Define the Evaluation Scenario

Create an AccuracyScenario to specify your evaluation parameters:

```
scenario = AccuracyScenario(
    client="lm_eval",          # Evaluation client
    datasets=[                # Target datasets
        "gsm8k_cot",
        "mmlu_flan_n_shot_generative_logical_fallacies",
    ],
    max_concurrent_requests=1, # Maximum concurrent requests
    timeout=3600,              # Timeout in seconds
    client_params={"limit": 200} # Client-specific parameters
)
```

## 3. Run the Evaluation

Execute the evaluation using run\_accuracy\_test:

```
# Run the test with a named scenario
results_collection = run_accuracy_test(
    server_config=server_config,
    named_scenarios={"mytest": scenario}
)

# Display results
print(results_collection)
```

This code will execute the evaluation on the specified datasets and return detailed performance metrics. The results include accuracy scores and other relevant metrics for each dataset.

## Tutorial: Disaggregated Inference on Trn2 [BETA]

### Overview

This tutorial will mainly cover how to run Disaggregated Inference (DI) 1P1D (1 prefill, 1 Decode) either on a single Trn2 instance (1P and 1D both are on same instance) or on 2 instances (1P and 1D are on separate instances). It will provide scripts that can setup both single and multi instance workflows. Next, the tutorial will demonstrate how to benchmark DI. Finally, we show how to benchmark non Disaggregated Inference (non-DI) continuous batching to compare results between DI vs. non-DI.

Read the [DI Developer Guide](#) for more detailed information.

---

**Note:** This tutorial was tested on trn2.48xlarge but its concepts are also be applicable to trn1.32xlarge.

---

## Set up and connect trn2.48xlarge instance

As a prerequisite, this tutorial requires that you have one or two Trn2 instances with Neuron SDK, Neuron vLLM and Elastic Fabric Adapter (EFA) enabled and installed. The Neuron Deep Learning AMI comes with Neuron dependencies and EFA enabled and installed so it is the recommended way to run this tutorial.

To set up a Trn2 instance using Deep Learning AMI with pre-installed Neuron SDK, see [Nx D Inference Setup Guide](#).

---

**Note:** Disaggregated Inference is only supported on Neuron instances with EFA enabled (trn1.32xlarge or trn2.48xlarge). EFA is still required even when running single instance as the KV cache transfer happens through EFA.

---

If you choose to manually install Nx D Inference follow the [EFA setup guide](#) to install and enable EFA.

If running multi-instance it is recommended to have shared storage between the two instances to avoid having to download, compile and save scripts twice. For more details, see documentation on mounting [EFS](#) or [FSX](#) filesystems.

After setting up an instance, use SSH to connect to the Neuron instance(s) using the key pair that you chose when you launched the instance.

After you are connected, activate the Python virtual environment that includes the Neuron SDK.

```
source ~/aws_neuronx_venv_pytorch_2_7_nxd_inference/bin/activate
```

Install the Neuron vLLM fork into the virtual environment(s).

```
git clone -b neuron-2.24-vllm-v0.7.2 https://github.com/aws-neuron/upstreaming-to-vllm.  
↪ git  
cd upstreaming-to-vllm  
pip install -r requirements-neuron.txt  
VLLM_TARGET_DEVICE="neuron" pip install -e .  
cd ..
```

For more information see [vLLM User Guide for Nx D Inference](#).

Run `pip list` to verify that the Neuron SDK is installed.

```
pip list | grep neuron
```

You should see Neuron packages including `neuronx-distributed-inference` and `neuronx-cc` and `vllm`.

## Download Dependencies

To use this sample, you must first download a [Llama-3.3-70B-Instruct](#) model checkpoint from Hugging Face to a local path on the Trn2 instance. Note that you may need access from Meta for model download. For more information, see [Downloading models](#) in the Hugging Face documentation.



## Compile the model

Compile the model for Neuron by using the following `compile.sh` script.

```
#!/bin/bash
# copy and paste me into a file called compile.sh
# then run chmod +x compile.sh

# Parse command line arguments
while [[ $# -gt 0 ]]; do
    case $1 in
        --tp-degree)
            TP_DEGREE="$2"
            shift 2
            ;;
        --batch-size)
            BATCH_SIZE="$2"
            shift 2
            ;;
        --model-path)
            MODEL_PATH="$2"
            shift 2
            ;;
        *)
            echo "Unknown parameter: $1"
            echo "Usage: $0 --tp-degree <value> --batch-size <value> --model-path <path>"
            exit 1
            ;;
    esac
done

export COMPILED_MODEL_PATH="di_traced_model_tp${TP_DEGREE}_b${BATCH_SIZE}/"

inference_demo \
    --model-type llama \
    --task-type causal-lm \
    run \
    --model-path $MODEL_PATH \
    --compiled-model-path $COMPILED_MODEL_PATH \
    --torch-dtype bfloat16 \
    --tp-degree $TP_DEGREE \
    --batch-size $BATCH_SIZE \
    --ctx-batch-size 1 \
    --tkg-batch-size $BATCH_SIZE \
    --is-continuous-batching \
    --max-context-length 8192 \
    --seq-len 8192 \
    --on-device-sampling \
    --fused-qkv \
    --global-topk 256 --dynamic \
    --top-k 50 --top-p 0.9 --temperature 0.7 \
    --do-sample \
    --sequence-parallel-enabled \
```

(continues on next page)

(continued from previous page)

```

--qkv-kernel-enabled \
--attn-kernel-enabled \
--mlp-kernel-enabled \
--cc-pipeline-tiling-factor 1 \
--pad-token-id 2 \
--logical-neuron-cores 2 \
--context-encoding-buckets 256 512 1024 2048 4096 8192 \
--token-generation-buckets 512 1024 2048 4096 8192 \
--apply-seq-ids-mask \
--enable-bucketing \
--prompt "test prompt" \
--save-sharded-checkpoint \
--attn-block-tkg-nki-kernel-enabled \
--attn-block-tkg-nki-kernel-cache-update \
--k-cache-transposed \
--async-mode \
--compile-only

```

The `--apply-seq-ids-mask` flag is required for DI because it tells Neuron to only update the KV cache of the current sequence ID to ensure KV cache integrity, and ultimately, accuracy.

## Multi-Instance

For multi-instance run:

```
./compile.sh --tp-degree 64 --batch-size 4 --model-path path/to/your/downloaded/model
```

## Single-Instance

For single-instance run:

```
./compile.sh --tp-degree 32 --batch-size 4 --model-path path/to/your/downloaded/model
```

We compile for `tp-degree=32` because 1 prefill server will take up half of the Neuron Cores cores while the decode server will take up the other half.

## Launch the Prefill and Decode Servers

We provide a script called `server.sh`, which you can use to launch prefill and decode servers.

`NEURON_RT_ASYNC_SENDRECV_EXPERIMENTAL_ENABLED=1` is currently required as DI is still in beta. `NEURON_RT_ASYNC_SENDRECV_BOOTSTRAP_PORT=45645` is required to tell the Neuron Runtime which port to use for KV Cache transfer communications. `NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS=2` enables *Asynchronous Runtime Support*

The `KVTransferConfig` provided to both servers on startup have key information. `kv_connector=NeuronConnector` lets vLLM know to use the Neuron implementation for KV cache transfer. `kv_role=producer` lets vLLM know that this server's job is to do prefill. `kv_role=consumer` lets vLLM know that this server's job is to do decode. `neuron_core_offset=n` lets vLLM know that the model is hosted starting on the `n`th Neuron Core.

```
#!/bin/bash
# copy and paste me into a file called server.sh
# then run chmod +x server.sh

#!/bin/bash

# Parse command line arguments
while [[ $# -gt 0 ]]; do
    case $1 in
        --tp-degree)
            TP_DEGREE="$2"
            shift 2
            ;;
        --batch-size)
            BATCH_SIZE="$2"
            shift 2
            ;;
        --model-path)
            MODEL_PATH="$2"
            shift 2
            ;;
        --compiled-model-path)
            COMPILED_MODEL_PATH="$2"
            shift 2
            ;;
        --send-ip)
            SEND_IP="$2"
            shift 2
            ;;
        --recv-ip)
            RECV_IP="$2"
            shift 2
            ;;
        *)
            echo "Unknown parameter: $1"
            echo "Usage: $0 --tp-degree <value> --batch-size <value> --model-path <path> \
--compiled-model-path <path> --send-ip <ip> --recv-ip <ip>"
            exit 1
            ;;
    esac
done

export NEURON_RT_ASYNC_SENDRECV_BOOTSTRAP_PORT=45645
export NEURON_RT_ASYNC_SENDRECV_EXPERIMENTAL_ENABLED=1
export NEURON_COMPILED_ARTIFACTS="$COMPILED_MODEL_PATH"
export NEURON_SEND_IP="$SEND_IP"
export NEURON_RECV_IP="$RECV_IP"
export NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS=2

if [ "$SEND" = "1" ]; then
    PORT=8100
    if [ "$SINGLE_INSTANCE" = "1" ]; then
```

(continues on next page)

(continued from previous page)

```

    export NEURON_RT_VISIBLE_CORES=0-31
fi
TRANSFER_CONFIG='{
    "kv_connector":"NeuronConnector",
    "kv_buffer_device":"cpu",
    "kv_role":"kv_producer",
    "kv_rank":0,
    "kv_parallel_size":2,
    "kv_buffer_size":2e11,
    "kv_ip":"'$NEURON_SEND_IP'",
    "neuron_core_offset": 0
}'

else
    PORT=8200
    if [ "$SINGLE_INSTANCE" = "1" ]; then
        NC_OFFSET=32
        export NEURON_RT_VISIBLE_CORES=32-63
    else
        NC_OFFSET=0
    fi
    TRANSFER_CONFIG='{
        "kv_connector":"NeuronConnector",
        "kv_buffer_device":"cpu",
        "kv_role":"kv_consumer",
        "kv_rank":1,
        "kv_parallel_size":2,
        "kv_buffer_size":2e11,
        "kv_ip":"'$NEURON_SEND_IP'",
        "neuron_core_offset": "'$NC_OFFSET'"
    }'
fi

python3 -m vllm.entrypoints.openai.api_server \
    --model "$MODEL_PATH" \
    --max-num-seqs "$BATCH_SIZE" \
    --max-model-len 8192 \
    --tensor-parallel-size "$TP_DEGREE" \
    --device neuron \
    --use-v2-block-manager \
    --override-neuron-config "{}" \
    --kv-transfer-config "$TRANSFER_CONFIG" \
    --port "$PORT"

```

You may need multiple terminals to run the following commands.

For multi-instance choose one instance to be your prefill instance and one instance to be your decode instance. Get the IP addresses of them by running `hostname -i` and use them in the commands below. Single instance can use 127.0.0.1 as the IP address since prefill and decode always run on the same instance.

## Multi-Instance

To launch a prefill server for multi-instance run:

```
SEND=1 ./server.sh --tp-degree 64 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp64_b4/ \
    --neuron-send-ip prefill_ip --neuron-recv-ip decode_ip
```

To launch a decode server open up a new tab and run:

```
./server.sh --tp-degree 64 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp64_b4/ \
    --neuron-send-ip prefill_ip --neuron-recv-ip decode_ip
```

## Single-Instance

To launch a prefill server for single-instance run:

```
SEND=1 SINGLE_INSTANCE=1 ./server.sh --tp-degree 32 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp32_b4/ \
    --neuron-send-ip 127.0.0.1 --neuron-recv-ip 127.0.0.1
↪ 1
```

To launch a decode server open up a new tab and run:

```
SINGLE_INSTANCE=1 ./server.sh --tp-degree 32 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp32_b4/ \
    --neuron-send-ip 127.0.0.1 --neuron-recv-ip 127.0.0.1
```

When you see the line `INFO: Uvicorn running on http://0.0.0.0:8100` (Press CTRL+C to quit) on your prefill and decode server tabs your servers are ready.

## Launch a Router (Proxy Server)

Both servers need to receive a request to run inference. The component that does this job is called the router as mentioned in *DI Developer Guide*. We offer an implementation of a router called the `neuron-proxy-server`. The `neuron-proxy-server` is an entrypoint in our fork of vLLM which launches a proxy server that will take a request and forward it to both the prefill and decode servers. It will then capture their responses and format them back to the user.

The implementation of the `neuron-proxy-server` can be found [here](#).

For multi-instance run the router as another process on your prefill instance. For single-instance run the router as another process on your Trn2.

A router can run on any instance that has a connection to both the prefill and decode nodes. For multi-instance 1P1D, it makes the most sense to have the router on the prefill node to reduce network latency.

Launch the proxy server by running:

```
pip install quart # only install one time
neuron-proxy-server --prefill-ip your_prefill_ip --decode-ip your_decode_ip --prefill-
↪port 8100 --decode-port 8200
```

The proxy server is ready when you see the line `INFO:hypercorn.error:Running on http://127.0.0.1:8000` (CTRL + C to quit)

## Test the DI Setup

Run a sanity check to see if your DI setup is working by sending a curl request to the `neuron-proxy-server`:

```
curl -s http://localhost:8000/v1/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "path/to/your/downloaded/model",
  "prompt": ["a tornado is a"],
  "max_tokens": 10,
  "temperature": 0
}'
```

A successful response looks like: `{"id": ... : [{"index":0,"text":" rotating column of air that forms during severe thunderstorms" ... } }`

The `neuron-proxy-server` also supports the streaming of responses. It can be tested by:

```
curl -s http://localhost:8000/v1/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "path/to/your/downloaded/model",
  "prompt": ["a tornado is a"],
  "max_tokens": 10,
  "temperature": 0,
  "stream": true
}'
```

## Benchmark the DI Setup

### Install LLMPerf

We will use `LLMPerf` to measure the performance.

LLMPerf will send requests to the `neuron-proxy-server` and capture data including Time To First Token, Inter Token Latency and throughput.

Install `llmperf` into the `aws_neuronx_venv_pytorch_2_7_nxd_inference` virtual environment.

For multi-instance LLMperf is only required to be installed on the prefill instance where you will run benchmarking.

```
git clone https://github.com/ray-project/llmperf.git
cd llmperf
pip install -e .
```

Once you have installed LLMPerf, apply the `neuron_perf.patch` as described in [LLM Inference Benchmarking guide](#).

Next use the `llmperf.sh` script to run benchmarks.

```
#!/bin/bash
# copy and paste me into a file called llmperf.sh
# then run chmod +x llmperf.sh

# Set environment variables
export OPENAI_API_BASE="http://localhost:8000/v1"
export OPENAI_API_KEY="mock_key"

python llmperf/token_benchmark_ray \
  --model=$MODEL_PATH \
  --tokenizer=$MODEL_PATH \
  --mean-input-tokens=1024 \
  --stddev-input-tokens=0 \
  --mean-output-tokens=100 \
  --stddev-output-tokens=10 \
  --max-num-completed-requests=200 \
  --timeout=1720000 \
  --num-concurrent-requests=4 \
  --results-dir=llmperf_results \
  --llm-api=openai \
  --additional-sampling-params '{"top_k": 50, "top_p": 0.9, "temperature": 0.7}'
```

Since the `llmperf.sh` script sends requests to localhost, it should be run on the same instance the router is running on.

In multi-instance that means as a separate process on your prefill instance. For single instance that means a separate process on your Trn2.

```
MODEL_PATH=path/to/your/downloaded/model ./llmperf.sh
```

This will run a total of 200 requests and your final output should have the line: `Completed Requests Per Minute: xx.xxxxxxx`. Scroll up to see metrics such as Inter Token Latency and Time To First Token.

## Benchmark a Non-DI Continuous Batching Setup for Comparison

To compare Disaggregated Inference against non-DI continuous batching we will run benchmarks without Disaggregated Inference.

First kill all DI servers. Then kill the `neuron-proxy-server`.

We will run the same compiled model as a singular server for non-DI benchmarks. For single instance non-DI benchmarking we will start one TP=32 server. For multi-instance non-DI benchmarking we will start one TP=64 server. This means you do not need your second (decode) instance for this step. Latency can be compared directly in DI vs non-DI benchmarks. You might need to adjust the throughput related metrics based on number of instances to compare apples-to-apples between DI and non-DI. In this case, Non-DI throughput should be doubled before comparing with DI as the non-DI benchmark uses half the amount of hardware.

Use the `baseline_server.sh` to launch a vLLM server without DI.

```
#!/bin/bash
# copy and paste me into a file called baseline_server.sh
```

(continues on next page)

(continued from previous page)

```

# then run chmod +x baseline_server.sh

#!/bin/bash

# Parse command line arguments
while [[ $# -gt 0 ]]; do
    case $1 in
        --tp-degree)
            TP_DEGREE="$2"
            shift 2
            ;;
        --batch-size)
            BATCH_SIZE="$2"
            shift 2
            ;;
        --model-path)
            MODEL_PATH="$2"
            shift 2
            ;;
        --compiled-model-path)
            COMPILED_MODEL_PATH="$2"
            shift 2
            ;;
        *)
            echo "Unknown parameter: $1"
            echo "Usage: $0 --tp-degree <value> --batch-size <value> --model-path <path> \
                --compiled-model-path <path>"
            exit 1
            ;;
    esac
done

export NEURON_COMPILED_ARTIFACTS="$COMPILED_MODEL_PATH"
export NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS=2

if [ "$SINGLE_INSTANCE" = "1" ]; then
    NEURON_RT_VISIBLE_CORES=0-31
fi

python3 -m vllm.entrypoints.openai.api_server \
    --model "$MODEL_PATH" \
    --max-num-seqs "$BATCH_SIZE" \
    --max-model-len 8192 \
    --tensor-parallel-size "$TP_DEGREE" \
    --device neuron \
    --use-v2-block-manager \
    --override-neuron-config "{}" \
    --port 8000

```



## Multi-Instance

Launch for multi-instance with:

```
./baseline_server.sh --tp-degree 64 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp64_b4/
```

## Single-Instance

Launch for single-instance with:

```
SINGLE_INSTANCE=1 ./baseline_server.sh --tp-degree 32 --batch-size 4 \
    --model-path path/to/your/downloaded/model \
    --compiled-model-path di_traced_model_tp32_b4/
```

Now we have a server launched with the same underlying model but with DI turned off.

Then on the same instance run llmperf which will now directly send requests to the server instead of going through a proxy:

```
MODEL_PATH=path/to/your/downloaded_model ./llmperf.sh
```

This will run a total of 200 requests and your final output should have the line: Completed Requests Per Minute: xx.xxxxxxx. Scroll up to see metrics such as Inter Token Latency and Time To First Token.

## Known Issues

ENC:kv\_store\_acquire\_file\_lock Failed to open kv store server lock file Permission denied usually means that another user on the system ran a DI workload and left behind a lock file that the current user does not have access to. The solution is to delete /tmp/nrt\_kv\_store\_server.lock file.

This section includes tutorials that you can follow to get started with the NeuronX Distributed (NxD) Inference library.

- [Tutorial: Deploying Llama3.1 405B \(Trn2\)](#)
- [Tutorial: Deploying Llama3.2 Multimodal Models](#)
- [Tutorial: Using Speculative Decoding to improve Llama-3.3-70B inference performance on Trn2 instances](#)
- [Tutorial: Scaling LLM Inference with Data Parallelism on Trn2](#)
- [Tutorial: Multi-LoRA serving for Llama-3.1-8B on Trn2 instances](#)
- [Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances](#)
- [Tutorial: Evaluating Accuracy of Llama-3.1-70B on Neuron using open source datasets](#)
- [nxdi-trn2-llama3.3-70b-apc-tutorial](#)
- [Tutorial: Disaggregated Inference on Trn2 \[BETA\]](#)

### 3.2.6 Application Notes

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### Introducing NeuronX Distributed (NxD) Inference

##### Table of contents

- *What are we introducing?*
- *How can I install NxD Inference library?*
- *I am currently using the Transformers NeuronX library for inference. How does the NxD Inference library affect me?*
- *I am currently using vLLM with Transformers NeuronX library for inference. Does NxD Inference library support vLLM ?*
- *What features and models are available in Transformers NeuronX (TNx) but not yet in NeuronX Distributed Inference?*
- *I currently use Hugging Face TGI serving engine for deploying and serving Large Language Models (LLMs) on Neuron. How does NxD Inference library affect me?*
- *I am new to Neuron and have inference workloads, what library should I use?*
- *Additional Resources*

#### What are we introducing?

Starting with the Neuron SDK 2.21 release, we are introducing NxD Inference, an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. NxD Inference is designed for optimized inference, enabling quick onboarding of PyTorch models with minimal changes. It features a modular architecture that facilitates easy integration of HuggingFace PyTorch models and is compatible with serving engines like vLLM.

Please see [NxD Inference](#) for NxD Inference overview and documentation.

#### How can I install NxD Inference library?

Please refer to [NxD Inference Setup Guide](#) for installation instructions.

#### I am currently using the Transformers NeuronX library for inference. How does the NxD Inference library affect me?

If you are using Transformers NeuronX (TNx) in production, you can continue doing so. However, if you are planning to onboard new models to Neuron for inference, NxD Inference offers several advantages to consider.

NxD Inference is designed to enable easy on-boarding of PyTorch models and comes with new features and enhanced support:

- **Hardware Support:** While TNx is not supported on Trn2, NxD Inference supports all platforms (Trn1, Inf2, and Trn2)

- **Simplified interface:** To simplify model development with NxD Inference, you write modeling code using PyTorch with standard Python, rather than using PyHLO as in TNx.
- **Easy Migration:** NxD Inference was designed to provide seamless migration from TNx, especially if you are using it with vLLM. You can migrate your existing TNx inference scripts using the [migration guide](#)
- **Enhanced Capabilities:** NxD Inference offers more comprehensive support for MoE models and multimodal models (Llama 3.2) compared to TNx
- **Future Development:** New inference features and support for advanced model architectures (like multi-modality/video models) will be focused on NxD Inference

### I am currently using vLLM with Transformers NeuronX library for inference. Does NxD Inference library support vLLM ?

Yes, NxD Inference library supports vLLM inference engine. Neuron vLLM integration in 2.21 release will start supporting both NxD Inference and Transformers NeuronX libraries. To use vLLM with NxD Inference library, you can refer to the [vLLM User Guide for NxD Inference](#).

### What features and models are available in Transformers NeuronX (TNx) but not yet in NeuronX Distributed Inference?

While NxD Inference supports most features and models available in TNx, there are some differences in current support that users should be aware of.

**Features that are not yet supported in NxD Inference:** The following TNx features aren't supported yet in the NxD Inference library.

- Multi-Node Inference support

**Models not part of NxD Inference Model Hub:** The following models are included in Transformers NeuronX but not currently in NxD Inference library:

- Bloom
- GPT2
- GPT-J
- GPT-NEOX

If you need to use these models with NxD Inference, we encourage you to follow the [onboarding models developer guide](#). The onboarding process in NxD Inference is more straightforward compared to TNx due to its PyTorch-based architecture.

### I currently use Hugging Face TGI serving engine for deploying and serving Large Language Models (LLMs) on Neuron. How does NxD Inference library affect me?

If you are currently using Hugging Face TGI serving engine to deploy models on Neuron, the introduction of NxD Inference library will not have any impact and you can continue to use your existing inference workloads. Hugging Face TGI integrates with Neuron SDK Inference libraries in a way that abstracts the underlying library for the users.

## I am new to Neuron and have inference workloads, what library should I use?

We recommend you use NxD Inference for your model inference workloads. To learn how to get started using NxD Inference, see the [NxD Inference](#) documentation

## Additional Resources

- [NxD Inference](#)
- [NxD Inference Overview](#)
- [NxD Inference Release Notes \(neuronx-distributed-inference\)](#)

This document is relevant for: Inf1, Inf2, Trn1, Trn2

## Parallelism Techniques for LLM Inference

### Table of contents

- [Overview](#)
- [Tensor Parallelism](#)
  - [How to Use Tensor Parallelism with NxD Inference](#)
- [Sequence Parallelism](#)
  - [How to Use Sequence Parallelism with NxD Inference](#)
- [Flash Decoding](#)
  - [How to Use Flash Decoding with NxD Inference](#)
- [Data Parallelism](#)
  - [How to Use Data Parallelism with NxD Inference](#)

## Overview

Large language models (LLMs) have grown exponentially in size in the past few years, requiring increasing accelerator memory to run the model. In order to effectively generate predictions from an LLM, it is often necessary to use one or more **parallelism techniques** to shard operations across multiple available accelerators. **Model parallelism**, such as tensor and sequence parallelism described in this document, can reduce memory requirements per NeuronCore by sharding the model across multiple cores. **Data parallelism**, on the other hand, enables higher throughput by sharding input data.

## Tensor Parallelism

Tensor parallelism is a technique in which a tensor is split into a number of chunks along the intermediate dimension, resulting in sharding not only model parameters but also intermediate activations. Tensor parallelism has relatively high communication volume and presents a synchronization point in forward pass, making it costly to scale beyond 1 node. When tensors are sharded across multiple EC2 instances, the collective communication at these synchronization points must occur through network interfaces like Elastic Fabric Adapter (EFA) instead of the faster chip-to-chip NeuronLink connections.

A basic transformer MLP block contains a single matrix multiplication (matmul) called the up-projection, which increases the dimensionality from the `hidden_size` to the `intermediate_size`, and a single output matmul called the down projection, which reduces the dimensionality back to the `hidden_size`, with a non-linear activation function in-between. In order to avoid running collective operations (synchronization point) after each matrix multiply, we defer collective to run after 2nd linear layer. To ensure correctness of non-linear activation function computation ( $f(x+y) \neq f(x) + f(y)$  for non-linear  $f$  like `silu` in `SwiGLU`), we split first linear layer along columns (`ColumnParallel`) and second linear layer along rows (`RowParallel`), then run an `AllReduce` collective operation at the end.

Modern transformer architectures use `SwiGLU` activation function, where the MLP block has 3 matrices, first up and gate projection and later a down projection. We can view up and gate projection as the same (referred to as first matrix multiply or first linear layer) in this context because they have the same sharding approach. In this case up and gate projection is column parallel, while down projection is row parallel.

In attention, we similarly split Q, K and V projections in column parallel fashion and use row parallel for final output (O) projection, then run `AllReduce` with input tensor size equal to `batch_size x sequence_length x hidden_size x per_element_bytes` bytes. Here, `per_element_bytes` depends on the numerical format of your tensors. When using `BF16`, for example, it would be 2. `AllReduce` input tensor size is the same for both MLP and attention blocks, resulting in two `AllReduce` operations with the same input size and output size as per `AllReduce` algorithm per transformer layer.

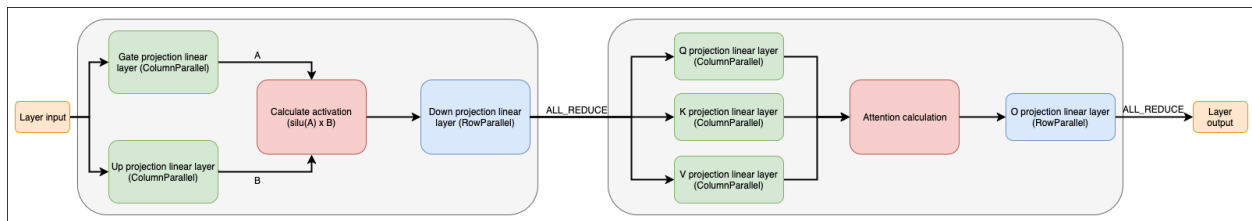


Fig. 3.6: Image visualizing transformer layer like llama3 with SwiGLU activation layer in MLP.

## How to Use Tensor Parallelism with NxN Inference

Tensor parallelism can be enabled by setting the `tp_degree` parameter in `NeuronConfig`. See [Tensor-parallelism support](#) for more detail.

Code example, defining `NeuronConfig`:

```
neuron_config = NeuronConfig(tp_degree=32)
```

See [Tensor Parallelism Overview](#) for a detailed reference of the concepts underlying tensor parallelism.

## Sequence Parallelism

One drawback of tensor parallelism is that it replicates attention/MLP layer norm and dropout operations across all NeuronCores. These operations are less compute intensive compared to linear layers, but still requires significant memory. These computations are independent along the sequence dimension, allowing us to shard across the sequence dimension. This requires AllGather in the transition from a sequence to a tensor parallel region and ReduceScatter in the transition from tensor to sequence parallel region during inference. Sequence parallelism is especially useful for longer sequences and usually used in conjunction with tensor parallelism.

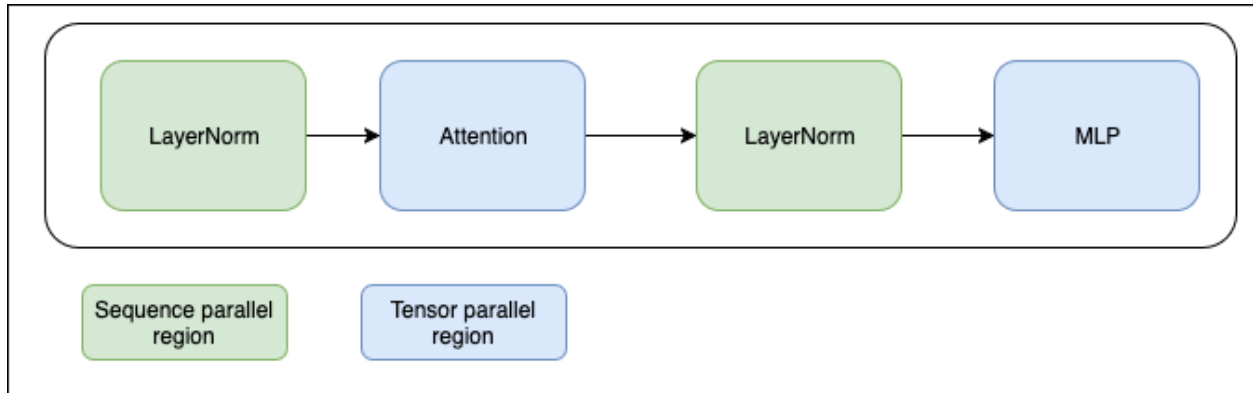


Fig. 3.7: Image visualizing how sequence and tensor parallelism intertwine in transformer layer like Llama 3.

## How to Use Sequence Parallelism with NxD Inference

Sequence parallelism can be enabled by setting the `sequence_parallel_enabled` parameter in `NeuronConfig`. See [Sequence Parallelism](#) for more detail.

Code example, defining `NeuronConfig`:

```
neuron_config = NeuronConfig(sequence_parallel_enabled=True)
```

## Flash Decoding

Flash decoding enables inference on long sequences by partitioning the KV cache. The technique is useful for long sequences and used in decoding phase. It is motivated by the fact that assuming KV caching, K and V memory footprint scales with sequence length, while Q has sequence length equal to 1 during decoding.

Flash decoding shards K and V, and at the start uses AllGather to gather all Q heads in each KV partition. Each KV partition computes partial softmax (also called log-sum-exp) which uses AllGather to compute log-sum-exp scaling factor and correction denominator after “local” attention computation (multiply Q and K, then apply the mask). Lastly, the algorithm performs ReduceScatter on attention results at the end.

## How to Use Flash Decoding with NxD Inference

Flash decoding can be enabled by setting the `flash_decoding_enabled` parameter in `NeuronConfig`. The technique is only supported with GQA (grouped query attention).

Code example, defining `NeuronConfig`:

```
neuron_config = NeuronConfig(flash_decoding_enabled=True)
```

## Data Parallelism

Data parallelism will replicate the model (same architecture, weights and hyperparameters) but will shard input data. By distributing the data across NeuronCores or even different instances, data parallelism reduces the total execution time of large batch size inputs using parallelization across sharded inputs instead of sequential execution. Compared to batch parallel where KV cache is sharded, each data parallel replica has its own individual cache for separate sequences.

Data parallelism works as standalone technique or can be used in conjunction with other model sharding techniques such as tensor parallelism. For example, Trn2 instances has 64 NeuronCores available when using default Logical NeuronCore configuration (LNC) set to 2, so you can use a tensor parallel degree of 16 and a data parallel degree of 4, resulting in four copies of the model, each with disjunct data partitioning and with each model sharded across 16 logical NeuronCores. Model replicas can run on the same instance or separate instances. Data parallelism doesn't introduce any additional collective operations during inference.

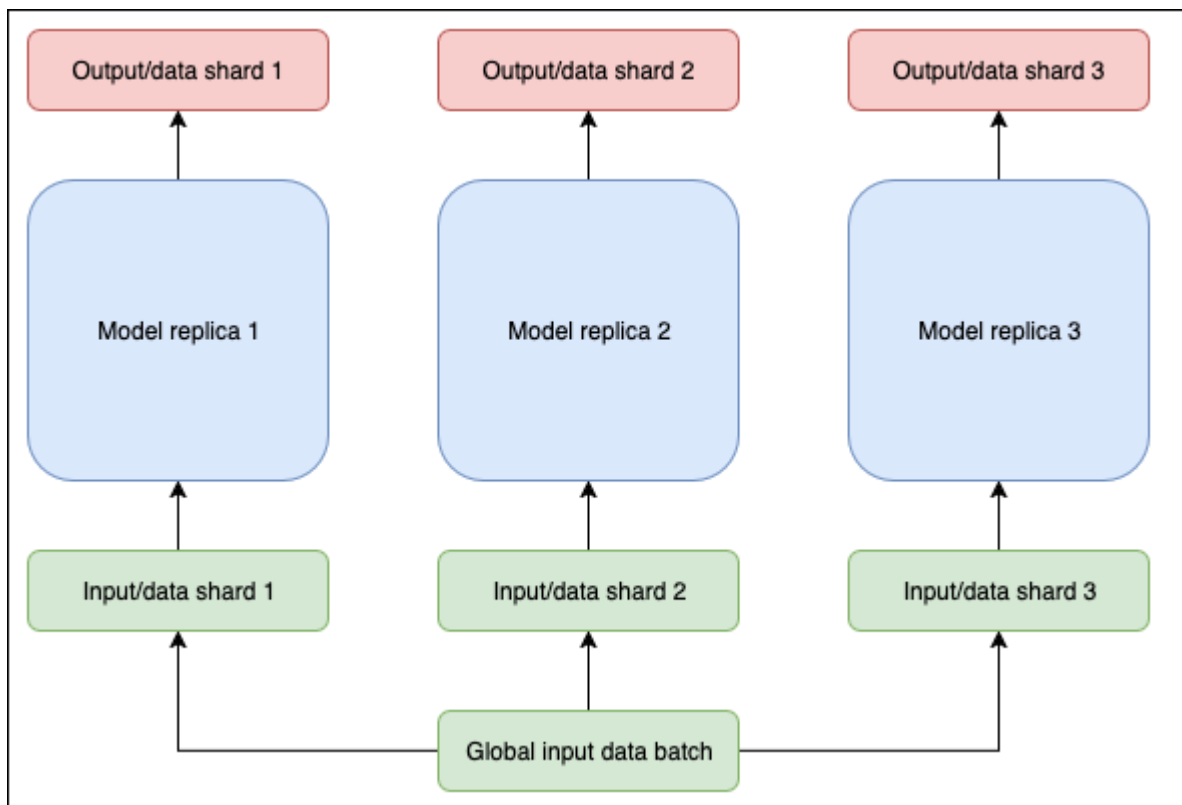


Fig. 3.8: Image visualizing how data parallelism shards inputs.

## How to Use Data Parallelism with NxD Inference

See *Tutorial: Scaling LLM Inference with Data Parallelism on Trn2* for detailed guidance on how to use vLLM to apply data parallelism along with tensor parallelism to increase model inference throughput in NxDI.

- *Introducing NeuronX Distributed (NxD) Inference*
- *Parallelism Techniques for LLM Inference*

### 3.2.7 NxD Inference Misc

This document is relevant for: Inf1, Inf2, Trn1, Trn2

#### NxD Inference Release Notes (neuronx-distributed-inference)

##### Table of contents

- *Neuronx Distributed Inference [0.4.7422] (Neuron 2.24.0 Release)*
- *Neuronx Distributed Inference [0.3.5591] (Neuron 2.23.0 Release)*
- *Neuronx Distributed Inference [0.2.0] (Beta) (Neuron 2.22.0 Release)*
- *Neuronx Distributed Inference [0.1.1] (Beta) (Neuron 2.21.1 Release)*
- *Neuronx Distributed Inference [0.1.0] (Beta) (Neuron 2.21 Release)*
- *Neuronx Distributed Inference [0.1.0] (Beta) (Trn2)*

This document lists the release notes for Neuronx Distributed Inference library.

#### Neuronx Distributed Inference [0.4.7422] (Neuron 2.24.0 Release)

Date: 06/24/2025

- Models
  - Qwen2.5 text models, which are tested on Trn1. Compatible models include:
    - \* Qwen2.5-0.5B-Instruct
    - \* Qwen2.5-7B-Instruct
    - \* Qwen2.5-32B-Instruct
    - \* Qwen2.5-72B-Instruct
- Features
  - Automatic Prefix Caching support (APC) through vLLM. APC improves efficiency by reusing KV cache from previous queries if the new query shares a prefix. APC can significantly improve TTFT based on how often different queries share the same prefixes. Performance gains are greater when requests have longer shared prefixes and when there is a higher frequency of prefix sharing across requests. For example, with Llama3.3 70B on Trn2, you can observe a 3.2x TTFT improvement with the math.math dataset (90% cache hit), a 1.6x TTFT improvement with a Sonnet dataset with 2K prompt length (25% cache hit), or no TTFT improvement with the HumanEval dataset (0% cache hit). For more information, see [nxdi-prefix-caching](#) and [nxdi-trn2-llama3.3-70b-apc-tutorial](#).



- Disaggregated Inference (DI) support through vLLM (Beta). Disaggregated Inference is also known as disaggregated serving, disaggregated prefill, or p/d disaggregation. DI separates the prefill and decode phase of inference onto different hardware resources. DI can improve inter token latency (ITL) by eliminating prefill stall in continuous batching, where decode is paused to perform prefill for a new incoming request. With DI, you can also scale prefill and decode resources independently to further improve performance. For more information, see [Disaggregated Inference \[BETA\]](#).
- Context parallelism in NeuronAttentionBase (Beta). Context parallelism distributes context processing across multiple NeuronCores. Context parallelism improves TTFT, particularly at higher sequence lengths where the number of KV heads is low. To use context parallelism, set `cp_degree` in `NeuronConfig`.
- Mixed-precision parameters in modeling code. This feature enables you to configure each module's dtype independently. To use mixed-precision parameters, set `cast_type="as-declared"` in `NeuronConfig`. Note: The default behavior (`cast_type="config"`) is to cast all parameters to the `torch_dtype` in `NeuronConfig`.
- Output logits when using on-device sampling. To output logits, enable `output_logits` in `NeuronConfig`. Note that this flag impacts performance and should only be used for debugging model logits.
- Other changes
  - Add support for PyTorch 2.7. This release includes support for PyTorch 2.5, 2.6, and 2.7.
  - Upgrade `transformers` requirement from v4.48 to v4.51.
  - Re-enable warmup on Trn2. NxD Inference disabled warmup on Trn2 in the previous release due to an issue that prevented certain model configurations from loading correctly. That issue is now fixed.
  - Update the behavior of the `attn_kernel_enabled` attribute in `NeuronConfig`, which configures whether to use the flash attention kernel. Previously, `True` meant to enable in all cases where supported, and `False` meant to auto-enable where beneficial (defaults to `False`). Now, `attn_kernel_enabled=False` disables the flash attention kernel in all cases. To use the previous auto-enable behavior, set `attn_kernel_enabled=None`. The default value for `attn_kernel_enabled` is now `None` to retain the same default behavior as before.
  - Enable `--verify-hlo` flag during compilation. Now, if an HLO is invalid, compilation will fail. Previously, in certain scenarios, the compiler would successfully compile invalid HLOs.
  - Update the flash attention kernel strategy to use the attention kernel on Trn2 in all cases where it's supported. This change fixes an issue where certain context lengths failed to trace.
  - Add `logical_nc_config` as an argument to the `build_module` and `build_function` test utilities, so you can use these utilities to test modules/functions for Trn2 using LNC2.
  - Other minor fixes and improvements.

## Known Issues and Limitations

### Increased Device Memory Usage for Certain Configurations

Certain model configurations require slightly more device memory than in previous releases. If your model used close to the maximum amount of device memory in previous releases, this increase could cause it to fail to load after you compile it with this release. This issue is most likely to affect Llama3.1-405B configurations that use a large number of buckets.

If this issue occurs, you will see the following error during model load.

```
ERROR   TDRV:log_dev_mem                               Failed to allocate 512.000MB_
↪(alignment: 4.000MB, usage: shared scratchpad) on ND14:NC 6
```

To avoid this error, reduce the number of buckets you use, or reduce the sequence lengths that you use in each bucket.

## Neuronx Distributed Inference [0.3.5591] (Neuron 2.23.0 Release)

Date: 05/20/2025

NxD Inference is now GA and out of beta in the Neuron 2.23 release.

### Features in this Release

- Features

- Shard-on-load for weight sharding is now enabled by default. With this change, end-to-end compile and load time is reduced by up to 70% when sharding weights. This change significantly reduces compile time by skipping weight sharding and serialization during compile, but may lead to increased load time. For example, for Llama 3.1 405B, end-to-end compile and load time is reduced from 40 minutes to 12 minutes. For best load performance, you can continue to serialize sharded weights by enabling `save_sharded_checkpoint` in `NeuronConfig`. For more information, see *NxD Inference Weights Sharding Guide*.
- Neuron Persistent Cache. NxD Inference now supports Neuron Persistent Cache, which caches compiled model artifacts to reduce compilation times. For more information, see *Neuron Persistent Cache*.
- Support for an attention block kernel for token generation. This kernel performs QKV projections, RoPE, attention, and output projections. You can use this kernel with Llama3-like attention on Trn2 to improve token gen performance. To use this kernel, enable `attn_block_tkg_nki_kernel_enabled` in `NeuronConfig`.
  - \* This kernel can also update the KV cache in parallel with each layer's attention compute to further improve performance. This functionality hides the latency of the KV cache update that is otherwise done for all layers at once at the end of each token generation iteration. To enable in-kernel KV cache updates, enable `attn_block_tkg_nki_kernel_cache_update` in `NeuronConfig`. When in-kernel KV cache updating is enabled, you can also enable `k_cache_transposed` to further improve the performance.
- Automatically extract `target_modules` and `max_lora_rank` from LoRA checkpoints. You no longer need to set these arguments manually.
- Support fused residual add in the QKV kernel. This feature improves the performance of context encoding at short sequence lengths. To use this feature, enable the `qkv_kernel_fuse_residual_add` flag in `NeuronConfig`.

- Backward incompatible changes

- Remove `set_async_mode(async_mode)` from `NeuronBaseForCausalLM`, as this feature didn't work as intended. Async mode cannot be enabled or disabled after the model is loaded. To enable async mode, set `async_mode=True` in `NeuronConfig`.

- Other changes

- Disable warmup for Trn2. This change avoids an issue that prevents certain model configurations from loading correctly. When warmup is disabled, you will see lower performance on the first few requests to the model. This change also affects initial performance for serving through vLLM. Warmup will work in many cases where it is now disabled, so you can try to reenble warmup by setting `skip_warmup=False` in `NeuronConfig`. Alternatively, you can manually warm up the model by sending a few requests to each bucket after loading the model.

- Fix an issue where when continuous batching and bucketing were enabled, NxDI padded each input to the largest sequence in the batch, rather than the next largest bucket for that input. This change improves performance when using continuous batching with bucketing, including through vLLM.
- Add a `num_runs` parameter to `benchmark_sampling`, so you can configure the number of runs to perform when benchmarking.
- Silence unimportant error messages during warmup.
- NeuronConfig now includes a `disable_kv_cache_tiling` flag that you can set to disable KV cache tiling in cases where it was previously enabled by default.
- Update the package version to include additional information in the version tag.
- Other minor fixes and improvements.

## Neuronx Distributed Inference [0.2.0] (Beta) (Neuron 2.22.0 Release)

Date: 04/03/2025

### Models in this Release

- Llama 3.2 11B (Multimodal)

### Features in this Release

- Multi-LoRA serving. This release adds support for multi-LoRA serving through vLLM by loading LoRA adapters at server startup. Multi-LoRA serving is currently supported for Llama 3.1 8B, Llama 3.3 70B, and other models that use the Llama architecture.
- Custom quantization. You can now specify which layers or modules in NxDI to quantize or keep in full precision during inference. To configure which layers or modules to skip during quantization, use the `modules_to_not_convert` and `draft_model_modules_to_not_convert` attributes in NeuronConfig.
- Models quantized through external libraries. NxDI now supports inference of models that are quantized externally using quantization libraries such as LLMCompressor.
- Async mode. This release adds support for async mode, which improves performance by asynchronously preparing the next forward call to a mode. To use async mode, enable the `async_mode` flag in NeuronConfig.
- CPU inference. You can now run models on CPU and compare against output on Neuron to debug accuracy issues. To use this feature, enable the `on_cpu` flag in NeuronConfig.
- Unit/module testing utilities. These common utilities include `build_module`, `build_function`, and `validate_accuracy`, which enable you to build a module or function and validate its accuracy on Neuron. You can use these utilities in unit/integration tests to verify your modeling code works correctly.
- Add support for models that use a custom `head_dim` value from InferenceConfig. This change enables support for models where `head_dim` isn't equivalent to `hidden_size` divided by `num_attention_heads`.
- Input capture hooks. When you call the `NeuronBaseForCausalLM` forward function, you can provide an `input_capture_hook` function that will be called with the model inputs as arguments.
- Runtime warmup. To improve the performance of the first request sent to a model, NxDI Inference now warms up the model during load. You can disable this behavior with the `skip_warmup` flag in NeuronConfig.

## Backward Incompatible Changes

- Fix the behavior of the `do_sample` sampling flag. Previously, NxDI used greedy sampling when `do_sample=True`, which was a bug because `do_sample=True` should result in multinomial sampling. If you use `do_sample=True` in a config where you intend to use greedy sampling, you must change it to `do_sample=False`. As part of this change, the default value for `do_sample` is now `False`.
- Enforce that tensors in a model's `state_dict` don't share memory with other tensors. This change can cause models to fail to load if their tensors share memory, which now results in an error: `RuntimeError: Error while trying to find names to remove to save state dict`. To fix this issue, apply `.clone().detach().contiguous()` to the model's `state_dict`, and re-shard the weights.
- Change the quantization `state_dict` keys from `weight_scale` to `scale` to match the Nx quantization scale keys and avoid any confusion. If you use quantization and have sharded weights from earlier versions of NxDI, you must re-shard the weights.
- If you use a model that skips quantization for certain modules (such as in Llama 3.1 405B FP8), you must now specify `modules_not_to_convert` to configure the modules that skip quantization.
- Validate when input size exceeds the model's maximum length (`max_context_length` or `max_length`). Nx Inference now throws a `ValueError` if given an input that's too large. To enable the previous behavior, where input is truncated to the maximum length, enable the `allow_input_truncation` flag in `NeuronConfig`.

## Other Changes

- Improve model performance by up to 50% (5-20% in most cases) by eliminating overheads in logging.
- Upgrade `transformers` from v4.45 to v4.48.
- Deprecate `NeuronConfig`'s `logical_neuron_cores` attribute and replace it with `logical_nc_config`. The LNC config is now automatically set from the `NEURON_LOGICAL_NC_CONFIG` environment variable if set.
- Deprecate `NeuronConfig`'s `trace_tokengen_model` attribute. This attribute is now determined dynamically based on other configuration attributes.
- Improve the performance of on-device sampling.
- When running Llama models with LNC2, the sharded flash attention kernel is now automatically enabled when context length is 256 or greater. Previously, this kernel was enabled for context length of 1024 or greater. This change improves performance at smaller context lengths.
- `NeuronConfig` now includes a `skip_sharding` flag that you can enable to skip weight sharding during model compilation. This option is useful in cases where you have already sharded weights, such as during iterative development, so you can iterate without re-sharding the weights each time you compile the model.
- `NeuronApplicationBase` now includes a `shard_weights` function that you can use to shard weights independent of compiling the model.
- Fix vanilla speculative decoding support for models with multiple EOS tokens.
- Other minor fixes and improvements.

## Known Issues and Limitations

- For some configurations that use continuous batching or vLLM, model warmup can cause `Numerical Error` during inference. If you encounter this error, set `skip_warmup=True` in `NeuronConfig` to disable warmup and avoid this issue. To disable warmup in vLLM, pass `"skip_warmup": true` in `override_neuron_config`. For more information about how to configure vLLM, see [vLLM Model Configuration](#).

```
RuntimeError: Failed to execute the model status=1003 message=Numerical Error
```

## Neuronx Distributed Inference [0.1.1] (Beta) (Neuron 2.21.1 Release)

Date: 01/14/2025

### Bug Fixes

- Fix minor issues with sampling params and add validation for sampling params.

## Neuronx Distributed Inference [0.1.0] (Beta) (Neuron 2.21 Release)

Date: 12/20/2024

### Features in this Release

NeuronX Distributed (NxD) Inference (`neuronx-distributed-inference`) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. Neuronx Distributed Inference includes a model hub and modules that users can reference to implement their own models on Neuron.

This is the first release of NxD Inference (Beta) that includes:

- Support for Trn2, Inf2, and Trn1 instances
- Support for the following model architectures. For more information, including links to specific supported model checkpoints, see [NxD Inference - Production Ready Models](#).
  - Llama (Text), including Llama 2, Llama 3, Llama 3.1, Llama 3.2, and Llama 3.3
  - Llama (Multimodal), including Llama 3.2 multimodal
  - Mistral (using Llama architecture)
  - Mixtral
  - DBRX
- Support for onboarding additional models.
- Compatibility with HuggingFace checkpoints and `generate()` API
- vLLM integration
- Model compilation and serialization
- Tensor parallelism
- Speculative decoding

- EAGLE speculative decoding
  - Medusa speculative decoding
  - Vanilla speculative decoding
- Quantization
- Dynamic sampling
- Llama3.1 405B Inference Example on Trn2
- Open Source Github repository: [aws-neuron/neuronx-distributed-inference](https://github.com/aws-neuron/neuronx-distributed-inference)

For more information about the features supported by NxDI, see *NxD Inference Features Configuration Guide*.

## Known Issues and Limitations

### Longer Load Times for Large Models

**Issue:** Users may experience extended load times when working with large models, particularly during weight sharding and initial model load. This is especially noticeable with models like Llama 3.1 405B.

**Root Cause:** These delays are primarily due to storage performance limitations.

**Recommended Workaround:** To mitigate this issue, we recommend that you store model checkpoints in high-performance storage options:

- **Instance store volumes:** On supported instances, instance store volumes offer fast, temporary block-level storage.
- **Optimized EBS volumes:** For persistent storage with enhanced performance.

By using these storage optimizations, you can reduce model load times and improve your overall workflow efficiency.

**Note:** Load times may still vary depending on model size and specific hardware configurations.

### Other Issues and Limitations

- Llama 3.2 11B (Multimodal) is not yet supported with PyTorch 2.5.
- The following model architectures are tested only on Trn1 and Inf2:
  - Llama (Multimodal)
- The following model architectures are tested only on Trn1:
  - Mixtral
  - DBRX
- The following kernels are tested only on Trn2:
  - MLP
  - QKV
- If you run inference with an prompt that is larger than the model's `max_context_length`, the model will generate incorrect output. In a future release, NxD Inference will throw an error in this scenario.
- Continuous batching (including through vLLM) supports batch size up to 4. Static batching supports larger batch sizes.
- To use greedy on-device sampling, you must set `do_sample` to `True`.

- To use FP8 quantization or KV cache quantization, you must set the `XLA_HANDLE_SPECIAL_SCALAR` environment variable to 1.

## Neuronx Distributed Inference [0.1.0] (Beta) (Trn2)

Date: 12/03/2024

### Features in this release

NeuronX Distributed (NxD) Inference (`neuronx-distributed-inference`) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances. Neuronx Distributed Inference includes a model hub and modules that users can reference to implement their own models on Neuron.

This is the first release of NxD Inference (Beta) that includes:

- Support for Trn2 instances
- Compatibility with HuggingFace checkpoints and `generate()` API
- vLLM integration
- Model compilation and serialization
- Tensor parallelism
- Speculative decoding
  - EAGLE speculative decoding
  - Medusa speculative decoding
  - Vanilla speculative decoding
- Quantization
- Dynamic sampling
- Llama3.1 405B Inference Example on Trn2
- Open Source Github repository: [aws-neuron/neuronx-distributed-inference](https://github.com/aws-neuron/neuronx-distributed-inference)

For more information about the features supported by NxDI, see *NxD Inference Features Configuration Guide*.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### Troubleshooting Guide for NxD Inference

This guide provides solutions for common issues encountered when using NxD Inference.

#### Table of contents

- *Accuracy Issues*
  - *Accuracy Degradation with Auto-Cast*
  - *Array indexing and in-place operations in Neuron*
- *Performance Issues*

- *Skip model warmup during inference*
- *Other Common Issues*
  - *Tensor Materialization During Tracing caused unexpected model behavior*
  - *Input Data Type Handling for int64/fp64 due to compiler dtype compatibility*

## Accuracy Issues

The primary methods for validating model accuracy on Neuron involve both token-by-token output matching and logit-level error analysis (relative or max absolute error) against a pre-calibrated GPU FP32 or CPU FP32 reference. When output deviations are observed, these can be systematically attributed to factors such as tokenizer/input discrepancies, amplification from large weight norms (high Lipschitz constants), quantization or precision loss, differences in operator implementation or kernel fusion, compiler optimization, or unintended hardware-level datatype casts.

When validating model accuracy on Neuron, it is important to recognize that predicting the exact output deviations from a high-precision reference (like CPU or GPU FP32) is theoretically NP-hard, due to the complex and nonlinear nature of large neural networks. Rather than attempting to anticipate every possible numerical difference, the recommended strategy is to systematically identify, localize, and diagnose deviations as they occur.

## Accuracy Degradation with Auto-Cast

**Issue:** You may observe accuracy degradation in model outputs when using the default auto-cast behavior of the Neuron compiler.

**Explanation:** By default, the Neuron compiler automatically casts certain operations to lower precision data types (BF16) to improve performance. While this works well for most cases, it can sometimes lead to accuracy issues, especially in operations involving integer-to-float conversions.

**Solution:** Use the `--auto-cast=none` compiler flag to disable automatic casting. This preserves the original precision of operations at the cost of some performance.

Example using `inference_demo`:

```
inference_demo --model-type llama --task-type causal-lm run \
  --model-path <path> \
  --compiled-model-path <path> \
  --torch-dtype bfloat16 \
  --tp-degree <value> \
  --batch-size <value> \
  --max-context-length <value> \
  --seq-len <value> \
  --on-device-sampling \
  --prompt "Your prompt here" \
  --compiler-args "--auto-cast=none"
```

Example using `NeuronConfig`:

```
from neuronx_distributed_inference.models.config import NeuronConfig

neuron_config = NeuronConfig(
    tp_degree=32,
    batch_size=1,
```

(continues on next page)



(continued from previous page)

```

max_context_length=1024,
seq_len=2048,
compiler_args="--auto-cast=None"
)

```

## Integer-to-Float Conversion Issues

**Issue:** Operations involving integer-to-float conversions (such as in rotary embeddings) may experience significant accuracy degradation when auto-cast is enabled.

**Explanation:** When integer values are converted to floating point and then automatically cast to lower precision (like BF16), the precision loss can be substantial. This is particularly problematic in operations like rotary embeddings where position IDs are converted to floating point for computing sin/cos values.

**Solution:** Use the `--auto-cast=None` compiler flag to prevent downcasting these operations. This is especially important for models that use rotary embeddings or similar position encoding mechanisms.

**Technical Details:** The issue occurs in operations like:

```

# Integer position IDs are converted to float for sin/cos computation
# Downcasting to BF16 here can cause significant precision loss
position_ids = position_ids.to(torch.bfloat16)
sin, cos = self.compute_sin_cos(position_ids)

```

## Memory Usage Considerations

**Note:** Using `--auto-cast=None` will increase memory usage as operations will use higher precision data types. Ensure your instance has sufficient memory when using this flag.

## Performance Impact

**Note:** Disabling auto-cast will typically result in slower inference. The exact performance impact depends on your model architecture and hardware configuration. Consider this trade-off when optimizing for accuracy.

## Array indexing and in-place operations in Neuron

**Issue:** When building attention masks, operations that combine array slicing with in-place modifications (e.g., `mask_i[: arx[0] * arx[1], :ntok] = 0`) can cause accuracy issues in Neuron. This is particularly problematic when the array indices are dynamically computed.

**Explanation:** The accuracy issue stems from two main factors:

1. Array Slicing with Dynamic Ranges:

```

# Problematic: Array slicing with dynamic range (arx[0] * arx[1])
mask_i[: arx[0] * arx[1], :ntok] = 0

```

- Uses computed indices to access specific portions of the tensor
- Dynamic ranges can lead to unpredictable memory access patterns

2. In-place Modifications:

```
# Problematic: Modifying tensor in-place
mask_i[...] = 0 # Direct modification of the original tensor
```

- Changes the original tensor's values directly
- Can cause issues with Neuron's memory management and optimization

**Solution:** Replace array slicing and in-place operations with element-wise operations:

```
# Instead of array slicing and in-place modification:
mask_i[: arx[0] * arx[1], :ntok] = 0 # Problematic

# Use element-wise operations:
arx_mask = (torch.arange(num_chunks, device=x.device) >= (arx[0] * arx[1])).to(dtype=x.
→dtype)
mask_i[:, :ntok] *= arx_mask.view(num_chunks, 1, 1) # Neuron-friendly
```

**Example:** File: test/unit/models/mlama/test\_vision\_encoder\_attention\_mask.py

```
# CPU version (problematic in Neuron):
def build_encoder_attention_mask_meta(x, ar, ntok, num_chunks, n_heads):
    masks = []
    for arx in ar:
        mask_i = torch.ones((num_chunks, x.shape[2], 1), dtype=x.dtype)
        mask_i[: arx[0] * arx[1], :ntok] = 0 # Problematic: array slicing + in-place
        # ...

# Neuron-friendly version:
def build_encoder_attention_mask(x, ar, ntok, num_chunks, n_heads):
    masks = []
    for arx in ar:
        mask_i = torch.ones((num_chunks, x.shape[2], 1), dtype=x.dtype, device=x.device)
        arx_mask = (torch.arange(num_chunks, device=x.device) >= (arx[0] * arx[1])).
→to(dtype=x.dtype)
        mask_i[:, :ntok] *= arx_mask.view(num_chunks, 1, 1) # Element-wise operation
        # ...
```

**Note:** This pattern applies to similar operations where array slicing and in-place modifications are used together. Consider using element-wise operations and avoiding in-place modifications for better Neuron compatibility.

## Performance Issues

### Skip model warmup during inference

**Issue:** You may observe slower performance for the first few inference requests, particularly on Trn2.

**Explanation:** By default, model warmup is disabled (`skip_warmup=True`) on Trn2 since warmup feature is not yet implemented for Trn2. This means the model needs to “warm up” naturally through actual inference requests, leading to slower performance during the initial requests.

**Solution:** There are approaches to ensure initial request performance:

1. Enable built-in warmup if your configuration supports it (on Inf2, Trn1):

```
neuron_config = NeuronConfig(
    tp_degree=32,
    batch_size=1,
    # skip_warmup=True is the default for Trn2 in release 2.23
    # skip_warmup=False is the default for Trn1, Inf2 in release 2.23
)
```

2. Implement manual warmup by sending dummy requests (on all instance types):

```
# Send a few dummy requests before serving real traffic
dummy_prompt = "This is a warmup request."
for _ in range(3): # Number of warmup iterations
    model.generate(
        prompt=dummy_prompt,
        max_new_tokens=32
    )
```

**Note:**

- When using vLLM for serving, the same initial performance impact applies if warmup is disabled.
- Use `--override-neuron-config '{"skip_warmup":false}'` to change the warmup setting

**Best Practice:**

- For production environments where initial latency is critical, test if your configuration supports built-in warmup.
- If built-in warmup isn't supported, implement manual warmup before serving real traffic.
- For development or non-latency-critical scenarios, the default configuration (warmup disabled) is sufficient.

## Other Common Issues

### Tensor Materialization During Tracing caused unexpected model behavior

**Issue:** Developers may inadvertently write code that forces tensor materialization during model tracing, leading to fixed computation paths and unexpected behaviors.

**Explanation:** When model logic depends on tensor values during the forward pass, the compiler may try to evaluate these values during tracing time. This “fixes” the computation path based on the initial values, resulting in a model that doesn't properly handle different runtime values.

Example of problematic code:

```
def forward(self, tensor):
    if tensor[0] == 1: # Forces tensor sync during tracing
        return tensor
    else:
        return tensor * 2
```

**Solution:** There are two debugging approaches to detect tensor materialization issues:

1. Enable warning messages:

```
import os
```

(continues on next page)

(continued from previous page)

```
# Set before model tracing
os.environ['PT_XLA_DEBUG_LEVEL'] = '2' # Will print warnings when tensor sync occurs
```

2. Force errors on tensor materialization:

```
import torch_xla

# Set before model tracing
torch_xla._XLAC._set_allow_execution(False) # Will raise an error if tensor sync is
↪ attempted
```

#### Best Practice:

- Avoid control flow that depends on tensor values during tracing. Instead, consider setting flags through configurations that should not change during runtime. See below example:

```
class TestModel(torch.nn.Module):
    def __init__(self, flag=1):
        super().__init__()
        # the flag should be pre-determined based on the model configuration
        # it should not be an input of the model during runtime
        self.flag = flag

    def forward(self, tensor):
        if self.flag:
            return tensor
        else:
            return tensor * 2
```

- If dynamic model path is required, consider using JIT inference (See: [Comparison of Traced Inference versus XLA Lazy Tensor Inference \(torch-neuronx\)](#))

### Input Data Type Handling for int64/fp64 due to compiler dtype compatibility

**Issue:** While you may be using 64-bit data types (int64/fp64) from tokenizers or other input sources, be aware that these are automatically converted to 32-bit types inside `ModelWrapper`.

**Explanation:** The Neuron compiler is optimized for 32-bit data types. To ensure optimal accuracy and compatibility, the model wrapper automatically converts 64-bit inputs (like those from Hugging Face tokenizers) to their 32-bit equivalents (int64 → int32, fp64 → fp32).

**Note:** No action is required from users as this conversion is handled automatically.

#### Best Practice:

- Continue using your tokenizers and input pipelines as normal
- Be aware that 64-bit inputs are automatically converted to 32-bit when using `ModelWrapper`
- If you're implementing custom pre-processing, using 32-bit types directly can be more efficient

This automatic conversion ensures consistent accuracy and compatibility with the Neuron compiler while maintaining ease of use with standard tokenizers and input pipelines.

- [NxD Inference Release Notes \(neuronx-distributed-inference\)](#)
- [Troubleshooting Guide for NxD Inference](#)

NxD Inference (where NxD stands for NeuronX Distributed) is an open-source PyTorch-based inference library that simplifies deep learning model deployment on AWS Inferentia and Trainium instances.

## NxDI Inference Overview and Setup

- *NxD Inference Overview*
- *NxD Inference Setup Guide*

## Developer Guide

- *NxD Inference Features Configuration Guide*
- *NxD Inference - Production Ready Models*
- *Onboarding models to run on NxD Inference*
- *vLLM User Guide for NxD Inference*
- *Testing modeling code with NxD Inference*
- *Migrating from NxD Core inference examples to NxD Inference*
- *Migrating from Transformers NeuronX to NeuronX Distributed(NxD) Inference*
- *LLM Inference Benchmarking guide*
- *Accuracy Evaluation of Models on Neuron Using Open Source Datasets*
- *Custom Quantization*
- *NxD Inference Weights Sharding Guide*
- *Disaggregated Inference [BETA]*

## API Reference Guide

- *NxD Inference API Reference*

## Tutorials

- *Tutorial: Deploying Llama3.1 405B (Trn2)*
- *Tutorial: Deploying Llama3.2 Multimodal Models*
- *Tutorial: Using Speculative Decoding to improve Llama-3.3-70B inference performance on Trn2 instances*
- *Tutorial: Scaling LLM Inference with Data Parallelism on Trn2*
- *Tutorial: Multi-LoRA serving for Llama-3.1-8B on Trn2 instances*
- *Tutorial: Using Speculative Decoding and Quantization to improve Llama-3.1-405B inference performance on Trn2 instances*
- *Tutorial: Evaluating Accuracy of Llama-3.1-70B on Neuron using open source datasets*
- *nxdi-trn2-llama3.3-70b-apc-tutorial*
- *Tutorial: Disaggregated Inference on Trn2 [BETA]*

## App Notes

- *Introducing NeuronX Distributed (NxD) Inference*
- *Parallelism Techniques for LLM Inference*

## Misc

- *NxD Inference Release Notes (neuronx-distributed-inference)*
- *Troubleshooting Guide for NxD Inference*

*This document is relevant for:* Inf2, Trn1, Trn2

## 3.3 NxD Core

NeuronX Distributed Core (NxD Core) is a package for supporting different distributed training/inference mechanism for Neuron devices. It would provide xla friendly implementations of some of the more popular distributed training/inference techniques. As the size of the model scales, fitting these models on a single device becomes impossible and hence we have to make use of model sharding techniques to partition the model across multiple devices. As part of this library, we enable support for Tensor Parallel sharding technique with other distributed library supported to be added in future.

*This document is relevant for:* Inf2, Trn1, Trn2

### 3.3.1 NeuronX Distributed Setup

*Install PyTorch Neuron on Trn1* to create a pytorch environment. It is recommended to work out of python virtual env so as to avoid package installation issues.

You can install the `neuronx-distributed` package using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.  
→amazonaws.com
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

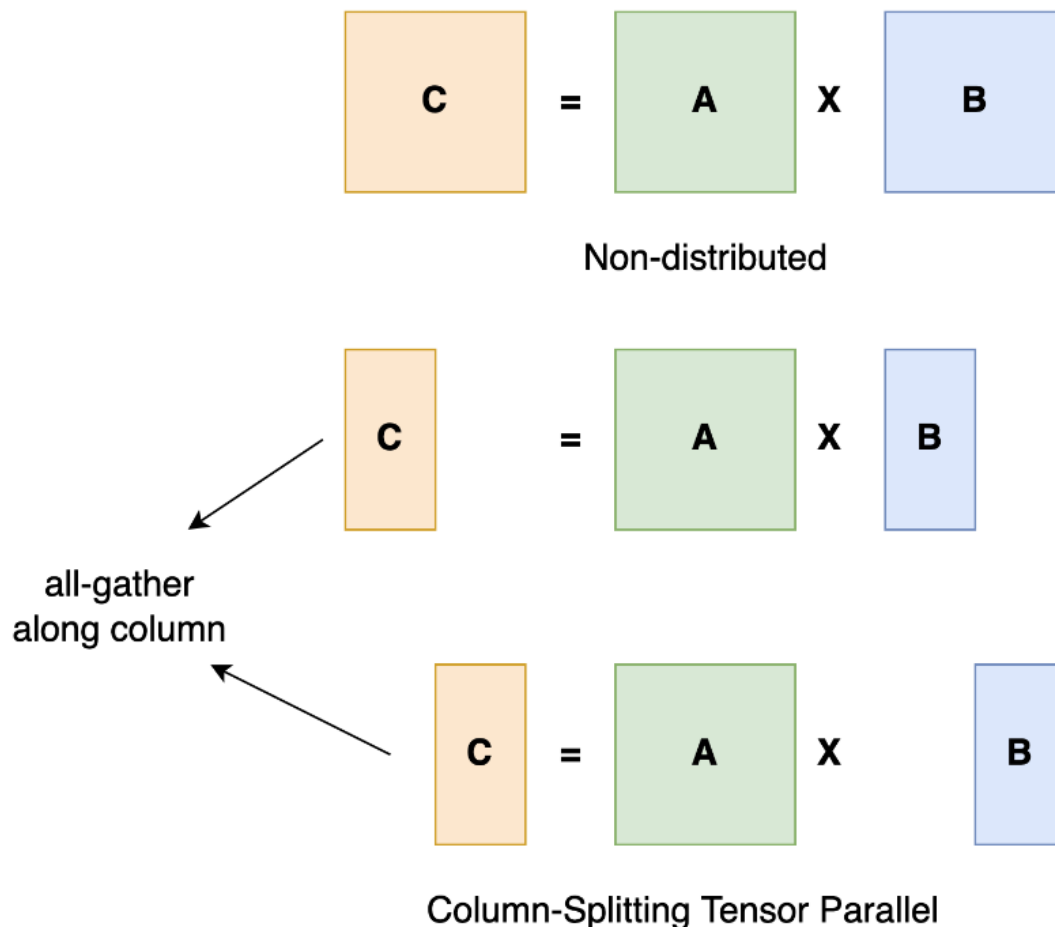
### 3.3.2 App Notes

*This document is relevant for:* Inf2, Trn1, Trn2

## Tensor Parallelism Overview

Tensor Parallelism is a technique in which a tensor is split into  $N$  chunks along a particular dimension such that each device only holds  $1/N$  chunk of the tensor. Computation is performed using this partial chunk so as to get partial output. These partial outputs are collected from all devices ensuring the correctness of the computation is maintained.

Taking a general matrix multiplication as an example, let's say we have  $C = AB$ . We can split  $B$  along the column dimension into  $[B_0 B_1 B_2 \dots B_n]$  and each device holds a column. We then multiply  $A$  with each column in  $B$  on each device, we will get  $[AB_0 AB_1 AB_2 \dots AB_n]$ . At this moment, each device still holds partial results, e.g. device rank 0 holds  $AB_0$ . To make sure the result is correct, we need to all-gather the partial result and concatenate the tensor along the column dimension. In this way, we are able to distribute the tensor over devices while making sure the computation flow remains correct.



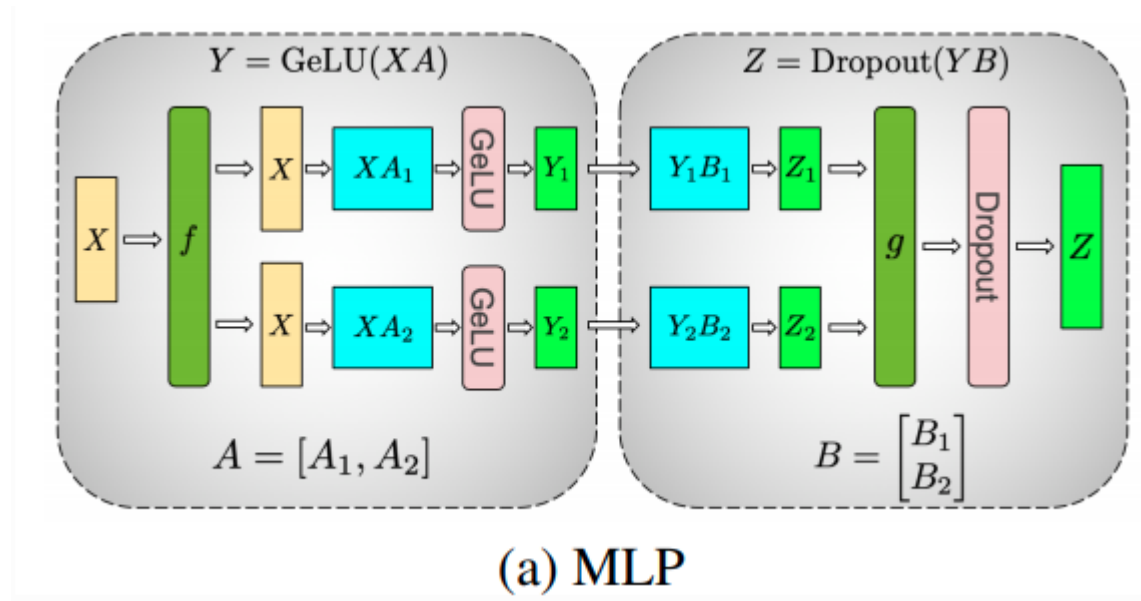
Tensor parallel illustration

Fig and TP explanation is borrowed from [https://colossalai.org/docs/concepts/paradigms\\_of\\_parallelism/#tensor-parallel](https://colossalai.org/docs/concepts/paradigms_of_parallelism/#tensor-parallel)

Similarly we can perform the partition along the row dimensions and create a RowParallel Linear layer. In RowParallel Linear layer, we partition the weight matrix along the row dimension. Let's say we have  $C = AB$ . We can split  $B$  along the row dimension into  $[B_0 B_1 B_2 \dots B_n]$  and each device holds a row. We then multiply each column of  $A$  on each device, we will get  $[A_0B_0 A_1B_1 A_2B_2 \dots A_nB_n]$ . At this moment, each device still holds partial results, e.g. device rank 0 holds  $A_0B_0$ . To make sure the result is correct, we need to all-reduce sum the partial result from all

devices to produce the final output.

Using this principle of sharded linear layers, we can construct MLPs of arbitrary depth until the need to operate on the whole output tensor, in which case we would have to construct the output but gathering it from all devices.



Here is an illustration from the Megatron-LM paper. In the above case, as you can see two linear layers are implemented using Column Parallel and Row Parallel linear layers, wherein the ColumnParallel Linear shards along the columns and then it is followed by RowParallel Linear layer which takes in parallel inputs (sharded outputs from ColumnParallelLinear). Consider the example shown in the above diagram,  $Z = (XA)B$ . In this case we split the first matrix multiplication over column dimension such that each device after first matrix multiplication holds partial result of  $Y_0=XA_0, Y_1=XA_1$  and so on. For the second matrix multiplication, we partition the weight matrix over row dimension and since the inputs are already columns sharded and we can multiply them to produce partial outputs. These outputs finally requires an all-reduce sum, since we want to sum up the single column\*row result.

Tensor Parallelism for Transformers:

A transformer block



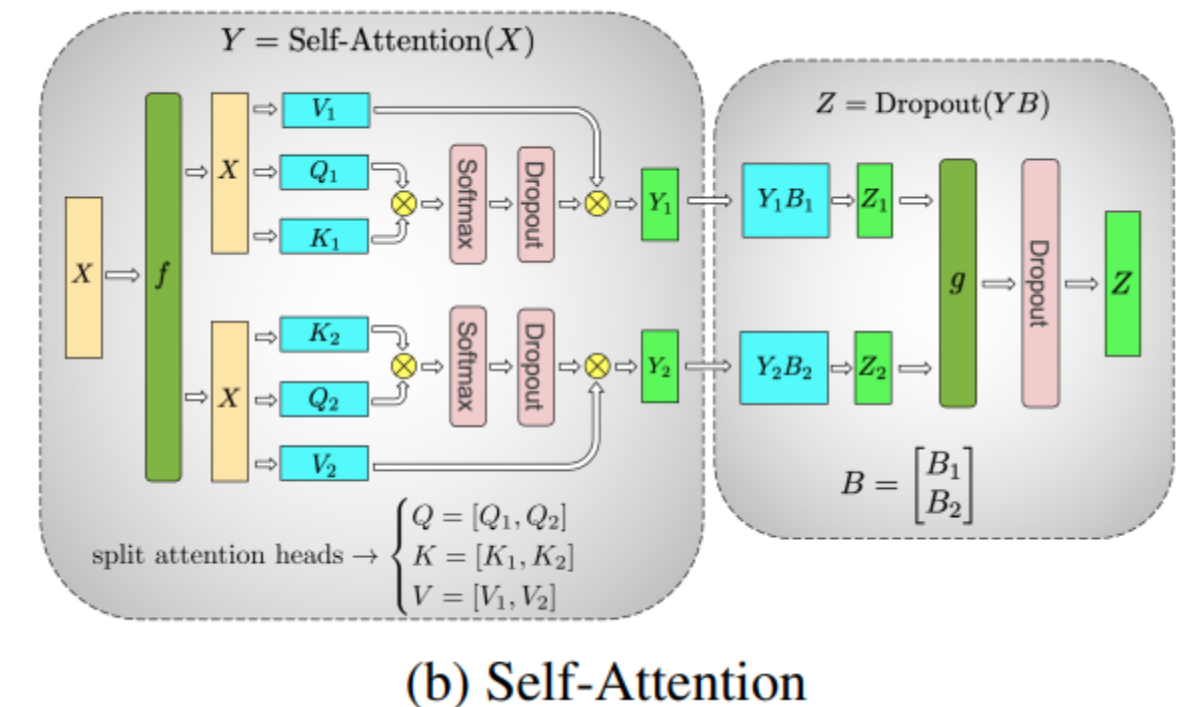


Fig: Taken from Megatron-LM paper.

As seen from the figure above, a simple self attention block has the QKV linear layer followed by MLP. Using the same Column and Row Parallel linear layers, we can partition the self-attention block across devices thereby reducing the memory footprint on each device, since each device now only holds partial parameters. This weight distribution strategy allows us to scale large model training across devices.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Pipeline Parallelism Overview

Pipeline parallelism is a technique used in deep learning model training to improve efficiency and reduce the training time of large neural networks. Currently NeuronxDistributed's pipeline parallelism is built specially for transformer based models, where each Neuron core will be assigned with a subset of transformer layers. Pipelining is a technique to achieve true parallelization in pipeline parallelism, by having the Neuron cores compute simultaneously on different data samples, and to overcome the performance loss due to sequential computation. When you use pipeline parallelism, training job is executed in a pipelined fashion over microbatches to maximize device usage.

## Model partitioning

In NeuronxDistributed, we use [Pytorch's FX](#) to trace the model and do partition on the FX IR. User simply needs to specify where to cut the pipeline stages, and our algorithm will cut the pipeline stages and assign the corresponding modules to each Neuron core automatically. Currently we require user to provide model partition decision but auto-partition will be supported in the future. Here is an example of simple partition with 5 linear layers

```
# original NN module
class MyModule(torch.nn.Module):
    def __init__(self):
```

(continues on next page)

(continued from previous page)

```

    super().__init__()
    self.linears = torch.nn.ModuleList([torch.nn.Linear(4, 4) for _ in range(5)])

    def forward(self, x):
        for lin in self.linears:
            x = lin(x)
        return x

m = MyModule()
gm = torch.fx.symbolic_trace(m)
print(gm)
"""
GraphModule(
  (linears): Module(
    (0): Linear(in_features=4, out_features=4, bias=True)
    (1): Linear(in_features=4, out_features=4, bias=True)
    (2): Linear(in_features=4, out_features=4, bias=True)
    (3): Linear(in_features=4, out_features=4, bias=True)
    (4): Linear(in_features=4, out_features=4, bias=True)
  )
)

def forward(self, x):
    linears_0 = getattr(self.linears, "0")(x); x = None
    linears_1 = getattr(self.linears, "1")(linears_0); linears_0 = None
    linears_2 = getattr(self.linears, "2")(linears_1); linears_1 = None
    linears_3 = getattr(self.linears, "3")(linears_2); linears_2 = None
    linears_4 = getattr(self.linears, "4")(linears_3); linears_3 = None
    return linears_4
"""

```

If user decide to cut the pipeline stage at the 3rd linear call, after partition there will be 2 submodules, where *submod\_0* contains first 3 linear layers and *submod\_1* contains last 2 linear layers.

```

After Split module
GraphModule(
  (submod_0): GraphModule(
    (linears_0): Linear(in_features=4, out_features=4, bias=True)
    (linears_1): Linear(in_features=4, out_features=4, bias=True)
    (linears_2): Linear(in_features=4, out_features=4, bias=True)
  )
  (submod_1): GraphModule(
    (linears_3): Linear(in_features=4, out_features=4, bias=True)
    (linears_4): Linear(in_features=4, out_features=4, bias=True)
  )
)

def forward(self, x):
    submod_0 = self.submod_0(x); x = None
    submod_1 = self.submod_1(submod_0); submod_0 = None
    return submod_1

```

## Pipeline Execution Schedule

Pipelining is based on splitting a mini-batch into microbatches, which are fed into the training pipeline one-by-one and follow an execution schedule defined by the library runtime. A microbatch is a smaller subset of a given training mini-batch. The pipeline schedule determines which microbatch is executed by which device for every time slot.

For example, depending on the pipeline schedule and the model partition, Neuron core  $i$  might perform (forward or backward) computation on microbatch  $b$  while Neuron core  $i+1$  performs computation on microbatch  $b+1$ , thereby keeping both Neuron cores active at the same time. An example taken from Megatron paper is showed as below

|          |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|
| Device 1 | 1 | 2 | 3 | 4 |   |   |   |   |   |   | 1 | 5 | 2 | 6 | 3 | 7 | 4 | 8 | 5 |   | 6 |   | 7 |   | 8 | 9 | 10 | 11 | 12 |    |    |    |    | 9  |    | 10 |
| Device 2 |   | 1 | 2 | 3 | 4 |   |   |   |   | 1 |   | 2 | 5 | 3 | 6 | 4 | 7 | 5 | 8 | 6 |   | 7 |   | 8 |   |   | 9  | 10 | 11 | 12 |    |    | 9  |    | 10 |    |
| Device 3 |   |   | 1 | 2 | 3 | 4 | 1 |   | 2 |   | 3 | 5 | 4 | 6 | 5 | 7 | 6 | 8 | 7 |   | 8 |   |   |   |   |   | 9  | 10 | 11 | 12 |    | 9  | 13 | 10 |    | 11 |
| Device 4 |   |   |   | 1 | 1 | 2 | 2 | 3 | 3 | 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 | 8 | 8 |   |   |   |   |   |   |   |    | 9  | 9  | 10 | 10 | 11 | 11 | 12 | 12 |    |

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## Activation Memory Reduction

There are three major contributors to high device memory utilization: *Parameters*, *Optimizer states* and *Activation Memory*. To reduce the size of parameter/optimizer states memory, one can use parallelism techniques like *Tensor-parallelism*, *Pipeline-parallelism* or *Zero1*. However, as the hidden size and sequence length increases, the size of the activation memory keeps growing linearly with hidden size and quadratically with sequence length.

The total activation memory without any parallelism comes to about:

$$\text{Activations memory per layer} = sbh \left( 34 + \frac{5as}{h} \right)$$

where,

- $a$ : Number of attention heads
- $b$ : microbatch size
- $h$ : hidden dimension size
- $s$ : sequence length

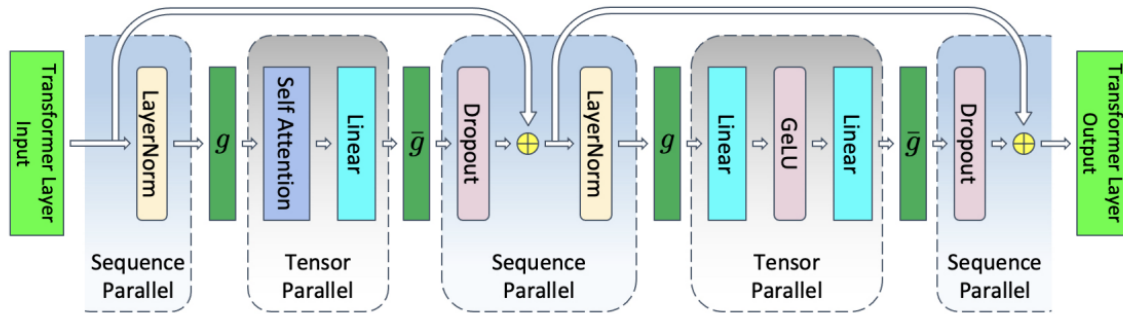
When we use tensor-parallelism, it not only helps to reduce the parameter and optimizer states size on each device, but it also helps to reduce the activation memory. For a transformer model, where we apply the tensor-parallel sharding on the attention block (more info [here](#)), the activation memory within the attention block also drops by a factor of tensor-parallel degree ( $t$ ). However, since the layernorms and dropouts (which are outside these attention blocks) are not parallelised and they are replicated on each device. These layernorms and dropouts are computationally inexpensive, however, they increase the overall activation memory on each device. Moreover, since we only parallelize within the attention block or within the MLP block ( $h \rightarrow 4h$  projection), the inputs to the QKV multiplies and the MLP are still unsharded. This overall adds to about  $10sbh$  of total activation memory. To reduce this activation memory, one can use 2 methods:

- [Sequence-Parallelism](#)
- [Activation Recomputation](#)

## Sequence Parallelism

Sequence-Parallelism was proposed by Shenggui and et.al . The authors propose to parallelize the compute along all the sequence dimension in an attempt to reduce increasing the memory pressure due to high sequence-lengths. Sequence-parallelism can be combined with tensor-parallelism to reduce the activation memory pressure due to increasing sequence-lengths.

Tensor-parallelism parallelizes the parts of the transformer which are computationally heavy, however, it leaves the layer-norms, dropouts and some MLP block intact. The activation memory for this block adds up to a factor of  $10sbh$ . Vijay Korthikanti et.al noticed that the compute in the non-tensor parallel region is independent in the sequence dimension. This property can be leveraged to shard the compute along the sequence dimension. The main advantage of sharding these non-tensor parallel block is reducing the activation memory. We can use the same tensor-parallel degree to partition, thereby reducing the overall activation memory by a factor of  $t$ . However, this partitioning comes at a cost. Since we are partitioning the non-tensor parallel region along sequence dimension, we have to collect the activations before we feed to the tensor-parallel block. This requires an introduction of **all-gather** collective operation which would gather the activations along the sequence dimension. Similarly, after the tensor-parallel block, since we would have to split the activations along the sequence dimension and distribute among the devices. Since, the tensor-parallel block in the transformer module already uses an all-reduce (Row-parallel linear layer used for MLP), we can replace the all-reduce operation with a **reduce-scatter** operation.



Ref: Reducing Activation Recomputation in Large Transformer Models

In the figure,  $g$  is all-gather operation and  $\bar{g}$  is the reduce-scatter operation.  $g$  and  $\bar{g}$  are conjugates and in the backward pass,  $\bar{g}$  becomes an all-gather operation and  $g$  becomes the reduce-scatter operation. At first glance, it appears that sequence-parallelism when combined with tensor-parallelism introduces an extra communication operation, however, in a ring all-reduce, the op is broken down into all-gather and reduce-scatter. Hence, the bandwidth required for sequence-parallelism is same as tensor-parallelism only. Hence, we are not losing out on compute but end up saving the activation memory per device. The final activation memory when sequence-parallelism is combined with tensor-parallelism:

$$\text{Activations memory per layer} = sbh \left( \frac{10}{t} + \frac{24}{t} + \frac{5as}{ht} \right) = \frac{sbh}{t} \left( 34 + \frac{5as}{h} \right)$$

## Activation Recomputation

The total required memory in the above equation can still be high as we increase the sequence length and hidden size. We would have to keep increasing the tensor-parallel degree to accommodate this requirement. Increasing the tensor-parallel degree might soon start producing diminishing returns as the model would start becoming bandwidth bottlenecked because of the extra collective communication operations. **Activation recomputation** can help to alleviate this problem. In this method, we recompute a part of the forward pass during the backward pass, thereby avoiding the need to save the activations during the forward pass. Activation recomputation is a trade-off between duplicate computation vs memory. It allows you to save on memory at the cost of extra recompute. This trade-off becomes valuable when we can fit larger models at the expense of recomputing forward pass activations.

Ideally one can recompute the entire forward pass, thereby resulting in an activation memory of  $2sbh$  per transformer layer. This method is called *Full-activation checkpointing*. This memory can further go down by a factor of  $t$  if we use tensor-parallelism. In the activation memory equation, we have a quadratic term of  $5abs^2$ . As the sequence length, this term will grow at a much faster rate. This quadratic term comes from the softmax computation. [Vijay Korthikanti et.al](#) propose *Selective activation checkpointing* where they only recompute the softmax and attention computation and thereby avoid saving the activations coming from softmax and attention computation. This completely gets rid of the quadratic term and brings down the activation memory per layer to  $34sbh/t$ . The LLama-7B example in [this tutorial](#) used selective activation checkpointing.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Context Parallelism Overview

Context parallelism (CP) is a technique used in deep learning model training to train large context models. CP parallelizes the processing of neural network activations across multiple devices by partitioning the input tensors along the sequence dimension. CP reduces the memory footprint and computational cost of processing long sequences. Unlike Sequence Parallelism (SP) that partitions the activations of specific layers, CP divides the activations of all layers.

The implementation of Context Parallelism in NxD leverages [Ring Attention](#). Ring Attention enables efficient communication between devices by organizing them in a ring topology, allowing tokens to attend to each other across devices without needing full attention computation on each device. This reduces memory overhead while extending the feasible context length beyond traditional transformer models.

For more details, refer to Context Parallelism in Megatron <[https://docs.nvidia.com/megatron-core/developer-guide/latest/api-guide/context\\_parallel.html](https://docs.nvidia.com/megatron-core/developer-guide/latest/api-guide/context_parallel.html)>\_

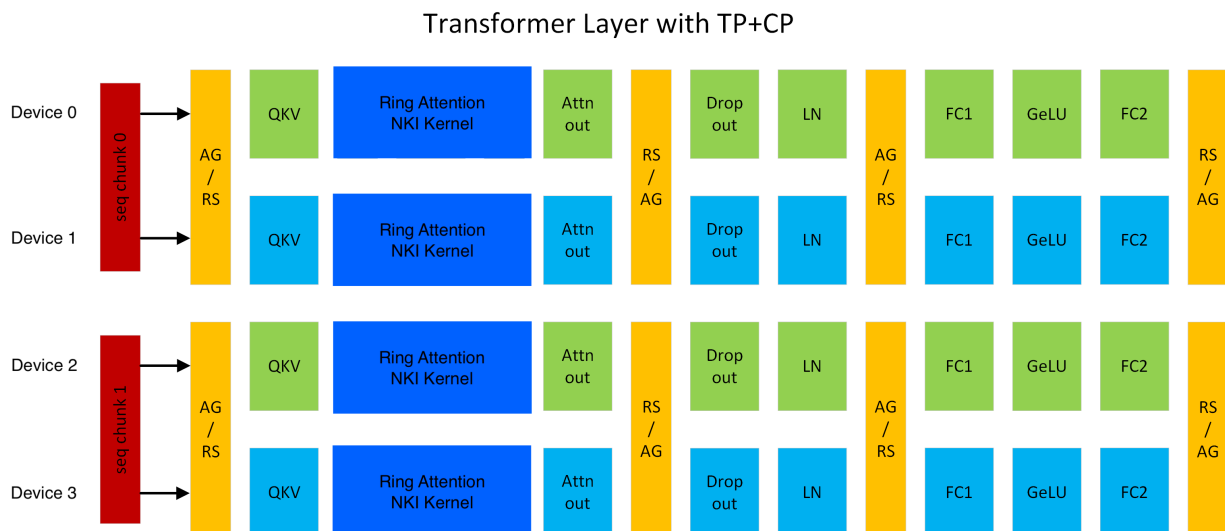


Fig: Context Parallelism in NxD (Figure adapted from [Megatron CP](#)). In NxD's TP implementation, we make use of All-Gather (AG), Reduce-Scatter (RS) collectives. Further CP is applied to all layers including LayerNorm (LN), Linear (LIN) and Fully-Connected (FC) layers. The figure shows a transformer layer running with TP2 and CP2. Assuming sequence length is 8K, each device processes 4K tokens. Device0 and Device2 form a CP group and exchange KV with each other; similarly, Device1 and Device3 form a CP group and exchange KV with each other. The collective communication to exchange KV is handled by NxD using approaches described in the [Ring Attention](#) paper.

*This document is relevant for:* Inf2, Trn1, Trn2

- [Tensor Parallelism Overview](#)

- *Pipeline Parallelism Overview*
- *Activation Memory Reduction*
- `context_parallelism_overview`

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### 3.3.3 API Reference Guide

*This document is relevant for:* Inf2, Trn1, Trn2

#### Distributed Strategies APIs

NeuronX Distributed Core (NxD Core) is XLA based library for distributed training and inference on Neuron devices. As part of this library, we support 3D parallelism: Tensor-Parallelism, Pipeline-Parallelism and Data-Parallelism. We also support Zero1 optimizer to shard the optimizer weights. To support tensor-parallelism on Neuron, we adopted the Apex Library built for CUDA devices. We modified the implementations to work with XLA. This document enlist the different APIs and modules provided by the library

##### Table of contents

- *Parallel Model State:*
  - *Initialize Model Parallelism:*
  - *Other helper APIs:*
- *Parallel Layers:*
  - *Parallel Embedding:*
  - *ColumnParallel Linear Layer:*
  - *RowParallel Linear Layer:*
  - *Padding Tensor-Parallel Layers*
  - *Loss functions:*
- *Pipeline parallelism:*
  - *Neuron Distributed Pipeline Model*
  - *Common used APIs*

#### Parallel Model State:

##### Initialize Model Parallelism:

```
def neuronx_distributed.parallel_state.initialize_model_parallel(  
    tensor_model_parallel_size=1,  
    pipeline_model_parallel_size=1,  
)
```

This module would initialize the distributed model training and allows users to set the number of tensor\_parallel world size.

Parameters:

- `tensor_model_parallel_size`: This should set the number of tensor parallel workers. Note the default value is set to 1
- `pipeline_model_parallel_size`: This should set the number of pipeline parallel workers. Note the default value is set to 1

### Other helper APIs:

- `neuronx_distributed.parallel_state.get_data_parallel_size()`: Returns the data parallel world size depending on the number of global workers and tensor parallel workers.
- `neuronx_distributed.parallel_state.get_tensor_model_parallel_size()`: Returns the tensor parallel world size.
- `neuronx_distributed.parallel_state.get_tensor_model_parallel_rank()`: Returns the rank of the worker within the tensor parallel group
- `neuronx_distributed.parallel_state.get_pipeline_model_parallel_size()`: Returns the pipeline parallel world size.
- `neuronx_distributed.parallel_state.get_pipeline_model_parallel_rank()`: Returns the rank of the worker within the pipeline parallel group
- `neuronx_distributed.parallel_state.get_data_parallel_rank()`: Returns the rank of the worker in the data parallel group.
- `neuronx_distributed.parallel_state.get_data_parallel_group(as_list=False)`: Returns the data parallel group after taking into account the tensor parallel size and the global world size. `as_list` argument when set to True, would return the group as a List[List] otherwise it would return a torch.distributed.group.
- `neuronx_distributed.parallel_state.get_tensor_model_parallel_group(as_list=False)`: Returns the tensor parallel group after taking into account the tensor parallel size and the global world size. `as_list` argument when set to True, would return the group as a List[List] otherwise it would return a torch.distributed.group.
- `neuronx_distributed.parallel_state.get_pipeline_model_parallel_group(as_list=False)`: Returns the pipeline parallel group after taking into account the pipeline parallel size and the global world size. `as_list` argument when set to True, would return the group as a List[List] otherwise it would return a torch.distributed.group.
- `move_model_to_device(model, device)`: This api moves the model to device by preserving tensor parallel attributes.

### Parallel Layers:

Majority of parameters within the transformer based model reside in the Embedding and Linear layers. Hence, to reduce the number of parameters on a single device because of these layers, we provided sharded Embedding and Linear layers.

### Parallel Embedding:

```
class neuronx_distributed.parallel_layers.ParallelEmbedding(  
    num_embeddings, embedding_dim, init_method=init.normal_,  
    dtype=torch.float32, device=None)
```

This module is intended to replace `torch.nn.Embedding`. In cases where the vocab size is too large, we can shard the Embedding table across workers. Note: The embedding table would be sharded across all the tensor-parallel workers.

Parameters:

- `num_embeddings` (int) : size of the dictionary of embeddings
- `embedding_dim` (int) : the size of each embedding vector
- `init_method`: (`torch.nn.init`) : Initialization function for the embedding weights.
- `dtype`: (dtype) : Datatype for the weights
- `device`: (`torch.device`) : Device to initialize the weights on. By default, the weights would be initialized on CPU

### ColumnParallel Linear Layer:

```
class neuronx_distributed.parallel_layers.ColumnParallelLinear(  
    input_size, output_size, bias=True, gather_output=True,  
    sequence_parallel_enabled=False, dtype=torch.float32, device=None)
```

This module would perform a Column wise partition of the weight matrix. Linear layer is defined as  $Y = XA + b$ , here  $A$  is parallelized along second dimension as  $A = [A_1, A_2 \dots A_p]$ . Note: This layer is designed to operate on 3-dimensional inputs.

Parameters:

- `input_size`: (int) : First dimension of the weight matrix
- `output_size`: (int) : Second dimension of the weight matrix
- `bias`: (bool): If set to True, bias would be added
- `gather_output`: (bool) : If true, call all-gather on output and make  $Y$  available to all Neuron devices, otherwise, every Neuron device will have its output which is  $Y_i = XA_i$
- **`sequence_parallel_enabled`: (bool)**  
[When sequence-parallel is enabled, it would] gather the inputs from the sequence parallel region and perform the forward and backward passes
- `dtype`: (dtype) : Datatype for the weights
- `device`: (`torch.device`) : Device to initialize the weights on. By default, the weights would be initialized on CPU



### RowParallel Linear Layer:

```
class neuronx_distributed.parallel_layers.RowParallelLinear(
    input_size, output_size, bias=True, input_is_parallel=False,
    sequence_parallel_enabled=False, dtype=torch.float32, device=False
)
```

The linear layer is defined as  $Y = XA + b$ . A is parallelized along its first dimension and X along its second. Note: This layer is designed to operate on 3-dimensional inputs.

Parameters:

- **input\_size**: (int) : First dimension of the weight matrix
- **output\_size**: (int) : Second dimension of the weight matrix
- **bias**: (bool) : If set to True, bias would be added
- **input\_is\_parallel**: (bool) : If true, we assume that the input is already split across the Neuron devices and we do not split again. This is useful when we have a ColumnParallel Layer just before the Row Parallel layer
- **sequence\_parallel\_enabled**: (bool) : When sequence-parallel is enabled, it would gather the inputs from the sequence parallel region and perform the forward and backward passes
- **dtype**: (dtype) : Datatype for the weights
- **device**: (torch.device) : Device to initialize the weights on. By default, the weights would be initialized on CPU

### Padding Tensor-Parallel Layers

```
def neuronx_distributed.parallel_layers.pad.pad_model(
    model, tp_degree, n_heads, wrapped_classes=(), pad_hook_fn=None)
```

Pads a generic model to function to a desired tensor parallelism degree by padding the number of attention heads. Returns the original model modified with padding. Uses 1-axis padding strategy: pads the sharded dim of the ParallelLinear layers to the size it would have been for the padded number of heads.

Parameters:

- **model** (torch.nn.Module) : model to be padded
- **tp\_degree** (int) : tensor parallel degree
- **n\_heads** (int)  
[the number of heads the given model to be padded has. This can] typically be found in the config
- **wrapped\_classes** (Tuple[any], \*optional\*, defaults to `()`)  
[tuple of classes] (and their submodules) which should be padded
- **pad\_hook\_fn** (Callable[any, float], \*optional\*, defaults to `None`)  
[a hook] function that is called whenever encountering a class to pad. Receives an instance of the class to pad and the `tgt_src_ratio` (`num_heads_padded / num_heads`) as its argument

Usage:

When modifying the Attention layer, typically you must divide by TP degree like so:

```
self.num_heads = neuronx_dist_utils.divide(self.num_heads, get_tensor_model_
↪parallel_size())
```

This line must be modified like so:

```
self.num_heads = neuronx_dist_utils.divide(
    self.num_heads + get_number_of_extra_heads(self.num_heads, get_tensor_model_
    ↪parallel_size()),
    get_tensor_model_parallel_size())
```

Then, after initializing the model, you must call this wrapper:

```
model = get_model(config=desired_config)
model = pad_model(model, tp_degree=32, desired_config.num_heads) # Use the model,
    ↪as desired after this point
```

You can specify a specific layer or class for your model to pad, so you aren't unnecessarily padding. Typically, this layer will be your Attention layer

```
model = pad_model(model, tp_degree=32, desired_config.num_heads, wrapped_
    ↪classes=[MyAttention])
```

You can also specify a `pad_hook_fn`, to be called whenever encountering an instance of `wrapped_class`, passing in said instance as a parameter, along with the `tgt_src_ratio` (`num_heads_padded / num_heads`).

```
def my_hook(attention_to_pad, tgt_src_ratio):
    attention_to_pad.split_size = int(model.split_size * tgt_src_ratio)
    model = pad_model(
        model,
        tp_degree=32,
        desired_config.num_heads,
        wrapped_classes=[MyAttention],
        pad_hook_fn=my_hook
    )
```

## Loss functions:

When you shard the final MLP layer using tensor-parallelism, instead of recollecting all the outputs from each TP rank, we can use the `ParallelCrossEntropy` loss function. This function would take the parallel logits produced by final parallel MLP and produce a loss by taking into account that the logits are sharded across multiple workers.

```
def neuronx_distributed.parallel_layers.loss_functions.parallel_cross_entropy(
    parallel_logits, labels, label_smoothing=0.0)
```

Parameters:

- `parallel_logits` (Tensor) : Sharded logits from the previous MLP
- `labels` (Tensor) : Label for each token. Labels should not be sharded, and the `parallel_cross_entropy` would take care of sharding the labels internally
- `label_smoothing` (float) : A float in [0.0, 1.0]. Specifies the amount of smoothing when computing the loss, where 0.0 means no smoothing

## Pipeline parallelism:

### Neuron Distributed Pipeline Model

```
class NxDPModel(
    module: torch.nn.Module,
    transformer_layer_cls: Optional[Any] = None,
    num_microbatches: int = 1,
    virtual_pipeline_size: int = 1,
    output_loss_value_spec: Optional[Union[Dict, Tuple]] = None,
    return_mb_loss: bool = False,
    broadcast_and_average_loss: bool = False,
    pipeline_cuts: Optional[List[str]] = None,
    input_names: Optional[List[str]] = None,
    leaf_module_cls: Optional[List[Any]] = None,
    autowrap_functions: Optional[Tuple[ModuleType]] = None,
    autowrap_modules: Optional[Tuple[Callable, ...]] = None,
    tracer_cls: Optional[Union[str, Any]] = None,
    param_init_fn: Optional[Any] = None,
    trace_file_path: Optional[str] = None,
    use_zero1_optimizer: bool = False,
    auto_partition: Optional[bool] = False,
    deallocate_pipeline_outputs: bool = False,
)
```

Parameters:

- **module:** Module to be distributed with pipeline parallelism
- **transformer\_layer\_cls:** The module class of transformer layers
- **num\_microbatches:** Number of pipeline microbatches
- **virtual\_pipeline\_size:** Virtual pipeline size if greater than 1 we will use the interleaved pipeline schedule.
- **output\_loss\_value\_spec:**  
The output\_loss\_value\_spec value can be specified to disambiguate which value in the output of *forward* is the loss value on which NxDPModel should apply backpropagation. For example, if your forward returns a tuple (loss, model\_out), you can specify output\_loss\_value\_spec=(True, False). Or, if your forward returns a dict {'loss': loss\_value, 'model\_out': model\_out}, you can specify output\_loss\_value\_spec={'loss': True, 'model\_out': False} referred from [this](#)
- **return\_mb\_loss:** Whether return a list of loss for all microbatches
- **broadcast\_and\_average\_loss:** Whether to broadcast loss to all PP ranks and average across dp ranks, when set to True return\_mb\_loss must be False
- **pipeline\_cuts:** A list of layer names that will be used to annotate pipeline stage boundaries
- **input\_names:** The input names that will be used for tracing, which will be the same as the model inputs during runtime.
- **leaf\_module\_cls:** A list of module classes that should be treated as leaf nodes during tracing. Note transformer layer class will be by default treat as leaf nodes.
- **autowrap\_modules: (symbolic tracing only)**  
Python modules whose functions should be wrapped automatically without needing to use fx.wrap(). reference [here](#)

- **autowrap\_functions:** (symbolic tracing only)

Python functions that should be wrapped automatically without needing to use `fx.wrap()`. reference [here](#)

- **tracer\_cls:** User provided tracer class for symbolic tracing. It can be “hf”, “torch” or any tracer class user created.

- **param\_init\_fn:**

Function used to initialize parameters. This is useful if user wants to use meta device to do delayed parameter initialization. `param_init_fn` should take a module as input and initialize the parameters that belongs to this module only (not for submodules).

- **use\_zero1\_optimizer:** Whether to use the zero1 optimizer. When setting to True the gradient average will be handed over.

- **auto\_partition:**

Boolean to indicate whether to use `auto_partition` for the model. When set to True, the pipeline cuts used as the pipeline stage boundaries to partition the model are automatically determined. When set to True, the `pipeline_cuts` parameter should not be set. The `pipeline_cuts` are chosen on the basis of the transformer layer names.

- **deallocate\_pipeline\_outputs:**

Whether to deallocate the pipeline outputs after send. After send the output tensor is only useful for its `‘grad_fn’` field, and not its `‘data’`.

## Common used APIs

```
NxDPPModel.run_train(**kwargs)
```

Train the model with PP schedule, which will run both forward and backward in a PP manner. The kwargs should be the same as the `input_names` provided to the trace function. Will output the loss that provided by user from `output_loss_value_spec`.

```
NxDPPModel.run_eval(**kwargs)
```

Eval the model with PP schedule, which will run forward only. The kwargs should be the same as the `input_names` provided to the trace function. Will output the loss that provided by user from `output_loss_value_spec`.

```
NxDPPModel.local_named_parameters(**kwargs)
```

The parameters that are local to this PP rank. This must be called after the model is partitioned.

```
NxDPPModel.local_named_modules(**kwargs)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Training APIs

### Table of contents

- *Neuronx-Distributed Training APIs:*
  - *Initialize NxD Core config:*
  - *Initialize NxD Core Model Wrapper:*
  - *Initialize NxD Core Optimizer Wrapper:*
  - *Enable LoRA fine-tuning:*
  - *Save Checkpoint:*
  - *Load Checkpoint:*
- *Modules:*
  - *GQA-QKV Linear Module:*
- *Checkpointing:*
  - *Save Checkpoint:*
  - *Load Checkpoint*
  - *Gradient Clipping:*
  - *Neuron Zero1 Optimizer:*
- *Neuron PyTorch-Lightning*
  - *Neuron Lightning Module*
  - *Neuron XLA Strategy*
  - *Neuron XLA Precision Plugin*
  - *Neuron TQDM Progress Bar*
  - *Neuron TensorBoard Logger*

### Neuronx-Distributed Training APIs:

In Neuronx-Distributed, we provide a series of APIs under *neuronx\_distributed* directly that helps user to apply optimizations in NxD Core easily. These APIs cover configuration, model/optimizer initialization and saving/loading checkpoint.

**Initialize NxD Core config:**

```
def neuronx_distributed.trainer.neuronx_distributed_config(
    tensor_parallel_size=1,
    pipeline_parallel_size=1,
    pipeline_config=None,
    optimizer_config=None,
    activation_checkpoint_config=None,
    pad_model=False,
    sequence_parallel=False,
    model_init_config=None,
    lora_config=None,
)
```

This method initializes NxD Core training config and initialize model parallel. This config maintains all optimization options of the distributed training, and it's a global config (the same for all processes).

Parameters:

- **tensor\_parallel\_size** (int) : Tensor model parallel size. Default: 1.
- **pipeline\_parallel\_size** (int) : Pipeline model parallel size. Default: 1.
- **pipeline\_config** (dict)
  - [Pipeline parallel config. For details please refer to] [pipeline parallel guidance](#). Default: None.
- **optimizer\_config** (dict) : Optimizer config. Default: {"zero\_one\_enabled": False, "grad\_clipping": True, "max\_grad\_norm": 1.0}.
  - Enable ZeRO-1 by setting zero\_one\_enabled to True.
  - Enable grad clipping by setting grad\_clipping to True.
  - Change maximum grad norm value by setting max\_grad\_norm.
- **activation\_checkpoint\_config** (str of torch.nn.Module)
  - [Activation checkpoint config.] accept value: "full", None, or any torch.nn.Module. When set to full, regular activation checkpoint enabled (every transformer layer will be re-computed). When set to None, activation checkpoint disabled. When set to any torch.nn.Module, selective activation checkpoint enabled, any provided module will be re-computed. Default: None.
- **pad\_model** (bool) : Whether to pad attention heads of model. Default: False.
- **sequence\_parallel** (bool) : Whether to enable sequence parallel. Default: False.
- **model\_init\_config** (dict) : Model initialization config. Default: {"sequential\_move\_factor": 11, "meta\_device\_init": False, "param\_init\_fn": None}.
- **lora\_config**: LoRA configuration. Default: None with LoRA disabled.
- **sequential\_move\_factor**: num of processes instantiating model on host at the same time.
  - This is done to avoid the host OOM. Range: 1-32.
- **meta\_device\_init**: whether to initialize model on meta device.
- **param\_init\_fn**: method that initialize parameters of modules, should be provided when param\_init\_fn is True.

### Initialize NxD Core Model Wrapper:

```
def neuronx_distributed.trainer.initialize_parallel_model(nxd_config, model_fn, *model_
↳args, **model_kwargs)
```

This method initialize NxD Core model wrapper, return a wrapped model that can be used as a regular `torch.nn.Module`, while has all the model optimizations in config applied. This wrapper is designed to hide the complexity of optimizations such as pipeline model parallel, so that users can simply use the wrapped model as the unwrapped version.

Parameters:

- `nxd_config` (dict): config generated by `neuronx_distributed_config`.
- `model_fn` (callable): user provided function to get the model for training.
- `model_args` and `model_kwargs`: arguments that will be passed to `model_fn`.

Model wrapper class and its methods:

```
class neuronx_distributed.trainer.model.NxDModel(torch.nn.Module):
    def local_module(self):
        # return the unwrapped local module

    def run_train(self, *args, **kwargs):
        # method to run one iteration, when pipeline parallel enabled,
        # user have to use this instead of forward+backward

    def named_parameters(self, *args, **kwargs):
        # only return parameters on local rank.
        # same for `parameters`, `named_buffers`, `buffers`

    def named_modules(self, *args, **kwargs):
        # only return modules on local rank.
        # same for `modules`, `named_children`, `children`
```

**Note:** As a short cut, users can call `model.config` or `model.dtype` from wrapped model if original model is hugging face transformers pre-trained model.

### Initialize NxD Core Optimizer Wrapper:

```
def neuronx_distributed.trainer.initialize_parallel_optimizer(nxd_config, optimizer_
↳class, parameters, **defaults)
```

This method initialize NxD Core optimizer wrapper, return a wrapped optimizer that can be used as a regular `torch.optim.Optimizer`, while has all the optimizer optimizations in config applied.

This optimizer wrapper is inherited from `torch.optim.Optimizer`. It takes in the `nxd_config` and configures the optimizer to work with different distributed training regime.

The `step` method of the wrapped optimizer contains necessary all-reduce operations and grad clipping. Other methods and variables work the same as the unwrapped optimizer.

Parameters:

- `nxd_config` (dict): config generated by `neuronx_distributed_config`.
- `optimizer_class` (Type[torch.optim.Optimizer]): optimizer class to create the optimizer.
- `parameters` (iterable): parameters passed to the optimizer.
- `defaults`: optimizer options that will be passed to the optimizer.

### Enable LoRA fine-tuning:

LoRA model wrapper

```
class LoRAModel(module, LoraConfig)
```

Parameters:

- `module` (torch.nn.Module): Module to be wrapped with LoRA
- `LoraConfig`: The LoRA configuration defined in `neuronx_distributed.modules.lora.LoraConfig`

The flags in `LoraConfig` to initialize LoRA adapter:

- `enable_lora` (bool): Enable LoRA fine-tuning.
- `lora_rank` (int): The rank of LoRA adapter. A small LoRA rank reduces the memory footprint during fine-tuning, but it may harm the model quality.
- `lora_alpha` (float): The alpha parameter for LoRA scaling, i.e., scaling LoRA weights against base model weights.
- `lora_dropout` (float): The dropout probability for LoRA layers.
- `bias` (str): Bias type for LoRA. Can be `none`, `all` or `lora_only`.
- `target_modules` (List[str]): The names of the modules that need LoRA.
- `use_rslora` (bool): If True, uses Rank-Stabilized LoRA, which sets the adapter scaling factor to `lora_alpha/math.sqrt(lora_rank)`.
- `init_lora_weights` (str): Weights initialization of LoRA adapter. Can be `default` (initialized with `torch.nn.init.kaiming_uniform_()`) or `gaussian` (initialized with `torch.nn.init.normal_()`).

### Usage:

We first define the LoRA configuration for fine-tuning. Suppose the target modules is `[q_proj, v_proj, k_proj]`, it indicates that LoRA will be applied to modules whose name includes any of the keywords. An example is

```
lora_config = neuronx_distributed.modules.lora.LoraConfig(
    enable_lora=True,
    lora_rank=16,
    lora_alpha=32,
    lora_dropout=0.05,
    bias="none",
    target_modules=["q_proj", "v_proj", "k_proj"],
)
```

You can enable LoRA fine-tuning like below



```
nxd_config = neuronx_distributed.neuronx_distributed_config(
    ...
    lora_config=lora_config,
)
model = neuronx_distributed.initialize_parallel_model(nxd_config, ...)
```

Then the NxD model will be initialized with LoRA adapter enabled.

### Save Checkpoint:

Method to save checkpoint, return None.

This method saves checkpoints for model, optimizer, scheduler and user contents sequentially. Model states are saved on data parallel rank-0 only. When ZeRO-1 optimizer is not turned on, optimizer states are also saved like this; while when ZeRO-1 optimizer is turned on, states are saved on all ranks. Scheduler and user contents are saved on master rank only. Besides, users can use `use_xser=True` to boost saving performance and avoid host OOM. It's achieved by saving tensors one by one simultaneously and keeping the original data structure. However, the resulted checkpoint cannot be loaded using load api of PyTorch. Users can also use `async_save=True` to further boost saving performance. It's achieved by saving tensors in separate processes along with computation. Setting `async_save` to true will result in more host memory being used, therefore increase the risk of application crash due to system ran out of memory.

```
def neuronx_distributed.trainer.save_checkpoint(
    path,
    tag="",
    model=None,
    optimizer=None,
    scheduler=None,
    user_content=None,
    num_workers=8,
    use_xser=False,
    num_kept_ckpts=None,
    async_save=False,
    zero1_optimizer=False
)
```

Parameters:

- `path` (str): path to save the checkpoints.
- `tag` (str): tag to save the checkpoints.
- `model` (torch.nn.Module): model to save, optional.
- `optimizer` (torch.optim.Optimizer): optimizer to save, optional.
- `scheduler`: scheduler to save, optional.
- `user_content`: user contents to save, optional.
- **`num_workers` (int): num of processes saving data on host at the same time.**  
This is done to avoid the host OOM, range: 1-32.
- **`use_xser` (bool): whether to use torch-xla serialization. When enabled, `num_workers` will be ignored and maximum num of workers will be used. Default: False.**
- `num_kept_ckpts` (int): number of checkpoints to keep on disk, optional. Default: None.
- `async_save` (bool): whether to use asynchronous saving method. Default: False.

- `zero1_optimizer` (bool): : whether the optimizer state is from a zero1 optimizer, used when optimizer is a dict

### Save LoRA Checkpoint:

NxD also uses `neuronx_distributed.trainer.save_checkpoint()` to save LoRA models, but it can set `save_lora_base` and `merge_lora` in `LoraConfig` to specify how to save LoRA checkpoint. There are three modes for LoRA checkpoint saving:

- `save_lora_base=False`, `merge_lora=False`: Save the LoRA adapter only.
- `save_lora_base=True`, `merge_lora=False`: Save both the base model and the LoRA adapter separately.
- `save_lora_base=True`, `merge_lora=True`: Merge the LoRA adapter into the base model and then save the base model.

Other than the adapter, NxD also needs to save the LoRA configuration file for LoRA loading. The configuration can be saved into the same checkpoint with the adapter, or saved as a separately json file.

- `save_lora_config_adapter` (bool): If False, save the configuration file as a separately json file.

Note that if LoRA configuration file is saved separately, it is named as `lora_adapter/adapter_config.json`.

A configuration example to save the LoRA adapter only is

```
lora_config = neuronx_distributed.modules.lora.LoraConfig(  
    ...  
    save_lora_base=False,  
    merge_lora=False,  
    save_lora_config_adapter=True,  
)
```

### Load Checkpoint:

Method to load checkpoint saved by `save_checkpoint`, return user contents if exists otherwise None. If `tag` not provided, will try to use the newest tag tracked by `save_checkpoint`.

Note that the checkpoint to be loaded must have the same model parallel degrees as in current use, and if ZeRO-1 optimizer is used, must use the same data parallel degrees.

```
def neuronx_distributed.trainer.load_checkpoint(  
    path,  
    tag=None,  
    model=None,  
    optimizer=None,  
    scheduler=None,  
    num_workers=8,  
    strict=True,  
)
```

Parameters:

- `path` (str): path to load the checkpoints.
- `tag` (str): tag to load the checkpoints.
- `model` (torch.nn.Module): model to load, optional.
- `optimizer` (torch.optim.Optimizer): optimizer to load, optional.

- `scheduler`: scheduler to load, optional.
- `num_workers` (int): num of processes loading data on host at the same time.

This is done to avoid the host OOM, range: 1-32. - `strict` (bool): whether to use strict mode when loading model checkpoint. Default: True.

### Load LoRA Checkpoint:

NxD loads LoRA checkpoints by setting flags in `LoraConfig`.

- `load_lora_from_ckpt`: Resumes the checkpoint process.
- `lora_save_dir`: Load LoRA checkpoint from the specified folder
- `lora_load_tag`: Load the LoRA checkpoint with the specified tag

An example is:

```
lora_config = LoraConfig(
    enable_lora=True,
    load_lora_from_ckpt=True,
    lora_save_dir=checkpoint_dir, # checkpoint path
    lora_load_tag=tag, # sub-directory under checkpoint path
)
nxd_config = nxd.neuronx_distributed_config(
    ...
    lora_config=lora_config,
)
model = nxd.initialize_parallel_model(nxd_config, ...)
```

The NxD model will be initialized with LoRA enabled and LoRA weights loaded. LoRA-related configurations are the same as the LoRA adapter checkpoint.

### Sample usage:

```
import neuronx_distributed as nxd

# create config
nxd_config = nxd.neuronx_distributed_config(
    tensor_parallel_size=8,
    optimizer_config={"zero_one_enabled": True, "grad_clipping": True, "max_grad_norm": 1.0},
)

# wrap model
model = nxd.initialize_parallel_model(nxd_config, get_model)

# wrap optimizer
optimizer = nxd.initialize_parallel_optimizer(nxd_config, AdamW, model.parameters(), lr=1e-3)

...
(training loop):
    loss = model.run_train(inputs)
    optimizer.step()
...
```

(continues on next page)

(continued from previous page)

```

# loading checkpoint (auto-resume)
user_content = nxd.load_checkpoint(
    "ckpts",
    model=model,
    optimizer=optimizer,
    scheduler=scheduler,
)
...
# saving checkpoint
nxd.save_checkpoint(
    "ckpts",
    nxd_config=nxd_config,
    model=model,
    optimizer=optimizer,
    scheduler=scheduler,
    user_content={"total_steps": total_steps},
)

```

## Modules:

### GQA-QKV Linear Module:

```

class neuronx_distributed.modules.qkv_linear.GQAQKVColumnParallelLinear(
    input_size, output_size, bias=True, gather_output=True,
    sequence_parallel_enabled=False, dtype=torch.float32, device=None, kv_size_
    multiplier=1, fuse_qkv=True)

```

This module parallelizes the Q,K,V linear projections using ColumnParallelLinear layers. Instead of using 3 different linear layers, we can replace it with a single QKV module. In case of GQA module, the number of Q attention heads are  $N$  times more than the number of K and V attention heads. The K and V attention heads are replicated after projection to match the number of Q attention heads. This helps to reduce the K and V weights and is useful especially during inference. However, in case of training these modules, it restricts the tensor-parallel degree that can be used, since the attention heads should be divisible by tensor-parallel degree. Hence, to mitigate this bottleneck, the *GQAQKVColumnParallelLinear* takes in a *kv\_size\_multiplier* argument. The module would replicate the K and V weights *kv\_size\_multiplier* times thereby allowing you to use higher tensor-parallel degree. Note: here instead of replicating the projection  $N/tp\_degree$  times, we end of replicating the weights *kv\_size\_multiplier* times. This would produce the same result, allow you to use higher *tp\_degree* degree, however, it would result in extra memory getting consumed.

Parameters:

- **input\_size:** (int) : First dimension of the weight matrix
- **output\_sizes:** (List[int]) : A list of second dimension of the Q and K/V weight matrix
- **bias:** (bool): If set to True, bias would be added
- **gather\_output:** (bool)  
[If true, call all-gather on output and make Y available to all] Neuron devices, otherwise, every Neuron device will have its output which is  $Y_i = XA_i$
- **sequence\_parallel\_enabled:** (bool)  
[When sequence-parallel is enabled, it would gather] the inputs from the sequence parallel region and perform the forward and backward passes

- `init_method: (torch.nn.init)`: Initialization function for the Q and K/V weights.
- `dtype: (dtype)`: Datatype for the weights
- **`device: (torch.device)`**  
[Device to initialize the weights on. By default, the weights] would be initialized on CPU
- `kv_size_multiplier: (int)`: Factor by which the K and V weights would be replicated along the first dimension
- `fuse_qkv: (bool)`: When `fuse_qkv` is enabled, a single fused tensor is used for QKV. By default, this parameter is True.

### Checkpointing:

These are set of APIs for saving and loading the checkpoint. These APIs take care of saving and loading the shard depending the tensor parallel rank of the worker.

### Save Checkpoint:

```
def neuronx_distributed.parallel_layers.save(state_dict, save_dir, save_serially=True,
↪ save_xser: bool=False, down_cast_bf16=False)
```

---

**Note:** This method will be deprecated, use `neuronx_distributed.trainer.save_checkpoint` instead.

---

This API will save the model from each tensor-parallel rank in the `save_dir`. Only workers with data parallel rank equal to 0 would be saving the checkpoints. Each tensor parallel rank would be creating a `tp_rank_ii_pp_rank_ii` folder inside `save_dir` and each ones saves its shard in the `tp_rank_ii_pp_rank_ii` folder. If `save_xser` is enabled, the folder name would be `tp_rank_ii_pp_rank_ii.tensors` and there will be a ref data file named as `tp_rank_ii_pp_rank_ii` in `save_dir` for each rank.

Parameters:

- `state_dict: (dict)`: Model state dict. Its the same dict that you would save using `torch.save`
- `save_dir: (str)`: Model save directory.
- `save_serially: (bool)`: This flag would save checkpoints one model-parallel rank at a time. This is particularly useful when we are checkpointing large models.
- `save_xser: (bool)`: This flag would save the model with torch xla serialization. This could significantly reduce checkpoint saving time when checkpointing large model, so it's recommended to enable `xser` when the model is large. Note that if a checkpoint is saved with `save_xser`, it needs to be loaded with `load_xser`, vice versa.
- `down_cast_bf16: (bool)`: This flag would downcast the `state_dict` to bf16 before saving.

## Load Checkpoint

```
def neuronx_distributed.parallel_layers.load(
    load_dir, model_or_optimizer=None, model_key='model', load_xser=False, sharded=True)
```

---

**Note:** This method will be deprecated, use `neuronx_distributed.trainer.load_checkpoint` instead.

---

This API will automatically load checkpoint depending on the tensor parallel rank. For large models, one should pass the model object to the load API to load the weights directly into the model. This could avoid host OOM, as the load API would load the checkpoints for one tensor parallel rank at a time.

Parameters:

- `load_dir`: (str) : Directory where the checkpoint is saved.
- `model_or_optimizer`: (torch.nn.Module or torch.optim.Optimizer): Model or Optimizer object.
- `model`: (torch.nn.Module or torch.optim.Optimizer): Model or Optimizer object, equivalent to `model_or_optimizer`
- `model_key`: (str) : The model key used when saving the model in the `state_dict`.
- `load_xser`: (bool) : Load model with torch xla serialization. Note that if a checkpoint is saved with `save_xser`, it needs to be loaded with `load_xser`, vice versa.
- `sharded`: (bool) : If the checkpoint is not sharded, pass False. This is useful (especially during inference) when the model is trained using a different strategy and you end up saving a single unsharded checkpoint. You can then load this unsharded checkpoint onto the sharded model. When this attribute is set to False, it is necessary to pass the model object. Note: The keys in the state-dict should have the same name as in the model object, else it would raise an error.

## Gradient Clipping:

With tensor parallelism, we need to handle the gradient clipping as we have to accumulate the total norm from all the tensor parallel ranks. This should be handled by the following API

```
def neuronx_distributed.parallel_layers.clip_grad_norm(
    parameters, max_norm, norm_type=2)
```

Parameters:

- `parameters` (Iterable[Tensor] or Tensor) : an iterable of Tensors or a single Tensor that will have gradients normalized
- `max_norm` (float or int) : max norm of the gradients
- `norm_type` (float or int) : type of the used p-norm. Can be 'inf' for infinity norm.

## Neuron Zero1 Optimizer:

In Neuronx-Distributed, we built a wrapper on the Zero1-Optimizer present in torch-xla.

```
class NeuronZero1Optimizer(Zero1Optimizer)
```

This wrapper takes into account the tensor-parallel degree and computes the grad-norm accordingly. It also provides two APIs: `save_sharded_state_dict` and `load_sharded_state_dict`. As the size of the model grows, saving the optimizer state from a single rank can result in OOMs. Hence, the api to `save_sharded_state_dict` can allow saving states from each data-parallel rank. To load this sharded optimizer state, there is a corresponding `load_sharded_state_dict` that allows each rank to pick its corresponding shard from the checkpoint directory.

```
optimizer_grouped_parameters = [
    {
        "params": [
            p for n, p in param_optimizer if not any(nd in n for nd in no_decay)
        ],
        "weight_decay": 0.01,
    },
    {
        "params": [
            p for n, p in param_optimizer if any(nd in n for nd in no_decay)
        ],
        "weight_decay": 0.0,
    },
]

optimizer = NeuronZero1Optimizer(
    optimizer_grouped_parameters,
    AdamW,
    lr=flags.lr,
    pin_layout=False,
    sharding_groups=parallel_state.get_data_parallel_group(as_list=True),
    grad_norm_groups=parallel_state.get_tensor_model_parallel_group(as_list=True),
)
```

The interface is same as `Zero1Optimizer` in torch-xla

```
save_sharded_state_dict(output_dir, save_serially = True)
```

---

**Note:** This method will be deprecated, use `neuronx_distributed.trainer.save_checkpoint` instead.

---

Parameters:

- `output_dir (str)` : Checkpoint directory where the sharded optimizer states need to be saved
- `save_serially (bool)`  
[Whether to save the states one data-parallel rank at a time. This is] especially useful when we want to checkpoint large models.

```
load_sharded_state_dict(output_dir, num_workers_per_step = 8)
```

---

**Note:** This method will be deprecated, use `neuronx_distributed.trainer.load_checkpoint` instead.

---

Parameters:

- `output_dir (str)`: Checkpoint directory where the sharded optimizer states are saved
- `num_workers_per_step (int)`: This argument controls how many workers are doing model load in parallel.

## Neuron PyTorch-Lightning

Neuron PyTorch-Lightning is currently based on Lightning version 2.1.0, and will eventually be upstreamed Lightning-AI code base

## Neuron Lightning Module

Inherited from `LightningModule`

```
class neuronx_distributed.lightning.NeuronLTModule(
    model_fn: Callable,
    nxd_config: Dict,
    opt_cls: Callable,
    scheduler_cls: Callable,
    model_args: Tuple = (),
    model_kwargs: Dict = {},
    opt_args: Tuple = (),
    opt_kwargs: Dict = {},
    scheduler_args: Tuple = (),
    scheduler_kwargs: Dict = {},
    grad_accum_steps: int = 1,
    log_rank0: bool = False,
    manual_opt: bool = True,
)
```

Parameters:

- `model_fn`: Model function to create the actual model
- `nxd_config`: Neuronx Distributed Config, output of `neuronx_distributed.neuronx_distributed_config`
- `opt_cls`: Callable to create optimizer
- `scheduler_cls`: Callable to create scheduler
- `model_args`: Tuple of args fed to model callable
- `model_kwargs`: Dict of keyworded args fed to model callable
- `opt_args`: Tuple of args fed to optimizer callable
- `opt_kwargs`: Dict of keyword args fed to optimizer callable
- `scheduler_args`: Tuple of args fed to scheduler callable
- `scheduler_kwargs`: Dict of keyworded args fed to scheduler callable
- `grad_accum_steps`: Grad accumulation steps



- `log_rank0`: Log at rank 0 (by default it will log at the last PP rank). Note that setting this to True will introduce extra communication per step hence causing performance drop
- `manual_opt`: Whether to do manual optimization, note that currently NeuronLTModule doesn't support auto optimization so this should always set to True

## Neuron XLA Strategy

Inherited from `XLAStrategy`

```
class neuronx_distributed.lightning.NeuronXLAStrategy(
    nxd_config: Dict = None,
    tensor_parallel_size: int = 1,
    pipeline_parallel_size: int = 1,
    save_load_xser: bool = True,
)
```

Parameters:

- `nxd_config`: Neuronx Distributed Config, output of `neuronx_distributed.neuronx_distributed_config`
- `tensor_parallel_size`: Tensor parallel degree, only needed when `nxd_config` is not specified
- `pipeline_parallel_size`: Pipeline parallel degree, only needed when `nxd_config` is not specified (Note that for now we only support TP with Neuron-PT-Lightning)
- `save_load_xser`: Set to True will enable save/load with xla serialization, for more context check [Save Checkpoint](#)

## Neuron XLA Precision Plugin

Inherited from `XLAPrecisionPlugin`

```
class neuronx_distributed.lightning.NeuronXLAPrecisionPlugin
```

## Neuron TQDM Progress Bar

Inherited from `TQDMProgressBar`

```
class neuronx_distributed.lightning.NeuronTQDMProgressBar
```

## Neuron TensorBoard Logger

Inherited from `TensorBoardLogger`

```
class neuronx_distributed.lightning.NeuronTensorBoardLogger(save_dir)
```

Parameters:

- `save_dir`: Directory to save the log files

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Inference APIs

### Table of contents

- *Model Trace:*
- *Trace Model Save/Load:*
  - *Save:*
  - *Load:*
  - *Parameters:*

### Model Trace:

We can use the tensor parallel layers to perform large model inference too. For performing inference, we can re-use the Parallel model built above for training and then use the trace APIs provided by the `neuronx_distributed` package to trace it for inference. One can use the following set of APIs for running distributed inference:

```
def neuronx_distributed.trace.parallel_model_trace(func, example_inputs, compiler_
↳workdir=None, compiler_args=None, inline_weights_to_neff=True, bucket_config=None, tp_
↳degree=1, max_parallel_compilations=None)
```

This API would launch tensor parallel workers, where each worker would trace its own model. These traced models would be wrapped with a single `TensorParallelModel` module which can then be used like any other traced model.

Parameters:

- **func** : Callable: This is a function that returns a `Model` object and a dictionary of states. The `parallel_model_trace` API would call this function inside each worker and run trace against them. Note: This differs from the `torch_neuronx.trace` where the `torch_neuronx.trace` requires a model object to be passed.
- **example\_inputs**: (`torch.Tensor` like) : The inputs that needs to be passed to the model. If you are using `bucket_config`, then this must be a list of inputs for each bucket model. This configuration is similar to `torch_neuronx.bucket_model_trace()`
- **compiler\_workdir**: `Optional[str, pathlib.Path]` : Work directory used by `[neuronx-cc]`. This can be useful for debugging and inspecting intermediary `[neuronx-cc]` outputs.
- **compiler\_args**: `Optional[Union[List[str], str]]` : List of strings representing `[neuronx-cc]` compiler arguments. See *Neuron Compiler CLI Reference Guide (neuronx-cc)* for more information about compiler options.
- **inline\_weights\_to\_neff**: `bool` : A boolean indicating whether the weights should be inlined to the NEFF. If set to `False`, weights will be separated from the NEFF. The default is `True`.
- **bucket\_config**: `torch_neuronx.BucketModelConfig` : The config object that defines bucket selection behavior. See `torch_neuronx.BucketModelConfig()` for more details.
- **tp\_degree**: (`int`) : How many devices to be used when performing tensor parallel sharding
- **max\_parallel\_compilations**: `Optional[int]` : If specified, this function will only trace these numbers of models in parallel, which can be necessary to prevent OOMs while tracing. The default is `None`, which means the number of parallel compilations is equal to the `tp_degree`.

**Trace Model Save/Load:****Save:**

```
def neuronx_distributed.trace.parallel_model_save(model, save_dir)
```

This API should save the traced model in `save_dir`. Each shard would be saved in its respective directory inside the `save_dir`. Parameters:

- `model`: (`TensorParallelModel`) : Traced model produced using the `parallel_model_trace` api. - `save_dir`: (`str`) : The directory where the model would be saved

**Load:**

```
def neuronx_distributed.trace.parallel_model_load(load_dir)
```

This API will load the sharded traced model into `TensorParallelModel` for inference.

**Parameters:**

- `load_dir`: (`str`) : Directory which contains the traced model.

*This document is relevant for: Inf2, Trn1, Trn2*

- *Distributed Strategies APIs*
- *Training APIs*
- *Inference APIs*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**3.3.4 Developer Guide**

*This document is relevant for: Inf2, Trn1, Trn2*

**Training Developer Guides**

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for Tensor Parallelism

### Training

For training models with tensor-parallelism, one would have to make few changes to their model/training script. Below we walk through the different changes one would have to make to shard the models across devices.

#### Creating DataLoader:

When we shard the model across devices using tensor parallelism, all the tensor parallel workers are operating on the same batch of data. Hence, to ensure that each tensor parallel worker is getting the same data, we make use of `DistributedSampler` as shown in the snippet below

```
def create_pretraining_dataset(
    input_file, max_pred_length, mini_batch_size, worker_init
):
    train_data = pretraining_dataset(
        input_file=input_file, max_pred_length=max_pred_length
    )
    # To distribute the data across different workers in the world,
    # we use the DistributedSampler. The num_replicas should be equal
    # to the data_parallel_world_size. Note: data_parallel_rank=0 can have
    # multiple tensor parallel ranks and each of these should get the same
    # data.
    train_sampler = DistributedSampler(
        train_data,
        num_replicas=parallel_state.get_data_parallel_world_size(),
        rank=parallel_state.get_data_parallel_rank(),
    )
    train_dataloader = DataLoader(
        train_data,
        sampler=train_sampler,
        batch_size=mini_batch_size,
        num_workers=0,
        worker_init_fn=worker_init,
        drop_last=True,
        pin_memory=True,
    )
    return train_dataloader
```

#### Creating Model:

One can create models by replacing the large linear layers with `ColumnParallel` and `RowParallel` Linear layers. In case of transformers, we have a good structure where the Attention block usually have linear projections for QKV and this is followed by a fully connected layer. Let's take a look at the example for the BERT model. We make the attention module of BERT model to use tensor parallel layers, thereby adding the ability to shard the model across devices.

```
class ParallelSelfAttention(transformers.models.bert.modeling_bert.BertSelfAttention):
    def __init__(self, config, position_embedding_type=None):
        super().__init__(config, position_embedding_type)
```

(continues on next page)

(continued from previous page)

```

self.query = ColumnParallelLinear(config.hidden_size,
                                   self.all_head_size,
                                   gather_output=False)
self.key = ColumnParallelLinear(config.hidden_size,
                                 self.all_head_size,
                                 gather_output=False)
self.value = ColumnParallelLinear(config.hidden_size,
                                   self.all_head_size,
                                   gather_output=False)
# Since we shard the number of attention heads across tensor parallel
# ranks, each rank would have a subset of heads, hence, we update
# the num_attention_heads here.
tp_size = parallel_state.get_tensor_parallel_size()
self.num_attention_heads = self.num_attention_heads // tp_size
self.all_head_size = self.all_head_size // tp_size

```

As seen we just had to swap out the linear layers with ColumnParallel Linear layers and the rest of the forward method of the attention layer can work as is. Note: In the above ColumnParallelLinear layer we are not gathering output from each rank, in other words, each ranks is working on its own shard. We can make gather\_output=True and that would gather output and you would get a full dim output. However, gathering output from all ranks would introduce an all-gather operation which can be expensive depending on the size of the tensor. In the case of attention module, we know that the SelfAttention block is followed by MLP block. Hence, we replace the linear layer there with a RowParallelLinear as shown below:

```

class ParallelSelfOutput(transformers.models.bert.modeling_bert.BertSelfOutput):
    def __init__(self, config):
        super().__init__(config)
        self.dense = RowParallelLinear(config.hidden_size,
                                       config.hidden_size,
                                       input_is_parallel=True)

```

As seen we just had to replace the dense layer here, and pass the input\_is\_parallel argument. This way, the RowParallelLinear should operator on partitions and get a collective result.

Making just the above two changes can help you partition good chunk of your model across multiple workers, thereby allowing models of larger size to be trained on a single instance. Note: Majority of the parameters of a transformer model are in these linear layers and hence partitioning these layers can help you scale.

### Final Training script:

Once the dataloader and model changes are done, we are ready to build the training script. Good news, you can use the same training loop as before for data-parallel training, and would need just the minor tweaks to get it all started.

```

from neuronx_distributed.parallel_layers import parallel_state, clip_grad_norm

neuronx_distributed.parallel_state.initialize_model_parallel(tensor_model_parallel_
↪size=2)
dataloader = create_pretraining_dataset(
    input_file, max_pred_length, mini_batch_size, worker_init)

model = YourNewlyBuiltParallelModel(config)

```

(continues on next page)

(continued from previous page)

```

# We have to move the model to device using this API, because when
# we move model to device using .to(device), the model parameter's
# attributes aren't preserved. This causes some of the tensor parallel
# attributes to be lost. Hence, this API takes care of preserving the
# tensor parallel attributes.
parallel_layers.move_model_to_device(model, device)

for inputs, labels in dataloader:
    output = model(*inputs)
    loss = loss_fn(output, labels)
    loss.backward()
    # Here we use clip_grad_norm from neuronx_distributed as that
    # can handle tensor parallel ranks
    clip_grad_norm(model.parameters(), max_norm)
    # For the optimizer step, we have to pass the data_parallel group
    xm.optimizer_step(
        optimizer,
        groups=parallel_state.get_data_parallel_group(as_list=True)
    )
    optimizer.zero_grad()
    scheduler.step()

```

Few things to take note of in the above code snippet: 1. We are initializing the model parallel with tensor parallel size of 2. This will shard the model across 2 devices. 2. We use the `move_model_to_device` API to move model to device. This is equivalent to doing `model.to(device)`. We need to explicitly call this API since some of the tensor-parallel attributes do not get copied over when we move the model to device using `model.to(device)`. 3. We are calling the `clip_grad_norm` from `parallel_layers`. This `clip_grad_norm` should take care of accumulating the `max_norm` from the tensor\_parallel ranks and producing the correct output. 4. We pass the `data_parallel_group` to the `optimizer_step`. If we don't pass the group, default would be all the workers in the world.

### Saving Model:

Once training is done, we want to save the model. This can be done easily by calling the `save` api from `neuronx_distributed.parallel_layers`. Here is an example:

```

neuronx_distributed.parallel_layers.save({
    'epoch': epoch,
    'model': model.state_dict(),
    'optimizer_state_dict': optimizer.state_dict(),
    'loss': loss,
    ...
}, PATH)

```

Note the `model` key used here, we need to provide the same key during model load.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for Pipeline Parallelism

### Training

For training models with pipeline-parallelism, user needs to make few changes to their model/training script. In the below steps, we walk through different changes user has to make to use pipeline parallelism. For general changes please refer to [tensor parallel guidance](#).

### Creating Model

To train with pipeline parallel, user needs to wrap their torch module with NeuronxDistributed's Pipeline Parallel model wrapper, i.e. `NxDPPModel`. Let's take a look at our Llama example:

```
# Create torch model
config.return_dict = False
model = transformers.LlamaForCausalLM(config)
# Create pipeline cuts
pipeline_cuts = create_partition(config, args)
# Apply model wrapper
model = NxDPPModel(
    model,
    transformer_layer_cls=LlamaDecoderLayer,
    num_microbatches=args.num_microbatches,
    virtual_pipeline_size=1,
    output_loss_value_spec=(True, False),
    input_names=["input_ids", "attention_mask", "labels"],
    pipeline_cuts=pipeline_cuts,
    trace_file_path=args.trace_file_path,
    leaf_module_cls=[LlamaRMSNorm.__name__],
    autowrap_modules=[mappings],
    use_zero1_optimizer=args.use_zero1_optimizer,
    deallocate_pipeline_outputs=False,
)
model.move_model_to_device()
```

We first create the model from the Hugging Face model config. If tensor parallel needs to be applied to model it must be done here before applying the pipeline parallel model wrapper. The next step is to create the partitions. Here is an example to evenly partition the layers for all stages:

```
def create_partition(config, args):
    """
    Evenly split the transformer layers between the PP ranks
    """
    assert config.num_hidden_layers % args.pipeline_parallel_size == 0
    num_layer_per_partition = config.num_hidden_layers // args.pipeline_parallel_size
    pipeline_cuts = []
    current_cut = num_layer_per_partition - 1
    for i in range(args.pipeline_parallel_size-1):
        pipeline_cuts.append(f"model.layers.{current_cut}")
        current_cut += num_layer_per_partition
    if torch.distributed.get_rank() == 0:
```

(continues on next page)

(continued from previous page)

```
print(f"pipeline_cuts {pipeline_cuts}")
return pipeline_cuts
```

Note that the pipeline cuts should be at the transformer layer module name, which in Llama model is indicated as `model.layers.i` where `i` is the layer index. Users have the option to either provide the pipeline cuts, or set `auto_partition` to `True` to automatically determine the pipeline cuts to use. After pipeline cuts are decided, pipeline model wrapper is applied. Let's take a deeper look into each input of the model wrapper

- `model`: The original Pytorch module, could be TPfied.
- `transformer_layer_cls=LlamaDecoderLayer`: The transformer layer class, we will use it for partition
- `num_microbatches=args.num_microbatches`: The number of microbatches we used for pipeline execution.
- `virtual_pipeline_size`: Virtual pipeline size if greater than 1 we will use the interleaved pipeline schedule.
- `output_loss_value_spec=(True, False)`: This tells `NxDPPModel` how to get the loss from the model output. In this case output is a tuple, where first value is loss and second value is something else. `NxDPPModel` will use loss to run backward and return loss as the output.
- `input_names=["input_ids", "attention_mask", "labels"]`: The model input names that we will use to run training. As our partition uses FX symbolic trace to trace the model, we will use these input names to create `concrete_args`. Usually this will be the same input as you will feed into model for the execution. For details please check [https://pytorch.org/docs/stable/fx.html#torch.fx.symbolic\\_trace](https://pytorch.org/docs/stable/fx.html#torch.fx.symbolic_trace)
- `pipeline_cuts=pipeline_cuts`: The pipeline cuts to decide the stages
- `leaf_module_cls=[LlamaRMSNorm.__name__]`: We can add some pytorch modules as leaf module so that FX symbolic trace won't trace it through. Here we mark the `LlamaRMSNorm` as one leaf module. If you hit any issue about tracing you can skip tracing that part by add the module as a leaf module here. The transformer layer module will be a leaf module by default.
- `autowrap_modules`: This serves as the same functionality to simplify FX tracing. User can provide a **python** module here and all the methods from this python module will not be traced.
- `use_zero1_optimizer`: When zero-1 optimizer is used, set this to `True`, so the PP model will understand that zero-1 optimizer will handle data parallel gradient averaging.
- **deallocate\_pipeline\_outputs:**  
Whether to deallocate the pipeline outputs after send. After send the output tensor is only useful for its `'grad_fn'` field, and not its `'data'`.

After applying model wrapper, `NxDPPModel` will partition the model based on the pipeline cuts. If the original model is not yet moved to device, we can call `model.move_model_to_device()` so that `NxDPPModel` will only move the local module to device.

### Runtime execution:

To use pipeline runtime, user simply needs to replace their original model call with `NxDPPModel.run_train`, rest will remain unchanged. Please note that the pipeline runtime will take care of both forward and backward call, so user will not need to explicitly make backward calls. The `NxDPPModel.run_train` call will return the loss that is achieved from `output_loss_value_spec`.



### Interleaved Pipeline-Parallelism:

To use interleaved pipeline parallel, one has to set `virtual_pipeline_size` greater than 1. The value of the `virtual_pipeline_size * pipeline_parallel_size` should be equal to the number of layers in the models. Interleave pipeline can help to reduce the pipeline bubble size and improve performance especially in cases when the number of microbatches per data-parallel rank is small. More information can be found [here](#)

### Mixed precision training

We support the torch autocast to do mixed precision, simply apply the context manager for the `NxDPPModel.run_train` call. Here is an example:

```
# replace loss, _ = model(input_ids, attention_mask, labels) with below
with torch.autocast(enabled=args.use_amp > 0, dtype=torch.bfloat16, device_type="cuda"):
    loss = model.run_train(
        input_ids=input_ids,
        attention_mask=attention_mask,
        labels=labels,
    )
```

### Things that require user attention:

#### Model initialization

When the model is large, it is easy to cause host OOM when full model is created on every Neuron core. We recommend 2 ways to deal with this situation:

#### Using torchdistx's deferred initialization

Pytorch's torchdistx package (<https://github.com/pytorch/torchdistx/tree/main>) provides easy way to do deferred initialization. If you have torchdistx installed, using deferred initialization is simple as below

```
from torchdistx import deferred_init
# Instead of model = LlamaForCausalLM(config)
model = deferred_init.deferred_init(LlamaForCausalLM, config)
```

The model weights will be initialized in fake tensor mode which will not consume memory. After applying the `NxDPPModel` model wrapper we will only materialize the weights that belong to the local module. Please be aware that the torchdistx package is not actively maintained by Meta, please use at your own risk.

## Using meta device for initialization

NeuronxDistributed also supports also offer a way to first create the model on meta device, then reinitialize it to host device with only the local modules. To create the model on meta device, follow the below example:

```
from neuronx_distributed.utils.model_utils import init_on_device
with init_on_device(torch.device("meta")):
    model = LlamaForCausalLM(config)
```

With `init_on_device(torch.device("meta"))` context manager, all model weights will be create to meta device, which will not consume host memory. Then during applying the PP model wrapper, user can pass the `param_init_fn` kwargs which can define how to reinit the parameter. Here is an example:

```
def init_weights(module):
    from neuronx_distributed.parallel_layers import ColumnParallelLinear,
    ↪RowParallelLinear, ParallelEmbedding
    if isinstance(module, (nn.Linear, Conv1D)):
        module.weight.data.normal_(mean=0.0, std=model_config.initializer_range)
        if module.bias is not None:
            module.bias.data.zero_()
    elif isinstance(module, nn.Embedding):
        module.weight.data.normal_(mean=0.0, std=model_config.initializer_range)
        if module.padding_idx:
            module.weight.data[module.padding_idx].zero_()
    elif isinstance(module, nn.LayerNorm):
        module.bias.data.zero_()
        module.weight.data.fill_(1.0)
    elif isinstance(module, (ParallelEmbedding, RowParallelLinear,
    ↪ColumnParallelLinear)):
        module.init_weight_cpu()
        if hasattr(module, "bias") and module.bias is not None:
            module.bias.data.zero_()

model = NxDPModel(..., param_init_fn=init_weights,...)
```

`param_init_fn` should take a module as input and initialize how the weight of that module should be initialized.

## Moving model to device

When user create the model it is usually either created on CPU, or using meta device/torchdistx for delayed parameter initialization. It is important to understand when the delayed parameter will be materialized and how/when to move model to device.

Once the `NxDPModel` wrapper is applied with the model together with the partition information, tracing and partition will happen immediately. After partition we will materialize the local module if torchdistx is used or `param_init_fn` is passed. So the returned model of `NxDPModel` wrapper will have local parameters on host device.

After model is wrapped with `NxDPModel` user can do things that are recommended to run on CPU, e.g. loading shared checkpoint. It is important to make sure to call `model.move_model_to_device()` before creating the optimizer, so that the optimizer can take the weights that are on the device. When using zero-1 optimizer, it is also required to use `model.local_parameters()` to create parameter groups so the optimizer can infer the right device information from parameter groups.

## Gradient checkpointing

Gradient checkpointing (or activation checkpointing) is a common method used in deep learning to reduce memory footprint by doing recomputation of forward computation. The common way to apply the gradient checkpointing on XLA device is to use the torch\_xla's [gradient checkpointing wrapper](#), which will apply an autograd function. However FX's symbolic tracing does not understand autograd function, and as a result the checkpointing information will be ignored if the checkpoint wrapper is traced during partition. To handle this case, user can manually re-apply gradient checkpoint after partition. Here we provide an example to checkpoint every transformer layer after partition.

```
from typing import Any, Dict, Iterator, Tuple
import torch.nn as nn

import torch
from torch_xla.utils.checkpoint import checkpoint as torch_checkpoint
from neuronx_distributed.parallel_layers.parallel_state import rmsg
from neuronx_distributed.utils.logger import get_logger
from torch.distributed.utils import _replace_by_prefix

logger = get_logger()

_CHECKPOINT_WRAPPED_MODULE = "mod"
_CHECKPOINT_PREFIX = _CHECKPOINT_WRAPPED_MODULE + "."

class CheckPointWrapper(torch.nn.Module):
    def __init__(self, mod) -> None:
        super().__init__()
        self.mod = mod
        # state_dict post hook to remove prefix to allow loading into a
        # non-checkpoint wrapped module.
        self._register_state_dict_hook(self._post_state_dict_hook)
        # load_state_dict pre-hook to allow loading back into
        # checkpoint-wrapped module.
        self._register_load_state_dict_pre_hook(
            self._pre_load_state_dict_hook, with_module=True
        )

    def forward(self, *args, **kwargs):
        ordered_args = list(args)
        for value in kwargs.values():
            ordered_args += [value]

        # Note: checkpoint cannot accept kwargs
        return torch_checkpoint(self.mod, *ordered_args, use_reentrant=True)

    def named_parameters(
        self,
        *args,
        **kwargs,
    ) -> Iterator[Tuple[str, torch.nn.Parameter]]:
        """
        Overrides :meth:`named_parameters()` to intercept parameter names and
        remove all occurrences of ``_CHECKPOINT_PREFIX``.

```

(continues on next page)

(continued from previous page)

```

"""
    for param_name, param in super().named_parameters(*args, **kwargs):
        updated_name = param_name.replace(_CHECKPOINT_PREFIX, "")
        yield updated_name, param

def named_modules(self, *args, **kwargs):
    for module_name, module in super().named_modules(*args, **kwargs):
        updated_name = module_name.replace(_CHECKPOINT_PREFIX, "")
        yield updated_name, module

@staticmethod
def _post_state_dict_hook(
    module: nn.Module,
    state_dict: Dict[str, Any],
    prefix: str,
    *args: Any,
) -> Dict[str, Any]:
    """
    _post_state_dict_hook() is called after the state_dict() of this
    FSDP module is executed. For ``checkpoint_wrapper``, it will strip
    checkpoint-wrapped module prefix so that this module can be loaded into
    non-checkpointed modules. It would still be able to be loaded into
    checkpoint-wrapped modules as this class adds the prefix back before
    loading the state_dict.
    """
    _replace_by_prefix(state_dict, f"{prefix}{_CHECKPOINT_PREFIX}", prefix)
    return state_dict

@staticmethod
def _pre_load_state_dict_hook(
    module: nn.Module,
    state_dict: Dict[str, Any],
    prefix: str,
    *args: Any,
) -> None:
    """
    ``_pre_state_dict_hook`` is called before ``self._load_from_state_dict()``
    is called. For ``checkpoint_wrapper``, it will add back the module
    prefix so that non-checkpointed modules can be loaded into
    checkpoint_wrapper modules properly.
    """
    _replace_by_prefix(state_dict, prefix, prefix + f"{_CHECKPOINT_PREFIX}")

def apply_checkpoint(dist_model, layers_to_checkpoint=None):
    checkpoint_wrapper_added = False
    if layers_to_checkpoint is not None and len(layers_to_checkpoint) == 0:
        raise RuntimeError(
            rmsg(f"invalid input layers_to_checkpoint {layers_to_checkpoint}, can't be_
↪ empty")
        )
    for name, module in dist_model.local_module.named_children():
        # checkpoint layers that are provided in input

```

(continues on next page)

(continued from previous page)

```

    # if layers not provide in input, then checkpoint if it is transformer layer
    if (layers_to_checkpoint and name in layers_to_checkpoint) or (
        not layers_to_checkpoint and type(module) == dist_model.transformer_layer_cls
    ):
        # add_module replaces old module with our own custom module.
        # https://pytorch.org/docs/stable/_modules/torch/nn/modules/module.html
        ↪ #Module.add_module
        dist_model.local_module.add_module(name, CheckPointWrapper(module))
        checkpoint_wrapper_added = True
    if layers_to_checkpoint is not None and not checkpoint_wrapper_added:
        logger.warning(
            rmsg(f"layers_to_checkpoint {layers_to_checkpoint} do not exist in the graph
        ↪ ")
        )
    elif layers_to_checkpoint is None and not checkpoint_wrapper_added:
        logger.warning(
            rmsg(
                f"During applying activation checkpointing, transformer_layer_cls {dist_
        ↪ model.transformer_layer_cls.__name__} can not be found in stage {dist_model.pipeline_
        ↪ parallel_rank}, skipping..."
            )
        )

model = NxDPModel(...)
# Will checkpoint every transformer layer
apply_checkpoint(model)

```

apply\_checkpoint function will try to apply gradient checkpointing to every transformer layer. Please note we have plan to add this functionality into NxDPModel in the future releases.

## Model tracing

It is important to understand that the model cannot be partitioned without tracing. The model tracing is currently done with FX's symbolic trace. There are [certain limitations for FX's symbolic trace](#). So in order to avoid any tracing issue, we would like to trace as less operations as possible, which means that we only want to trace the structure of the model, and cut the pipeline stages on the transformer layers, we do not care how exactly the computations are in the model. By default, we will mark all transformer layers as leaf nodes, so that the tracer will not trace inside these layers. If you have some module that might cause tracing problem, you can try to mark them as leaf nodes as well. Our previous example also marks the *LlamaRMSNorm* as leaf module for Llama model.

## Special treatment for Hugging Face models

Hugging Face offers FX support for many of its models. We will detect if user is using a Hugging Face model (by checking if the model class is `transformers.PreTrainedModel`), and if so we will use the Huggingface's FX tracer to do the symbolic trace. The Hugging Face's tracer has implementation of many functionalities to help tracing, for details please refer to [here](#). However, please be aware that Hugging Face's tracer will check if the model class name belongs to one of the Hugging Face models. So if you create your model class based on some Huggingface model class, it is important to maintain the same class name. Below is an example:

```
from transformers.models.llama.modeling_llama import LlamaForCausalLM as L
↳ LlamaForCausalLMHF

# Keep the same class name as original one
class LlamaForCausalLM(LlamaForCausalLMHF):
    ...
```

## Auto partition

Setting the `auto_partition` parameter to `True` means that the transformer layers are automatically partitioned by evenly splitting the transformer layers between the PP ranks. If the transformer layers are not evenly divisible by the PP ranks, the remaining layers are distributed to the latter pipeline ranks. The partitions are created on the basis of the transformer layer names. The transformer layer names are determined by recursively traversing the original torch module to find the layer names of modules that are of the `transformer_layer_cls` type in the model. If the user does not want to partition the model in this way, they can set the partitions to use by specifying the `pipeline_cuts`. Note that the pipeline cuts should be at the transformer layer module name, which in the Llama model is given by `model.layers.i` where `i` is the layer index.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for Activation Memory reduction

### Sequence Parallelism

To combine sequence parallelism with tensor-parallelism, one needs to follow the steps below:

#### Model changes for Tensor-Parallel block:

For tensor-parallelism, we replace the linear layers with `ColumnParallel` and `RowParallel` Linear layers as mentioned [here](#). To enable sequence-parallel, we need to pass the `sequence_parallel_enabled` for the `ColumnParallel` and `RowParallel` linear layers. Setting this argument to `true`, the `ColumnParallel` and `RowParallel` Linear layers will introduce the `all-gather` and `reduce-scatter` operations for gathering and distributing the activations along the sequence dimension.

```
from transformers.models.gpt_neox.modeling_gpt_neox import GPTNeoXAttention

class class GPTNeoXAttentionNx1D(GPTNeoXAttention):
    def __init__(self, config):
        super().__init__(config)
        ....
        self.query_key_value = ColumnParallelLinear(
            config.hidden_size,
            3 * config.hidden_size,
            stride=3,
            gather_output=False,
            init_method=init_method,
            sequence_parallel_enabled=self.config.sequence_parallel_
↳ enabled,
```

(continues on next page)

(continued from previous page)

```

        )
    self.dense = RowParallelLinear(
        config.hidden_size,
        config.hidden_size,
        input_is_parallel=True,
        init_method=init_method,
        sequence_parallel_enabled=self.config.sequence_parallel_enabled,
    )
    ....

```

### Model changes for Non-Tensor-Parallel block:

In a transformer module, the non-tensor parallel block contains mainly the Layer-Norm modules. Since we partition the computation along the sequence dimension for the layer-norm, we need to sum up the gradients along the sequence dimension for the Layer-norm. To help us do that, we use the Layer-norm provided from `neuronx-distributed.parallel_layers.layer_norm`. The Layer-norm in `neuronx-distributed` should use the same forward and backward as `torch.nn.LayerNorm`, however, it just marks the weights as sequence-parallel weights. This tagging allows us to look for weights with sequence-parallel tagging and reduce those gradients along the tensor-parallel degree. Hence we need to add the following two changes:

```

from transformers.models.gpt_neox.modeling_gpt_neox import GPTNeoXLayer
from neuronx_distributed.parallel_layers import layer_norm

class GPTNeoXLayerNxD(GPTNeoXLayer):
    def __init__(self, config):
        super().__init__(config)
        ...
        self.input_layernorm = layer_norm.LayerNorm(
            config.hidden_size,
            eps=config.layer_norm_eps,
            sequence_parallel_enabled=config.sequence_parallel_
↪enabled
        )
        self.post_attention_layernorm = layer_norm.LayerNorm(
            config.hidden_size,
            eps=config.layer_norm_eps,
            sequence_parallel_enabled=config.sequence_
↪parallel_enabled
        )

```

Once we replace the layernorm with `neuronx-distributed`'s layernorm, it will [mark the weights](#) as sequence-parallel weights. Note: If your model is using RMSNorm or any other layer that parallelizes in the sequence-dimension, you can mark the weights as sequence-parallel weights by using the following code:

```
setattr(param, "sequence_parallel_enabled", sequence_parallel_enabled)
```

Once marked, we then use this attribute when we compute gradients for layer-norm. We need to add the following code before our `optimizer.step`:

```

def allreduce_sequence_parallel_gradients(optimizer):
    """ All-reduce layernorm parameters across model parallel nodes when sequence_

```

(continues on next page)

(continued from previous page)

```

↪ parallelism is used.
    Modified from megatron-lm:
    https://gitlab-master.nvidia.com/ADLR/megatron-lm/-/blob/
↪ 3f91f09bb2ab32f9904b47f46f19d2fc3f518ed8/megatron/training.py#L425
    """
    from neuronx_distributed.parallel_layers.mappings import reduce_from_tensor_model_
↪ parallel_region
    grads = []
    for param_group in optimizer.__getstate__()['param_groups']:
        for group, params in param_group.items():
            if group == 'params':
                for p in params:
                    if isinstance(p, torch.Tensor) and p.grad is not None:
                        sequence_parallel_param = getattr(p, 'sequence_parallel_enabled',
↪ False)
                        if sequence_parallel_param:
                            grads.append(p.grad.data)
    for grad in grads:
        reduce_from_tensor_model_parallel_region(grad)

```

As seen in the above code, we reduce the gradients from all tensor parallel devices. This is because the compute is divided along the sequence dimension across all the devices participating in the tensor parallel group. For reference implementation, check the [GPTNeoX-20B modeling code](#).

### Transposing the activations:

Sequence-parallelism implementation requires the sequence dimension to be the 0th dimension whereas the tensor-parallel region requires the sequence dimension to be the first dimension. All our model implementation keeps the sequence dimension as 1st dimension and batch dimension as 0th dimension. Hence, to accommodate sequence parallelism, we need to insert a few transpose operations at the following places:

1. Before we start looping through all the layers, we need to transpose the sequence and batch dimension. We also need to partition the inputs along the sequence dimensions such that each tp-rank gets a part. This can be done as:

```

from neuronx_distributed.parallel_layers.mappings import scatter_to_sequence_parallel_
↪ region
# Nx D Core code change: sequence parallel uses seq_len as the 0-th dim
if self.config.sequence_parallel_enabled:
    hidden_states = hidden_states.transpose(0, 1).contiguous()
    hidden_states = scatter_to_sequence_parallel_region(hidden_states)

```

2. Since the attention block requires the sequence dimension to be 1st dimension, we transpose the output of QKV projection and then transpose it back before the final MLP of the attention block.

```

# Within the attention module
qkv = self.query_key_value(hidden_states)

if config.sequence_parallel_enabled:
    qkv = qkv.transpose(0, 1)
...

```

(continues on next page)



(continued from previous page)

```
attn_output = attn_output.transpose(0,1)
attn_output = self.dense(attn_output)
```

3. Finally before returning the final output, we need to put all the partial activations along the sequence dimension back together. This can be done as follows:

```
from neuronx_distributed.parallel_layers.mappings import gather_from_sequence_parallel_
    ↪ region
if self.config.sequence_parallel_enabled:
    hidden_states = gather_from_sequence_parallel_region(hidden_states, to_model_
    ↪ parallel=False)
    hidden_states = hidden_states.transpose(0, 1).contiguous()

return BaseModelOutputWithPast(
    last_hidden_state=hidden_states,
    past_key_values=presents,
    hidden_states=all_hidden_states,
    attentions=all_attentions,
)
```

These are the only major changes required to add sequence-parallelism on top of tensor-parallelism. Note: Sequence-parallelism uses the same tensor-parallel group. For reference implementation, follow [GPTNeoX-20B model script](#).

## Activation Recomputation

As seen in the [App notes on Activation Memory Recomputation](#) we can reduce the activation memory by recomputing few operations from the forward pass during the backward run. To replay some of the compute, we can use the `torch.utils.checkpoint.checkpoint`. To use this API, we need to put the compute, we want to replay, inside a function which can be passed to the `checkpoint` API. This API takes care of maintaining the RNG seed, not saving the activations and also inserting the forward recompute during the gradient computation.

To enable selective activation checkpointing for the attention block, we can simply pass the attention block to the `checkpoint` api as follows:

```
if config.selective_activation_checkpointing_is_enabled:
    attn_output = torch.utils.checkpoint.checkpoint(self._attn, query, key, value,
    ↪ attention_mask, head_mask)
else:
    attn_output = self._attn(query, key, value, attention_mask, head_mask)
```

Note: To use `torch.utils.checkpoint`, it is mandatory to use `-O1` compiler flag. If this is not enabled, the Neuron compiler would eliminate the duplicate recompute as an optimization and hence you would not see any memory gains.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for save/load checkpoint

This document will introduce how to use `nxd.save_checkpoint` and `nxd.load_checkpoint` to save and load checkpoint for distributed model training. This two methods handle all checkpoint in a single method: model, optimize, learning rate scheduler and any user contents.

Model states are saved on data parallel rank-0 only. When ZeRO-1 optimizer is not turned on, optimizer states are also saved like this; while when ZeRO-1 optimizer is turned on, states are saved on all ranks. Scheduler and user contents are saved on master rank only.

For a complete api guide, refer to [API GUIDE](#).

### Save checkpoint:

A sample usage:

```
nxd.save_checkpoint(
    args.checkpoint_dir, # checkpoint path
    tag=f"step_{total_steps}", # tag, sub-directory under checkpoint path
    model=model,
    optimizer=optimizer,
    scheduler=lr_scheduler,
    user_content={"total_steps": total_steps, "batch_idx": batch_idx, "cli_args": args.__
    dict__},
    use_xser=True,
    async_save=True,
)
```

Users can choose to not save every thing. For example, model states only:

```
nxd.save_checkpoint(
    args.checkpoint_dir, # checkpoint path
    tag=f"step_{total_steps}", # tag, sub-directory under checkpoint path
    model=model,
    use_xser=True,
    async_save=True,
)
```

To only keep several checkpoints (e.g. 5), just use `num_kept_ckpts=5`.

### Load checkpoint:

A sample usage, note that if no user contents detected, it will return `None`:

```
user_content = nxd.load_checkpoint(
    args.checkpoint_dir, # checkpoint path
    tag=f"step_{args.loading_step}", # tag
    model=model,
    optimizer=optimizer,
    scheduler=lr_scheduler,
)
```

Leave `tag` not provided, this loading method will try to automatically resume from the latest checkpoint.

```

user_content = nxd.load_checkpoint(
    args.checkpoint_dir, # checkpoint path
    model=model,
    optimizer=optimizer,
    scheduler=lr_scheduler,
)

```

### ZeRO-1 Optimizer State Offline Conversion:

ZeRO-1 optimizer checkpoint are sharded states stored for each rank. When user want to load ZeRO-1 optimizer states with different cluster setting (e.g. with DP degree changed), they can run the offline ZeRO-1 optimizer checkpoint conversion tool. This tool supports conversion from sharded states to full states, from full to sharded, and from sharded to sharded.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for Neuron-PT-Lightning

### Training

For training models with Neuron-PT-Lightning, user needs to make few changes to their model/training script. In this document we explain how we can train a model using Tensor Parallelism (TP), Data Parallelism (DP) and Zero-1.

First, let's start with the model changes. Please follow the guidelines here ([tensor parallel guidance](#)) for building the model with tensor-parallelism enabled and setting up training dataset.

Next, let's walkthrough how we can build the training loop with Neuron-PT-Lightning APIs

### Configure NeuronLTModule

NeuronxDistributed overrides [LightningModule](#) with built-in support for Neuron device. User needs to inherit from NeuronLTModule

```

class NeuronLlamaLTModule(NeuronLTModule):
    def training_step(self, batch, batch_idx):
        ...

```

Within LTModule, user needs to override the following methods `training_step` At this moment NeuronLTModule only support [manual optimization](#), so user needs to define forward, backward and optimization steps

```

def training_step(self, batch, batch_idx):
    xm.mark_step() # Isolate forward+backward graph
    for logger in self.trainer.loggers:
        logger.print_step = -1
    self.should_print = False
    outputs = self.model(
        input_ids=batch["input_ids"],
        attention_mask=batch["attention_mask"],

```

(continues on next page)

(continued from previous page)

```

        labels=batch["labels"],
    )
    loss = outputs.loss / self.grad_accum_steps
    loss.backward()
    self.averaged_loss += loss.detach()
    xm.mark_step() # Isolate forward+backward graph
    if not self.automatic_optimization and (batch_idx +1) % self.grad_accum_steps == 0:
        self.should_print = True
        loss_div = self.averaged_loss / self.trainer.strategy.data_parallel_size
        loss_reduced = xm.all_reduce(
            xm.REDUCE_SUM,
            loss_div,
            groups=parallel_state.get_data_parallel_group(as_list=True),
        )
        loss_reduced_detached = loss_reduced.detach()
        self.averaged_loss.zero_()
        optimizer = self.optimizers()
        scheduler = self.lr_schedulers()
        optimizer.step()
        optimizer.zero_grad()
        scheduler.step()
        xm.mark_step() # Isolate Optimization step graph

        # Setup items for logging
        self.loss = loss_reduced_detached
    return loss

```

`configure_optimizers` Configure optimizer and lr\_scheduler

```

def configure_optimizers(self):
    param_groups = self.get_param_groups_by_weight_decay()
    optimizer = initialize_parallel_optimizer(
        self.nxd_config, self.opt_cls, param_groups, **self.opt_kwargs
    )
    optimizer.zero_grad()
    scheduler = self.scheduler_cls(optimizer, *self.scheduler_args, **self.scheduler_
↪kwargs)
    return (
        [optimizer],
        [
            {
                "scheduler": scheduler,
            }
        ],
    )

```

`on_train_batch_end` Customized behaviour at the end of each training batch, like logging

```

def on_train_batch_end(self, *args, **kwargs):
    if self.should_print:
        if not self.automatic_optimization:
            self.log(

```

(continues on next page)

(continued from previous page)

```

        "loss",
        self.loss.detach().cpu().item() if self.loss is not None else torch.
→ zeros(1, device="cpu", requires_grad=False),
        prog_bar=True,
    )
    self.log(
        "global_step",
        self.global_step,
        prog_bar=True,
        on_step=True,
        on_epoch=True,
    )
    for logger in self.trainer.loggers:
        logger.print_step = self.global_step

```

Note that NeuronLTModule has a built-in function of `get_param_groups_by_weight_decay` for common use case as shown in snippet below, users can also override with their own `param_groups` generation.

```

def get_param_groups_by_weight_decay(self):
    """Get param groups. Customers can override this to have their own way of weight_
→ decay"""
    param_optimizer = list(self.model.named_parameters())
    no_decay = ["bias", "LayerNorm"] # gamma/beta are in LayerNorm.weight

    optimizer_grouped_parameters = [
        {
            "params": [p for n, p in param_optimizer if not any(nd in n for nd in no_
→ decay)],
            "weight_decay": 0.01,
        },
        {
            "params": [p for n, p in param_optimizer if any(nd in n for nd in no_decay)],
            "weight_decay": 0.0,
        },
    ]
    return optimizer_grouped_parameters

```

## Configure DataModule

Create a LightningDataModule for data loading/sampling

```

class NeuronLightningDataModule(LightningDataModule):
    def __init__(
        self,
        dataloader_fn: Callable,
        data_dir: str,
        batch_size: int,
        data_args: Tuple = (),
        data_kwargs: Dict = {},
    ):
        super().__init__()

```

(continues on next page)

(continued from previous page)

```

        self.dataloader_fn = dataloader_fn
        self.data_dir = data_dir
        self.batch_size = batch_size
        self.data_args = data_args,
        self.data_kwargs = data_kwargs

    def setup(self, stage: str):
        pass

    def train_dataloader(self):
        return self.dataloader_fn(
            self.data_dir,
            self.batch_size,
            self.trainer.strategy.data_parallel_size,
            self.trainer.strategy.data_parallel_rank,
            *self.data_args,
            **self.data_kwargs
        )

```

## Update Training Script

For detailed introduction to each api/class, check [api guide](#)

## Create NeuronLTModule and DataModule

```

model = NeuronLlamaLTModule(
    model_fn = LlamaForCausalLM,
    nxd_config = nxd_config,
    model_args = (model_config,),
    opt_cls = optimizer_cls,
    scheduler_cls = configure_scheduler,
    opt_kwargs = {
        "lr": flags.lr,
    },
    scheduler_args = (flags.warmup_steps, flags.max_steps),
    grad_accum_steps = flags.grad_accum_usteps,
    manual_opt = True,
)

dm = NeuronLightningDataModule(
    create_llama_pretraining_dataset,
    flags.data_dir,
    flags.batch_size,
    data_args = (flags.seed,),
)

```

### Add Strategy, Plugins, Callbacks

```
strategy = NeuronXLAStrategy(
    nxd_config = nxd_config
)
plugins = []
plugins.append(NeuronXLAPrecisionPlugin())
callbacks = []
callbacks.append(NeuronTQDMProgressBar())
```

### Create Trainer and Start Training

```
trainer = Trainer(
    strategy = strategy,
    max_steps = flags.steps_this_run,
    plugins = plugins,
    enable_checkpointing = flags.save_checkpoint,
    logger = NeuronTensorBoardLogger(save_dir=flags.log_dir),
    log_every_n_steps = 1,
    callbacks = callbacks,
)
trainer.fit(model=model, datamodule=dm)
```

### Checkpointing

To enable checkpoint saving, add [ModelCheckpoint](#) to the callbacks

```
callbacks.append(
    ModelCheckpoint(
        save_top_k = flags.num_kept_checkpoint,
        monitor="global_step",
        mode="max",
        every_n_train_steps = flags.checkpoint_freq,
        dirpath = flags.checkpoint_dir,
    )
)
```

To load from specific checkpoint, add `ckpt_path=ckpt_path` to `trainer.fit`

```
trainer.fit(model=model, datamodule=dm, ckpt_path=ckpt_path)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for model and optimizer wrapper

Model and optimizer wrapper are useful tools to wrap original model and optimizer while keep the API unchanged. We recommend to always use model and optimizer wrappers, it's helpful to apply optimizations and hide the complexity from the optimizations. Users need to care about the implementation details of the optimization, just use the wrappers as you normally use `torch.nn.Module` and `torch.optim.Optimizer`.

For a complete api guide, refer to [API GUIDE](#).

### Create training config:

To use model and optimizer wrapper, we need to create `neuronx_distributed` config firstly.

A sample config use tensor parallel, pipeline parallel, ZeRO-1 optimizer, sequence parallel and activation checkpointing:

```
nxd_config = nxd.neuronx_distributed_config(
    tensor_parallel_size=args.tensor_parallel_size,
    pipeline_parallel_size=args.pipeline_parallel_size,
    pipeline_config={
        "transformer_layer_cls": LlamaDecoderLayer,
        "num_microbatches": args.num_microbatches,
        "output_loss_value_spec": (True, False),
        "input_names": ["input_ids", "attention_mask", "labels"],
        "pipeline_cuts": pipeline_cuts,
        "trace_file_path": args.trace_file_path,
        "param_init_fn": None,
        "leaf_module_cls": [LlamaRMSNorm.__name__],
        "autowrap_modules": [mappings],
        "use_zero1_optimizer": args.use_zero1_optimizer > 0,
        "use_optimizer_wrapper": True,
    },
    optimizer_config={
        "zero_one_enabled": args.use_zero1_optimizer > 0,
        "grad_clipping": True,
        "max_grad_norm": 1.0,
    },
    sequence_parallel=args.use_sequence_parallel,
    activation_checkpoint_config=CoreAttention if args.use_selective_checkpoint > 0 else
    ↪ "full",
    model_init_config=model_init_config,
)
```



### Use model wrapper:

When we wrap a model with model wrapper, we need to implement a model getter function. The model getter function will be called to initialize model on CPU and then model will be moved to XLA device serially. Then, let's pass `nxd_config`, model getter function and its inputs to method `initialize_parallel_model`:

```
model = nxd.initialize_parallel_model(nxd_config, get_model, config)
```

If pipeline parallel is enabled, to run a training iteration, user must use `run_train`, it handles pipeline partitioned forward and backward in it:

```
loss = model.run_train(*inputs)
```

Otherwise, users can use either `run_train` or:

```
loss = model(*inputs)
loss.backward()
```

To access the wrapped model:

```
model.local_module()
```

Model wrapper also has short cuts to access some common fields of hugging face transformers model;

```
model.dtype # get model's dtype
model.config # get model's config
model.name_or_path # get model's name or path
```

### Use optimizer wrapper:

When we wrap an optimizer with optimizer wrapper, we need `nxd_config`, original optimizer class and its inputs (parameters and optimizer arguments):

```
optimizer = nxd.initialize_parallel_optimizer(
    nxd_config, torch.optim.AdamW, param_groups, lr=args.lr, betas=(args.beta1, args.
    ↪ beta2), weight_decay=args.weight_decay
)
```

One useful feature is that user can access grad norm value from wrapped optimizer directly:

```
# It's a XLA tensor
optimizer.grad_norm
```

Note that if optimizer has not been executed or `grad_clipping` is disable, access `grad_norm` will get `None`.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Developer guide for LoRA finetuning

This document will introduce how to enable model finetuning with LoRA.

For a complete api guide, refer to [API](#).

### Enable LoRA finetuning:

We first set up LoRA-related configurations:

```
lora_config = nxd.modules.lora.LoraConfig(  
    enable_lora=True,  
    lora_rank=16,  
    lora_alpha=32,  
    lora_dropout=0.05,  
    bias="none",  
    lora_verbose=True,  
    target_modules=["q_proj", "v_proj", "k_proj"],  
    save_lora_base=False,  
    merge_lora=False,  
)
```

The default target modules for different model architectures can be found in [model.py](#).

We then initialize NxD model with LoRA enabled:

```
nxd_config = nxd.neuronx_distributed_config(  
    ...  
    lora_config=lora_config,  
)  
model = nxd.initialize_parallel_model(nxd_config, ...)
```

### Save LoRA checkpoint

Users can save the LoRA adapter with

```
nxd.save_checkpoint(  
    checkpoint_dir_str=checkpoint_dir, # checkpoint path  
    tag=tag, # sub-directory under checkpoint path  
    model=model  
)
```

Because `save_lora_base=False` and `merge_lora=False`, only the LoRA adapter is saved under `checkpoint_dir/tag/`. We can also set `merge_lora=True` to save the merged model, i.e., merging LoRA adapter into the base model.

**Load LoRA checkpoint:**

A sample usage:

```
lora_config = LoraConfig(
    enable_lora=True,
    load_lora_from_ckpt=True,
    lora_save_dir=checkpoint_dir, # checkpoint path
    lora_load_tag=tag, # sub-directory under checkpoint path
)
nxd_config = nxd.neuronx_distributed_config(
    ...
    lora_config=lora_config,
)
model = nxd.initialize_parallel_model(nxd_config, ...)
```

The NxD model will be initialized with LoRA enabled and LoRA weights loaded. LoRA-related configurations are the same as the LoRA adapter checkpoint.

*This document is relevant for: Inf2, Trn1, Trn2*

- [Developer guide for Tensor Parallelism](#)
- [Developer guide for Pipeline Parallelism](#)
- [Developer guide for Activation Memory reduction](#)
- [Developer guide for save/load checkpoint](#)
- [Developer guide for Neuron-PT-Lightning](#)
- [Developer guide for model and optimizer wrapper](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**Inference Developer Guide**

*This document is relevant for: Inf2, Trn1, Trn2*

**Developer guide for Neuronx-Distributed Inference**

Neuronx-Distributed (NxD Core) provides fundamental building blocks that enable you to run advanced inference workloads on AWS Inferentia and Trainium instances. These building blocks include parallel linear layers that enable distributed inference, a model builder that compiles PyTorch modules into Neuron models, and more.

Neuron also offers Neuronx-Distributed (NxD) Inference, which is a library that provides optimized model and module implementations that build on top of NxD Core. We recommend that you use NxD Inference to run inference workloads and onboard custom models. For more information about NxD Inference, see [NxD Inference Overview](#).

For examples of how to build directly on NxD Core, see the following:

- [Llama 3.2 1B inference sample](#)
- [T5 3B inference tutorial \[html\]](#) [\[notebook\]](#)

*This document is relevant for: Inf2, Trn1, Trn2*

- [Developer guide for Neuronx-Distributed Inference](#)

*This document is relevant for:* Inf2, Trn1, Trn2

- [Training Developer Guides](#)
- [Inference Developer Guide](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### 3.3.5 Tutorials for NeuronX Distributed

*This document is relevant for:* Inf2, Trn1, Trn2

#### Training Tutorials

*This document is relevant for:* Inf2, Trn1, Trn2

#### Training with Tensor Parallelism

Keeping the above changes made in [Developer guide](#), let's now run an end-to-end training with tensor-parallelism. This section is adopted from [BERT pretraining tutorial](#) which used data-parallel training to scale the throughput. In this section we modify that tutorial to showcase the use of tensor-parallelism which should enable us to scale the size of the model.

Setting up environment:

For this experiment, we will use a trn1-32xl machine with the storage set to 512GB at least. Follow the instructions mentioned here: [Install PyTorch Neuron on Trn1](#). It is recommended to work out of python virtual env so as to avoid package installation issues.

We also have to install the neuronx-distributed package using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.  
↪amazonaws.com
```

Make sure the transformers version is set to 4.26.0 (Note: If you have transformers-neuronx in your environment, you need to uninstall it to avoid a conflict with the transformers version.)

Let's download the scripts and datasets for pretraining.

```
mkdir -p ~/examples/tp_dp_bert_hf_pretrain  
cd ~/examples/tp_dp_bert_hf_pretrain  
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/  
↪training/tp_dp_bert_hf_pretrain/tp_dp_bert_large_hf_pretrain_hdf5.py  
wget https://raw.githubusercontent.com/aws-neuron/neuronx-distributed/master/examples/  
↪training/tp_dp_bert_hf_pretrain/requirements.txt  
python3 -m pip install -r requirements.txt
```

Next let's download the tokenizer and the sharded datasets:

```
mkdir -p ~/examples_datasets/  
pushd ~/examples_datasets/  
aws s3 cp s3://neuron-s3/training_datasets/bert_pretrain_wikicorpus_tokenized_hdf5/bert_  
↪pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar . --no-sign-request
```

(continues on next page)

(continued from previous page)

```
tar -xf bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar
rm bert_pretrain_wikicorpus_tokenized_hdf5_seqlen128.tar
popd
```

At this point, you are all set to start training

#### Running training

We first pre-compile the graphs using the `neuron_parallel_compile`. This process is similar to one discussed in the [BERT pretraining tutorial](#). Let's run the command below:

```
cd ~/examples/tp_dp_bert_hf_pretrain
export XLA_DOWNCAST_BF16=1
neuron_parallel_compile torchrun --nproc_per_node=32 \
tp_dp_bert_large_hf_pretrain_hdf5.py \
--tensor_parallel_size 8 \
--steps_this_run 10 \
--batch_size 64 \
--grad_accum_usteps 64 |& tee compile_log.txt
```

This script uses a tensor-parallel size of 8. This will automatically set the data-parallel degree to 4 (32 workers / tensor\_parallel\_size). Once the graphs are compiled we can now run training and observe our loss go down. To run the training, we just use the above command but without `neuron_parallel_compile`.

```
XLA_DOWNCAST_BF16=1 torchrun --nproc_per_node=32 \
tp_dp_bert_large_hf_pretrain_hdf5.py \
--tensor_parallel_size 8 \
--steps_this_run 10 \
--batch_size 64 \
--grad_accum_usteps 64 |& tee training_log.txt
```

You would notice that the throughput is lower when you run the `dp_bert_large_hf_pretrain_hdf5.py`. This is expected as the number of data-parallel workers have gone down (from 32 to 4). However, if you open `neuron-top` in another terminal, you should see the memory utilization per core for this script is lower than the `dp_bert_large_hf_pretrain_hdf5.py`. Since the memory requirement has gone down, you can scale the size of model either by increasing the number of layers/attention heads/hidden sizes.

The loss curve should match to the loss curve we would get from the data\_parallel counterpart.

#### Known Issues:

1. Currently the checkpoints dumped during training are sharded and users would have to write a script to combine the checkpoints themselves. This should be fixed in the future release

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Training GPT-NeoX 6.9B with Tensor Parallelism and ZeRO-1 Optimizer

In this section, we showcase to pretrain a GPT-NeoX 6.9B model by using tensor parallelism and zero-1 optimizer in the neuronx-distributed package. Please refer to the [Neuron Samples repository](#) to view the files in this tutorial.

### Setting up environment:

For this experiment, we will use a ParallelCluster with at least four trn1-32xl compute nodes. [Train your model on ParallelCluster](#) introduces how to setup and use a ParallelCluster. We need first to create and activate a python virtual env on the head node of the ParallelCluster. Next follow the instructions mentioned here: [Install PyTorch Neuron on Trn1](#) to install neuron python packages.

We also need to install and clone the neuronx-distributed package using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.
↪amazonaws.com
git clone git@github.com:aws-neuron/neuronx-distributed.git
```

Let's download the scripts for pretraining.

```
cd ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/tp_dp_gpt_neox_6.
↪9b_hf_pretrain/
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/adamw_
↪fp32_optim_params.py ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/get_
↪dataset.py ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/
↪requirements.txt ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/tp_dp_gpt_neox_
↪20b_hf_pretrain/modeling_gpt_neox_nxd.py ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/tp_dp_gpt_neox_
↪20b_hf_pretrain/utils.py ./
python3 -m pip install -r requirements.txt
```

Next let's download and pre-process the dataset:

```
python3 get_dataset.py
```

At this point, you are all set to start training.

### Running training

We first pre-compile the graphs using the neuron\_parallel\_compile. Let's run the command below:

```
sbatch --exclusive \
--nodes 4 \
--wrap="srun neuron_parallel_compile bash $(pwd)/tp_dp_gpt_neox_6.9b_hf_pretrain.sh"
```

This script uses a tensor-parallel size of 8. This will automatically set the zero-1 sharding degree to 16 ( $4 * 32$  workers / tensor\_parallel\_size). Once the graphs are compiled we can now run training and observe our loss goes down. To run the training, we just the above command but without neuron\_parallel\_compile.

```
sbatch --exclusive \
--nodes 4 \
--wrap="srun bash $(pwd)/tp_dp_gpt_neox_6.9b_hf_pretrain.sh"
```

### ZeRO-1 Optimizer

The training script uses ZeRO-1 optimizer, where the optimizer states are partitioned across the ranks so that each rank updates only its partition. Below shows the code snippet of using ZeRO-1 optimizer in training script:

```
from neuronx_distributed.optimizer import NeuronZero1Optimizer

optimizer = NeuronZero1Optimizer(
    optimizer_grouped_parameters,
    AdamW_FP32OptimParams,
    lr=flags.lr,
    pin_layout=False,
    sharding_groups=parallel_state.get_data_parallel_group(as_list=True),
    grad_norm_groups=parallel_state.get_tensor_model_parallel_group(as_list=True),
)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Training GPT-NeoX 20B with Tensor Parallelism and ZeRO-1 Optimizer

In this section, we showcase to pretrain a GPT-NeoX 20B model by using the sequence parallel optimization of tensor parallelism in the neuronx-distributed package. Please refer to the [Neuron Samples repository](#) to view the files in this tutorial.

This GPT-NeoX 20B tutorial differs from the [GPT-NeoX 6.9B tutorial](#) in the following ways:

- sequence parallel optimization has been applied
- parallel cross entropy has been applied
- the model size has been increased from 6.9B to 20B
- the TP degree has been increased from 8 to 32

Setting up environment is same as the [GPT-NeoX 6.9B tutorial](#).

**Let's download the scripts for pretraining:**

```
cd ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/tp_dp_gpt_neox_20b_
  hf_pretrain/
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/adamw_
  fp32_optim_params.py ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/get_
  dataset.py ./
ln -sf ~/neuronx-distributed/examples/training/tp_dp_gpt_neox_hf_pretrain/common/
  requirements.txt ./
python3 -m pip install -r requirements.txt
```

Next let's download and pre-process the dataset:

```
python3 get_dataset.py
```

At this point, you are all set to start training.

### Running training

We first pre-compile the graphs using the `neuron_parallel_compile`. Let's run the command below:

```
sbatch --exclusive \  
--nodes 4 \  
--cpus-per-task 128 \  
--wrap="srun neuron_parallel_compile bash $(pwd)/tp_dp_gpt_neox_20b_hf_pretrain.sh"
```

This script uses a tensor-parallel size of 32. This will automatically set the zero-1 sharding degree to 4 ( $4 * 32$  workers / `tensor_parallel_size`). Once the graphs are compiled we can now run training and observe our loss goes down. To run the training, we just the above command but without `neuron_parallel_compile`.

```
sbatch --exclusive \  
--nodes 4 \  
--cpus-per-task 128 \  
--wrap="srun bash $(pwd)/tp_dp_gpt_neox_20b_hf_pretrain.sh"
```

### Sequence Parallel

We made the following model level modifications to enable sequence parallel:

- turn on `sequence_parallel_enabled` of `ColumnParallelLinear` and `RowParallelLinear` in `GPTNeoXAttention` and `GPTNeoXMLP`;
- replace torch `LayerNorm` in `GPTNeoXLayer` and `GPTNeoXModel` with neuronx-distributed `LayerNorm` with `sequence_parallel_enabled` turned on;
- dimension transposition of intermediate states in the forward function of `GPTNeoXAttention`.
- dimension transposition and collective communication of intermediate states in the forward function of `GPTNeoXModel`.

In the training script level, we enable:

- all-reduce sequence parallel gradients at the gradient accumulation boundary.

Please check `modeling_gpt_neox_nxd.py` and `tp_dp_gpt_neox_20b_hf_pretrain.py` for details.

### Parallel Cross Entropy

To enable parallel cross entropy, we made the following model level modifications:

- replace the `CrossEntropyLoss` with neuronx-distributed `parallel_cross_entropy` in the forward function of `GPTNeoXForCausalLM`.
- use `ColumnParallelLinear` for the `embed_out` layer in `GPTNeoXForCausalLM`.

Please check `modeling_gpt_neox_nxd.py` for details.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Training Llama3.1-8B, Llama3-8B and Llama2-7B with Tensor Parallelism and ZeRO-1 Optimizer

In this section, we showcase how to pre-train Llama3.1-8B, Llama3 8B and Llama2 7B model on four Trn1.32xlarge instances using the Neuron Distributed library. We will use AWS ParallelCluster to orchestrate the training jobs. To train the Llama model in this example, we will apply the following optimizations using the Neuron Distributed library:

1. *Tensor Parallelism*
2. *Sequence Parallel*
3. *Selective checkpointing*



#### 4. ZeRO-1

##### Setting up environment:

For this experiment, we will use AWS ParallelCluster with at least four Trn1.32xlarge compute nodes. [Train your model on ParallelCluster](#) introduces how to setup and use a ParallelCluster. To setup the packages on the headnode of the ParallelCluster, follow the instructions mentioned here: [Install PyTorch Neuron on Trn1](#).

We also need to install and clone the neuronx-distributed package inside the virtual env using the following commands:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.
↪amazonaws.com
git clone git@github.com:aws-neuron/neuronx-distributed.git
```

Let's download the scripts for pretraining:

1. Navigate to a directory to hold our experiments

```
cd ~/neuronx-distributed/examples/training/llama/tp_zero1_llama_hf_pretrain
```

2. Link the training scripts for our experiments

```
ln -sf ~/neuronx-distributed/examples/training/llama/training_utils.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/modeling_llama_nxd.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/get_dataset.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/requirements.txt ./
```

If you want to pre-train Llama3.1 8B, you would need to run the following steps -

```
chmod +x tp_zero1_llama3_8B_hf_pretrain.sh
cp ./8B_config_llama3.1/config.json ./8B_config_llama3
ln -sf 8B_config_llama3.1/config.json ./
```

If you want to pre-train Llama3 8B, you would need to run the following steps -

```
chmod +x tp_zero1_llama3_8B_hf_pretrain.sh
ln -sf 8B_config_llama3/config.json ./
```

If you want to pre-train Llama2 7B, run the following steps -

```
chmod +x tp_zero1_llama2_7B_hf_pretrain.sh
ln -sf 7B_config_llama2/config.json ./
```

3. Installing the additional requirements

```
python3 -m pip install -r requirements.txt
```

To tokenize the data, we must request the tokenizer from hugging face and meta by following the instructions at the following link: [HuggingFace Llama 3 8B Model](#).

Use of the Llama models is governed by the Meta license. In order to download the model weights and tokenizer, please visit the above website and accept their License before requesting access. After access has been granted, you may use the following python3 script along with your own hugging face token to download and save the tokenizer.

Run the following from ~/examples/tp\_zero1\_llama\_hf\_pretrain directory:

```

from huggingface_hub import login
from transformers import AutoTokenizer

login(token='your_own_hugging_face_token')

tokenizer = AutoTokenizer.from_pretrained('meta-llama/Meta-Llama-3-8B')
# For llama2 uncomment line below
# tokenizer = AutoTokenizer.from_pretrained('meta-llama/Llama-2-7b-hf')

tokenizer.save_pretrained(".")

```

For Llama3.1/Llama3, make sure your `~/examples/tp_zero1_llama_hf_pretrain` directory has the following files:

```
'./tokenizer_config.json', './special_tokens_map.json', './tokenizer.json'
```

For Llama2, you just copy the `tokenizer.model` to the `~/examples/tp_zero1_llama_hf_pretrain` directory. Next let's download and pre-process the dataset:

```
python3 get_dataset.py --llama-version 3 # change the version number to 2 for Llama-2,
↪models
```

*Note:* In case you see an error of the following form when downloading data: `huggingface_hub.utils._validators.HFValidationError: Repo id must be in the form 'repo_name' or 'namespace/repo_name': '/home/ubuntu/examples/tp_zero1_llama_hf_pretrain'`. Use `'repo_type'` argument if needed. This could be because of a stale cache. Try deleting the cache using:

```
sudo rm -rf /home/ubuntu/.cache/
```

At this point, you are all set to start training. The below tutorial uses Llama3 8B as an example. To run Llama2 7B, simply change the script from `tp_zero1_llama3_8B_hf_pretrain.sh` to `tp_zero1_llama2_7B_hf_pretrain.sh`

## Running training

By this step, the `ParallelCluster` is all setup for running experiments. Before we run training, we first pre-compile the graphs using the [neuron\\_parallel\\_compile](#). Let's run the command below:

```

sbatch --exclusive \
--nodes 4 \
--cpus-per-task 128 \
--wrap="srun neuron_parallel_compile bash $(pwd)/tp_zero1_llama3_8B_hf_pretrain.sh"

```

This script uses a tensor-parallel size of 8. This will automatically set the zero-1 sharding degree to 16 ( $4 * 32$  workers / tensor\_parallel\_size).

*Note:* You can use any number of nodes in this case, would just need to adjust the number of nodes in the above slurm command accordingly. Also, the number of nodes used in `parallel_compile` command should be same as the actual training run. This is because, as the number of nodes change, the data-parallel degree would change too. This would result in more workers participating in operations like *gradient all-reduce* which would result in new graphs getting created.

Once the graphs are compiled we can now run training and observe our loss goes down. To run the training, we just run the above command but without `neuron_parallel_compile`.

```

sbatch --exclusive \
--nodes 4 \
--cpus-per-task 128 \
--wrap="srun bash $(pwd)/tp_zero1_llama3_8B_hf_pretrain.sh"

```

## Performance:

To achieve better performance, the script applies few techniques:

### *Sequence Parallelism and Selective Activation Checkpointing*

As explained in the [Activation Memory Recomputation Doc](#), both *Sequence Parallelism* and *Selective activation checkpointing* can help with activation memory reduction thereby allowing us to fit bigger models with less number of devices. Please refer to [Activation Memory Reduction Developer Guide](#) on how to enable sequence parallel and selective activation checkpointing.

### *Coalescing Q, K, V layers:*

We coalesced parallel matrix multiply to improve throughput:

- We coalesced query, key and value into one matrix multiply
- We coalesced gate\_proj and up\_proj into one matrix multiply

Please check `modeling_llama_nxd.py` and `tp_dp_gpt_neox_20b_hf_pretrain.py` for details. *Note:* Because we coalesced the layers above, the [pretrained checkpoint provided here](#) cannot be loaded out of the box for fine-tuning, and would require preprocessing. The Q,K,V layers and the gate\_proj and up\_proj layers need to be coalesced in the checkpoint before loading.

### *Logging:*

Currently for better performance we log loss values every 10 steps. Logging frequently will result in frequent syncs between device and CPU which are expensive. Hence, it is recommended to do less frequent logging if possible.

### *Flash Attention:*

We're introducing flash attention function for better performance/memory efficiency. Currently it's enabled by default, to disable it set `--use_flash_attention 0`

## Checkpointing:

Currently by default, the checkpoint is saved at the end of training. You can modify that behaviour by saving the checkpoint after every *N steps* inside the training loop:

```

from neuronx_distributed.parallel_layers import checkpointing
if global_step % every_n_steps_checkpoint == 0:
    state_dict = {
        "model": model.state_dict(),
        "global_step": global_step,
        "epoch": epoch,
        "scheduler": scheduler.state_dict()
    }
    checkpointing.save(state_dict, flags.output_dir)
    optimizer.save_sharded_state_dict(flags.output_dir)

```

Here we have to save the model `state_dict` using the `checkpointing.save` API and the optimizer `state_dict` using the `optimizer.save_sharded_state_dict`. This is because, currently, `checkpointing.save` API only saves on data-parallel rank 0, while in case of Zero1 Optimizer, the optimizer states are distributed across all data-parallel ranks. Hence, we use Zero1 Optimizer's save API to save the optimizer states.

*Time to save a checkpoint:*

Checkpoint save time can vary depending on what location the checkpoint is saved. If the checkpoint is saved in the `home` directory, the checkpointing time can be higher. The same time can be reduce by 4x if the checkpoint is dumped to FSX file system.

By default, `checkpoint.save` API allows one tensor-parallel rank at a time to save the checkpoint. This is done in order to avoid HOST OOM. When all tensor-parallel ranks try to save at the same time, they would end up copying weights to CPU at the same time. This can result in HOST OOM. *Note:* Since, we use `XLA_DOWNCAST_BF16` flag for BF16 training, even though the weights on device are on bf16, the weights on CPU are copied in FP32 format. In case, you want to avoid this typecasting from BF16 to FP32 when copying weights from device to CPU for checkpoint saving, you can pass `down_cast_bf16=True` to the `checkpointing.save` API as follows:

```
from neuronx_distributed.parallel_layers import checkpointing
if global_step % every_n_steps_checkpoint == 0:
    state_dict = {
        "model": model.state_dict(),
        "global_step": global_step,
        "epoch": epoch,
        "scheduler": scheduler.state_dict()
    }
    checkpointing.save(state_dict, flags.output_dir, down_cast_bf16=True)
```

This should not only reduce the HOST memory pressure when saving weights, but at the same time reduce model checkpointing time by half. *Note:* We are saving checkpoint in sharded format, wherein each tensor-parallel rank is saving one shard. To deploy these pretrained models, one would have to combine these shards by loading them and concatenating the tensor-parallel layers together. (We are working on a checkpoint conversion script that combines the shards into a single checkpoint)

In addition to the above method, if we want to speed up checkpoint saving for the model further, we can do so by:

```
from neuronx_distributed.parallel_layers import checkpointing
if global_step % every_n_steps_checkpoint == 0:
    state_dict = {
        "model": model.state_dict(),
        "global_step": global_step,
        "epoch": epoch,
        "scheduler": scheduler.state_dict()
    }
    checkpointing.save(state_dict, flags.output_dir, down_cast_bf16=True, save_xser=True)
```

The `save_xser` uses torch-xla's `xser.save` to save the tensors serially. This API will copy one tensor at a time to the disk. This will allow all the ranks to save the checkpoint at the same time. This speeds up checkpoint saving especially for large models as all ranks are saving at the same time. Moreover, the risk of HOST OOM is completely eliminated because only one tensor is copied to CPU at a time.

*Note:* If we use `save_xser` to save the checkpoint, we would have to pass `load_xser` to the `checkpoint.load` API. Also, if you use `save_xser`, the checkpoint folder would contain a `.pt` file for each tensor instead of a single `.pt` for the entire `state_dict`. To read this checkpoint in your checkpoint conversion script, you would have to use `xser.load` API instead of `torch.load` to load the checkpoint. The `xser.load` should load the serialized checkpoint and return the full `state_dict`.

Finally, to speed up optimizer saving time, you can increase the number of workers saving at the same time. This can

be done as follows:

```
if global_step % every_n_steps_checkpoint == 0:
    ...
    optimizer.save_sharded_state_dict(flags.output_dir, num_workers_per_step=32)
```

By default, `num_workers_per_step` is set to 8.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Training Llama-3.1-70B, Llama-3-70B or Llama-2-13B/70B with Tensor Parallelism and Pipeline Parallelism

In this section, we showcase to pretrain Llama 3.1, Llama3 70B and Llama2 13B/70B model by using the tensor parallel, pipeline parallel, sequence parallel, activation checkpoint as well as constant mask optimization in the `neuronx-distributed` package.

Setting up environment:

For this experiment, we will use a `ParallelCluster` with at least 32 trn1-32xl compute nodes. [Train your model on ParallelCluster](#) introduces how to setup and use a `ParallelCluster`.

We also need to install the `neuronx-distributed` package using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.
amazonaws.com
git clone git@github.com:aws-neuron/neuronx-distributed.git
```

Let's download the scripts for pretraining:

```
cd ~/neuronx-distributed/examples/training/llama/tp_pp_llama_hf_pretrain
ln -sf ~/neuronx-distributed/examples/training/llama/lr.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/training_utils.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/convert_checkpoints.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/get_dataset.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/modeling_llama_nxd.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/requirements.txt ./
```

If you want to pre-train Llama3.1 70B, you would need to run the following steps -

```
chmod +x run_llama3_70B_tp_pp.sh
ln -sf 70B_config_llama3.1/config.json ./
```

If you want to pre-train Llama3 70B, you would need to run the following steps -

```
chmod +x run_llama3_70B_tp_pp.sh
ln -sf 70B_config_llama3/config.json ./
```

For llama2 13B, you would need to run the following steps -

```
chmod +x run_llama2_13B_tp_pp.sh
ln -sf 13B_config_llama2/config.json ./
```

If you want to pre-train Llama2 70B, you would need to run the following steps -

```
chmod +x run_llama2_70B_tp_pp.sh
ln -sf 70B_config_llama2/config.json ./
```

The below tutorial uses Llama3.1 70B as an example. To run Llama2 70B or 13B, simply change the script from `run_llama3_70B_tp_pp.sh` to `run_llama2_70B_tp_pp.sh` or `run_llama2_13B_tp_pp.sh`.

First, let's get all the needed dependencies

```
python3 -m pip install -r requirements.txt
```

To tokenize the data, we must request the tokenizer from hugging face and meta by following the instructions at the following link: [HuggingFace Llama 3 8B Model](#).

Use of the Llama models is governed by the Meta license. In order to download the model weights and tokenizer, please visit the above website and accept their License before requesting access. After access has been granted, you may use the following python3 script along with your own hugging face token to download and save required tokenizer.

Run the following from `~/examples/tp_pp_llama_hf_pretrain` directory:

```
from huggingface_hub import login
from transformers import AutoTokenizer

login(token='your_own_hugging_face_token')

tokenizer = AutoTokenizer.from_pretrained('meta-llama/Meta-Llama-3-8B')
# For llama2 uncomment line below
# tokenizer = AutoTokenizer.from_pretrained('meta-llama/Llama-2-7b-hf')

tokenizer.save_pretrained(".")
```

For Llama3.1/Llama3, make sure your `~/examples/tp_pp_llama2_hf_pretrain` directory has the following files:

```
'./tokenizer_config.json', './special_tokens_map.json', './tokenizer.json'
```

For Llama2, you can just copy the `tokenizer.model` to the `~/examples/tp_pp_llama2_hf_pretrain` directory.

Next let's download and pre-process the dataset:

```
python3 get_dataset.py --llama-version 3 # change the version number to 2 for Llama-2,
↪models
```

In case you see an error of the following form when downloading data: `huggingface_hub.utils._validators.HFValidationError: Repo id must be in the form 'repo_name' or 'namespace/repo_name': '/home/ubuntu/examples/tp_pp_llama2_hf_pretrain'. Use 'repo_type' argument if needed.` This could be because of a stale cache. Try deleting the cache using:

```
sudo rm -rf /home/ubuntu/.cache/
```

In case you see an error of the following form when downloading data: ``NotImplementedError: Loading a dataset cached in a LocalFileSystem is not supported.`` Try upgrading pip:

```
pip install -U datasets
```

At this point, you are all set to start training.

Running training

We first pre-compile the graphs using the `neuron_parallel_compile`. Let's run the command below:

```

sbatch --exclusive \
--nodes 32 \
--cpus-per-task 128 \
--wrap="srun neuron_parallel_compile bash $(pwd)/run_llama3_70B_tp_pp.sh"

```

This script uses a tensor-parallel size of 8, pipeline-parallel size of 8 To run the training, we just use the above command but without `neuron_parallel_compile`.

```

sbatch --exclusive \
--nodes 32 \
--cpus-per-task 128 \
--wrap="srun bash $(pwd)/run_llama3_70B_tp_pp.sh"

```

To achieve better performance, the script applies few techniques:

#### *Sequence Parallelism and Selective Activation Checkpointing*

As explained in the [Activation Memory Recomputation Doc](#), both *Sequence Parallelism* and *Selective activation checkpointing* can help with activation memory reduction thereby allowing us to fit bigger models with less number of devices. Please refer to [Activation Memory Reduction Developer Guide](#) on how to enable sequence parallel and selective activation checkpointing.

#### *GQAQKVColumnParallelLinear Layer:*

In LLama 70B GQA module, the K and V attention heads are 8 whereas Q has 64 attentions heads. Since the number of attention heads should be divisible by `tensor_parallel_degree`, we would end up using a `tp_degree` of 8. Hence to fit a 70B model, we would have to use a higher pipeline-parallel degree. Using higher pipeline-parallel degree works well when the global batch size is very high, however, as the data-parallel degree increases at higher cluster size, the batch size per node decreases. This would result in higher [pipeline bubble](#) thereby reducing performance. To mitigate this issue, one can use the *GQAQKVColumnParallelLinear* layer with the `kv_size_multiplier` set to 4. This would repeat the KV heads and make them 32. This would allow doing tensor-parallelism using `tp_degree` of 32. This reduces the activation memory per device and thereby eventually allows using a pipeline parallel degree of 4. This can be enabled by passing the argument:

```

torchrun $DISTRIBUTED_ARGS run_llama_nxd.py \
... \
--qkv_linear 1 \
--kv_replicator 4 \
--tb_dir $tb_dir |& tee $LOG_PATH/log

```

The above changes are already included in the `run_llama_70b_tp_pp.sh`. For Llama13B model we only do 8-way tensor parallelism so we do not need this change.

#### *Fusing Q,K,V layers:*

In the *GQAQKVColumnParallelLinear*, the parallel matrix multiply is coalesced to improve throughput. Currently it's enabled by default. To disable it, set `--fuse_qkv 0`

*Note:* Because the layers above are coalesced, ensure that any pretrained checkpoint loaded for fine-tuning has the q,k,v layers coalesced. Otherwise, preprocessing is required to fuse these layers in the checkpoint. Follow this [Checkpoint Conversion Guide](#) and set `--fuse_qkv` to coalesce the layers in the checkpoint.

#### *Flash Attention:*

We're introducing flash attention function for better performance/memory efficiency. Currently it's enabled by default, to disable it set `--use_flash_attention 0`

*Save/Load Checkpoint* (refer to [API GUIDE](#) for more context about checkpoint APIs):

To enable checkpoint saving, add the following flags to `run_llama_70b_tp_pp.sh`:

- `--checkpoint_freq` Number of steps to save a checkpoint, set to -1 to disable saving checkpoint, should set as -1 when pre-compiling graph
- `--checkpoint_dir` Direction to save the checkpoint
- `--num_kept_checkpoint` Number of checkpoints to save, older checkpoint will be deleted manually, set to -1 to keep all saved checkpoints.
- `--save_load_xser` Save with torch xla serialization to reduce time saving, it's recommended to enable xser for significantly faster save/load
- `--async_checkpoint_saving` Whether to use asynchronous checkpoint saving to reduce saving time.

To enable checkpoint loading, add the following flags to `run_llama_70b_tp_pp.sh`:

- `--loading_step` Step to retrieve checkpoint from, set to -1 to disable checkpoint loading. Set to `latest_if_exists` to load the latest checkpoint under `checkpoint_dir`.
- `--checkpoint_dir` Direction to load the checkpoint from
- `--save_load_xser` load with torch xla serialization to reduce time saving, it's recommended to enable xser for significantly faster save/load. Note that if the checkpoint is saved with xser, it can only be loaded with xser, vice versa.

Load pretrained model:

We also provide option to load from pretrained HF model. Before loading, convert the full model to sharded model with `convert_checkpoints.py`:

```
python3 convert_checkpoints.py --tp_size <tp_size> --pp_size <pp_size> --n_layers
↪ <number_of_layers> --input_dir <path_to_full_model> --output_dir <sharded_model_path>
↪ --convert_from_full_model
```

And add `--pretrained_weight_dir <sharded_model_path>` flag to `run_llama_70b_tp_pp.sh`

Convert sharded model to full model with `convert_checkpoints.py`:

```
python3 convert_checkpoints.py --tp_size <tp_size> --pp_size <pp_size> --n_layers
↪ <number_of_layers> --input_dir <sharded_model_dir> --output_dir <full_model_dir> --
↪ convert_to_full_model --kv_size_multiplier <kv_size_multiplier> --config config.json --
↪ qkv_linear True --load_xser True
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Training Llama-2-7B/13B/70B using Tensor Parallelism and Pipeline Parallelism with Neuron PyTorch-Lightning

In this section, we showcase to pretrain a Llama2 7B/13B/70B with Tensor Parallelism and Pipeline Parallel using Neuron PyTorch-Lightning APIs, please refer to [Llama2 7B Tutorial](#), [Llama2 13B/70B Tutorial](#) and Neuron PT-Lightning Developer Guide for more context.



## Setting up environment:

For this experiment, we will use AWS ParallelCluster with at least four trn1.32xlarge compute nodes(at least 32 nodes are needed for 13B/70B model size). [Train your model on ParallelCluster](#) introduces how to setup and use a ParallelCluster. To setup the packages on the headnode of the ParallelCluster, follow the instructions mentioned here: [Install PyTorch Neuron on Trn1](#).

We also need to install the neuronx-distributed package inside the virtual env using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.
↪amazonaws.com
git clone git@github.com:aws-neuron/neuronx-distributed.git
```

Let's download the scripts for pretraining:

1. Navigate to a directory to hold our experiments

```
cd ~/neuronx-distributed/examples/training/llama/lightning
```

2. Link the training scripts for our experiments

```
ln -sf ~/neuronx-distributed/examples/training/llama/get_dataset.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/lr.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/modeling_llama_nxd.py ./
ln -sf ~/neuronx-distributed/examples/training/llama/requirements.txt ./
ln -sf ~/neuronx-distributed/examples/training/llama/requirements_ptl.txt ./
ln -sf ~/neuronx-distributed/examples/training/llama/training_utils.py ./
```

If you want to pre-train Llama 7B, you would need to run the following steps -

```
chmod +x run_llama_7b_tp_ptl.sh
mkdir 7B_config_llama2
cp ~/neuronx-distributed/examples/training/llama/tp_zero1_llama_hf_pretrain/7B_config_
↪llama2/config.json ./7B_config_llama2
ln -sf 7B_config_llama2/config.json ./
```

If you want to pre-train Llama 13B, you would need to run the following steps -

```
chmod +x run_llama_13b_tp_pp_ptl.sh
mkdir 13B_config
cp ~/neuronx-distributed/examples/training/llama/tp_pp_llama_hf_pretrain/13B_config_
↪llama2/config.json ./13B_config
ln -sf 13B_config/config.json ./
```

If you want to pre-train Llama 70B, you would need to run the following steps -

```
chmod +x run_llama_70b_tp_pp_ptl.sh
mkdir 70B_config
cp ~/neuronx-distributed/examples/training/llama/tp_pp_llama_hf_pretrain/70B_config_
↪llama2/config.json ./70B_config
ln -sf 70B_config/config.json ./
```

3. Installing the additional requirements and giving the right permissions to our shell script

```
python3 -m pip install -r requirements.txt
python3 -m pip install -r requirements_ptl.txt # Currently we're supporting Lightning_
↪ version 2.1.0
```

Next, we tokenize our dataset. Note: To tokenize the data, we must request the tokenizer from *HuggingFace* and *Meta* by following the instructions at the following link: [HuggingFace Llama 2 7B Model](#). Use of the Llama 2 model is governed by the Meta license. In order to download the model weights and tokenizer, please visit the above website and accept their License before requesting access. After access has been granted, you may use the download scripts provided by Meta to download the model weights and tokenizer to your cluster.

Once you have downloaded the tokenizer and model weights, you can copy the `tokenizer.model` to the `~/examples/llama2_lightning` directory.

Next let's download and pre-process the dataset:

```
python3 get_dataset.py --llama-version 2
```

Note: In case you see an error of the following form when downloading data: `huggingface_hub.utils._validators.HFValidationError: Repo id must be in the form 'repo_name' or 'namespace/repo_name': '/home/ubuntu/examples/llama2_lightning'. Use 'repo_type' argument if needed.` This could be because of a stale cache. Try deleting the cache using:

```
sudo rm -rf /home/ubuntu/.cache/
```

At this point, you are all set to start training.

## Training Llama2-7B with Tensor Parallelism

By this step, the `ParallelCluster` is all setup for running experiments. Before we run training, we first pre-compile the graphs using the `neuron_parallel_compile`. Let's run the command below:

```
sbatch --exclusive \
--nodes 4 \
--cpus-per-task 128 \
--wrap="srun neuron_parallel_compile bash $(pwd)/run_llama_7b_tp_ptl.sh"
```

This script uses a tensor-parallel size of 8. This will automatically set the zero-1 sharding degree to 16 ( $4 * 32$  workers / `tensor_parallel_size`).

Note: You can use any number of nodes in this case, would just need to adjust the number of nodes in the above slurm command accordingly. Also, the number of nodes used in `parallel_compile` command should be same as the actual training run. This is because, as the number of nodes change, the data-parallel degree would change too. This would result in more workers participating in operations like *gradient all-reduce* which would result in new graphs getting created.

Once the graphs are compiled we can now run training and observe our loss goes down. To run the training, we just run the above command but without `neuron_parallel_compile`.

```
sbatch --exclusive \
--nodes 4 \
--cpus-per-task 128 \
--wrap="srun bash $(pwd)/run_llama_7b_tp_ptl.sh"
```

## Training Llama2-13B/70B with Tensor Parallelism and Pipeline Parallelism

Here we use Llama70B as an example. To run 13B, simply change the script from `run_llama_70b_tp_pp.sh` to `run_llama_13b_tp_pp.sh`. Before we run training, we first pre-compile the graphs using the [neuron\\_parallel\\_compile](#). Let's run the command below:

Pre-compiling

```
sbatch --exclusive \
--nodes 32 \
--cpus-per-task 128 \
--wrap="srun neuron_parallel_compile bash $(pwd)/run_llama_70b_tp_pp_ptl.sh"
```

This script uses a tensor-parallel size of 8, pipeline-parallel size of 8. To run the training, we just use the above command but without `neuron_parallel_compile`.

```
sbatch --exclusive \
--nodes 4 \
--cpus-per-task 128 \
--wrap="srun bash $(pwd)/run_llama_7b_tp_ptl.sh"
```

### Checkpointing:

To enable checkpoint saving, add following flags to `run_llama_7b_tp_ptl.sh`/ `run_llama_13b_tp_pp.sh` / `run_llama_70b_tp_pp.sh`: \* `--save_checkpoint` Add this flag to enable checkpoint saving \* `--checkpoint_freq` Number of steps to save a checkpoint \* `--checkpoint_dir` Direction to save the checkpoint \* `--num_kept_checkpoint` Number of checkpoints to save, older checkpoint will be deleted manually, set to -1 to keep all saved checkpoints \* `--save_load_xser` load with torch xla serialization to reduce time saving, it's recommended to enable xser for significantly faster save/load. Note that if the checkpoint is saved with xser, it can only be loaded with xser, vice versa.

To enable checkpoint loading, add following flags to `run_llama_7b_tp_ptl.sh`/ `run_llama_13b_tp_pp.sh` / `run_llama_70b_tp_pp.sh`: \* `--resume_ckpt` \* `--load_step` Step to retrieve checkpoint from \* `--checkpoint_dir` Direction to load the checkpoint from \* `--save_load_xser` load with torch xla serialization to reduce time saving, it's recommended to enable xser for significantly faster save/load. Note that if the checkpoint is saved with xser, it can only be loaded with xser, vice versa.

*This document is relevant for: Inf2, Trn1, Trn2*

- [Training with Tensor Parallelism](#)
- [Training GPT-NeoX 6.9B with Tensor Parallelism and ZeRO-1 Optimizer](#)
- [Training GPT-NeoX 20B with Tensor Parallelism and ZeRO-1 Optimizer](#)
- [Training Llama3.1-8B, Llama3-8B and Llama2-7B with Tensor Parallelism and ZeRO-1 Optimizer](#)
- [Training Llama-3.1-70B, Llama-3-70B or Llama-2-13B/70B with Tensor Parallelism and Pipeline Parallelism](#)
- [Training Llama-2-7B/13B/70B using Tensor Parallelism and Pipeline Parallelism with Neuron PyTorch-Lightning](#)
- [llama2\\_7b\\_tp\\_zero1\\_ptl\\_finetune\\_tutorial](#)
- [llama3\\_8b\\_tp\\_ptl\\_lora\\_finetune\\_tutorial](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Inference Tutorials

### T5 inference with Tensor Parallelism

This is an extension to the [t5 inference tutorial](#). Here we will use NeuronxDistributed to improve the inference performance using tensor parallelism.

This tutorial has the following main sections:

1. Install dependencies
2. Plug in NeuronxDistributed layers into T5
3. Compile the T5 model
4. Run distributed inference with beam search

This Jupyter notebook should be run on a Inf2 instance (inf2.24xlarge) or Trn1 instance (trn1.32xlarge)

The tutorial works for t5 and flan-t5 models. In this notebook we will run distributed inference with flan-t5-xl.

### Install dependencies

The code in this tutorial is written for Jupyter Notebooks. To use Jupyter Notebook on the Neuron instance, you can use this [guide](#).

Run the notebook by cloning aws-neuron-sdk

```
git clone https://github.com/aws-neuron/aws-neuron-sdk.git
cd aws-neuron-sdk/src/examples/pytorch/neuronx_distributed/t5-inference/
```

Once done execute `t5-inference-tutorial.ipynb`

It is recommended to go through the [t5 inference tutorial](#) before you start this tutorial. In addition to the dependencies in the [t5 inference tutorial](#), we need to install neuronx-distributed.

This tutorial requires the following pip packages:

- torch-neuronx
- neuronx-cc
- transformers
- optimum-neuron
- neuronx-distributed

Most of these packages will be installed when configuring your environment using the Trn1/Inf2 [setup guide](#). The additional dependencies must be installed here:

```
[ ]: ! pip install --upgrade transformers==4.33.1 optimum-neuron neuronx_distributed --extra-
    ↪index-url https://pip.repos.neuron.amazonaws.com

[ ]: # Pull the latest version of the compiler
    ! pip install --upgrade neuronx-cc>=2.11 --no-deps

[ ]: # Lets update numpy to a newer version
    ! pip install --upgrade "numpy>=1.22.2,<2" --no-deps
```

## Plug in NeuronxDistributed layers into T5

We extend the huggingface's T5 model to use the NeuronxDistributed parallel layers. To do so, we simply swap linear layers in T5LayerSelfAttention, T5LayerCrossAttention, and T5LayerFF definitions with ColumnParallelLinear and RowParallelLinear. We also need to swap the Embedding layer with ParallelEmbedding.

Let us take the example of T5Attention. The `attention block` has q, k, v, and o linear layers. The multi-head attention block uses q, k and v to compute the attention scores. The attention scores are then passed through o to compute the attention block output. So let us swap q, k and v layers with ColumnParallelLinear and o with RowParallelLinear. Having RowParallelLinear following a ColumnParallelLinear is a performance optimization. The attention scores computed with q, k and v are already split across Neuron devices. The row parallel layer can use this shared output directly. The embedding layer is simply swapped with the ParallelEmbedding.

```
class ParallelAttention(T5Attention):
    def __init__(self, config: T5Config, has_relative_attention_bias=False):
        super().__init__(config, has_relative_attention_bias)
        # Per attention head and per partition values
        world_size = parallel_state.get_tensor_model_parallel_size()
        self.num_attention_heads_per_partition = divide(self.n_heads, world_size)
        self.hidden_size_per_partition = self.num_attention_heads_per_partition * self.
        ↪key_value_proj_dim

        # Mesh TensorFlow initialization to avoid scaling before softmax
        self.q = ColumnParallelLinear(self.d_model,
                                     self.inner_dim,
                                     bias=False,
                                     gather_output=False)
        self.k = ColumnParallelLinear(self.d_model,
                                     self.inner_dim,
                                     bias=False,
                                     gather_output=False)
        self.v = ColumnParallelLinear(self.d_model,
                                     self.inner_dim,
                                     bias=False,
                                     gather_output=False)
        self.o = RowParallelLinear(self.inner_dim,
                                   self.d_model,
                                   bias=False,
                                   input_is_parallel=True)

        if self.has_relative_attention_bias:
            self.relative_attention_bias = ParallelEmbedding(self.relative_attention_num_
            ↪buckets, self.n_heads)
            self.n_heads = self.num_attention_heads_per_partition

        ...
```

You can find the all modified T5 layers defined in `t5_model_layers.py`.

Once we have the modified T5 layers, we can plug in the T5Attention and T5LayerFF into the pretrained model. Here is how you do that.

```
def load_pretrained_with_parallel_attn(model_name):
```

(continues on next page)

(continued from previous page)

```

model = T5ForConditionalGeneration.from_pretrained(model_name, torch_dtype="auto")

# Parallel implementation of Attention modules.
from t5_model_layers import ParallelSelfAttention, ParallelFF, ParallelCrossAttention

for index, block in enumerate(model.decoder.block):
    if index == 0:
        block.layer[0] = ParallelSelfAttention(model.config,
                                                has_relative_attention_bias=True)
    else:
        block.layer[0] = ParallelSelfAttention(model.config)
        block.layer[1] = ParallelCrossAttention(model.config)
        block.layer[2] = ParallelFF(model.config)
# Load the weights into the parallel layers
neuronx_distributed.parallel_layers.load(model_name + ".pt", model, sharded=False)

return model

```

## Compile the parallel T5 model

Let us set some model parameters.

```

[ ]: model_name = "google/flan-t5-xl"
max_length = 128
num_beams = 4
tp_degree = 8 # tensor parallelism degree

```

Download and save the model that we want to trace.

```

[ ]: import torch
from transformers import T5ForConditionalGeneration

model = T5ForConditionalGeneration.from_pretrained(model_name, torch_dtype="auto")
torch.save({"model": model.state_dict()}, model_name.split("/")[-1] + ".pt")
model.config.use_cache = True

```

To run HuggingFace T5 models on Neuron, we need to make a couple of changes. Let us reuse the code from the [t5 inference tutorial](#) which makes T5 compatible with Neuron. For your convenience, the code copied into [wrapper.py](#) and [t5\\_models.py](#). This notebook will import these files.

The only change made to this code is that we use `neuronx_distributed.trace` instead of `torch_neuronx.trace`.

Let us trace the encoder and decoder.

```

[ ]: import t5_models
import neuronx_distributed
import time

# This can take up to 20 minutes
encoder_compile_start_time = time.time()
traced_encoder = t5_models.parallel_trace_encoder(model_name, max_length, num_beams, tp_
    degree)

```

(continues on next page)

(continued from previous page)

```
print("Encoder compilation time {}".format(time.time() - encoder_compile_start_time))

neuronx_distributed.trace.parallel_model_save(traced_encoder, "TracedParallelEncoder.pt")
```

```
[ ]: # This can take up to 15 minutes
decoder_compile_start_time = time.time()
traced_decoder = t5_models.parallel_trace_decoder(model, model_name, num_beams, max_
↳length, tp_degree)
print("Decoder compilation time {}".format(time.time() - decoder_compile_start_time))

neuronx_distributed.trace.parallel_model_save(traced_decoder, "TracedParallelDecoder.pt")
```

### Inference with the traced parallel T5 model

With the traced model, let us try using beam search for inference.

```
[ ]: import neuronx_distributed
from wrapper import T5Wrapper
from transformers import T5Tokenizer

num_return_sequences = 4

traced_encoder = neuronx_distributed.trace.parallel_model_load("TracedParallelEncoder.pt
↳")
traced_decoder = neuronx_distributed.trace.parallel_model_load("TracedParallelDecoder.pt
↳")

tokenizer = T5Tokenizer.from_pretrained(model_name)
model = T5Wrapper.from_pretrained(model_name)

model.encoder = traced_encoder
model.decoder = traced_decoder
setattr(model.encoder, 'main_input_name', 'input_ids') # Attribute required by beam_
↳search

output = model.parallel_infer(tokenizer=tokenizer,
                             prompt="translate English to German: Lets eat good food.",
                             max_length=max_length,
                             num_beams=num_beams,
                             num_return_sequences=num_return_sequences,
                             device="xla")

results = [tokenizer.decode(t, skip_special_tokens=True) for t in output]

print('Results:')
for i, summary in enumerate(results):
    print(i + 1, summary)
```

Results:

```
1 Lassen Sie uns gutes Essen essen.
2 Lassen Sie uns gut essen.
3 Lassen Sie uns gutes Essen zu essen.
4 Lassen Sie uns gutes Essen zu sich nehmen.
```

## Benchmarking

Let us benchmark the per token decoder latency

```
[ ]: # Let us install NeuronPerf. We will use it to measure the performance.
! pip install neuronperf --extra-index-url=https://pip.repos.neuron.amazonaws.com

[ ]: import os
import neuronperf as npf

d_model = model.config.d_model
model_dir = "TracedParallelDecoder.pt"
decoder_run_count = 128

def load_fn(model_path, **kwargs):
    return neuronx_distributed.trace.parallel_model_load(model_path)

# NeuronPerf can't see tp_degree at the moment, so just expose all cores
def env_setup_fn(*_):
    del os.environ["NEURON_RT_VISIBLE_CORES"]

def benchmark():

    # Create some sample inputs for the decoder
    decoder_input_ids = torch.ones((num_beams, 1), dtype=torch.int64)
    decoder_attention_mask = torch.ones((num_beams, max_length), dtype=torch.int32)
    encoder_attention_mask = torch.ones((num_beams, max_length), dtype=torch.int64)
    encoder_hidden_states = torch.ones((num_beams, max_length, d_model), dtype=torch.
↪float32)
    beam_idx = torch.arange(0, num_beams, dtype=torch.int64)
    beam_scores = torch.zeros((num_beams,), dtype=torch.float)

    inputs = (decoder_input_ids,
              decoder_attention_mask,
              encoder_hidden_states,
              encoder_attention_mask,
              beam_idx,
              beam_scores)

    reports = npf.benchmark(
        load_fn,
        model_dir,
        [inputs],
        batch_sizes=1,
        n_models=1,
```

(continues on next page)



(continued from previous page)

```

        max_infers=decoder_run_count,
        workers_per_model=1, # no bottleneck on model inputs, so 1 is fine
        env_setup_fn=env_setup_fn,
        multiprocessing=False,
    )

    report = reports[0]

    # let's update throughput to be tokens / second and add a new record
    latency_in_s = report["latency_ms_avg"] / 1000
    tokens_per_s = decoder_run_count / latency_in_s
    report["throughput_avg"] = tokens_per_s

    # display and save results
    npf.print_reports(reports, cols=["throughput_avg", "latency_ms_p50", "latency_ms_p99",
    ↪])
    print(f"Results saved to: {npf.write_json(reports[0])}")

benchmark()

```

Now let's benchmark inference as a whole including sampling.

```

[ ]: import os
import torch
import neuronx_distributed
import neuronperf as npf

from transformers import T5Tokenizer
from wrapper import T5Wrapper

tokenizer = T5Tokenizer.from_pretrained(model_name)

generated_token_count = 0

class Wrapper(torch.nn.Module):
    def __init__(self,
                  traced_encoder,
                  traced_decoder):
        super().__init__()
        self.model = T5Wrapper.from_pretrained(model_name)
        self.model.encoder = traced_encoder
        self.model.decoder = traced_decoder
        setattr(self.model.encoder, 'main_input_name', 'input_ids') # Attribute_
    ↪required by beam search

    def forward(self, *inputs):
        input_ids = inputs[0]['input_ids']
        attention_mask = inputs[0]['attention_mask']
        return self.model.parallel_infer(input_ids=input_ids,
                                         attention_mask=attention_mask,
                                         max_length=max_length,
                                         num_beams=num_beams,

```

(continues on next page)

(continued from previous page)

```

num_return_sequences=num_return_sequences)

def load_fn(filename, **kwargs):
    traced_encoder = neuronx_distributed.trace.parallel_model_load(filename +
    ↪ "TracedParallelEncoder.pt")
    traced_decoder = neuronx_distributed.trace.parallel_model_load(filename +
    ↪ "TracedParallelDecoder.pt")
    return Wrapper(traced_encoder, traced_decoder)

# NeuronPerf can't see tp_degree at the moment, so just expose all cores
def env_setup_fn(*_):
    del os.environ["NEURON_RT_VISIBLE_CORES"]

def preprocess_fn(inputs):

    encoding = []
    for text in inputs:
        batch_encoding = tokenizer(text,
                                   max_length=max_length,
                                   truncation=True,
                                   padding='max_length',
                                   return_tensors="pt")

        input_ids = batch_encoding['input_ids']
        attention_mask = batch_encoding['attention_mask']
        encoding.append({"input_ids": input_ids,
                        "attention_mask": attention_mask})

    return encoding

def postprocess_fn(outputs):
    output = [tokenizer.decode(seq) for seq in outputs]
    global generated_token_count
    generated_token_count = len(outputs[0])
    return output

def benchmark():
    inputs = ["summarize: The Inflation Reduction Act lowers prescription drug costs,
    ↪ health care costs, and energy costs. It's the most aggressive action on tackling the
    ↪ climate crisis in American history, which will lift up American workers and create
    ↪ good-paying, union jobs across the country. It'll lower the deficit and ask the ultra-
    ↪ wealthy and corporations to pay their fair share. And no one making under $400,000 per
    ↪ year will pay a penny more in taxes."]
    reports = npf.benchmark(
        load_fn,
        "", # Model dir
        [inputs],
        batch_sizes=1,
        n_models=1,
        max_infers=5,
        max_duration=0, # sampling can take a while, so let's not timeout
        workers_per_model=1,
        env_setup_fn=env_setup_fn,
        preprocess_fn=preprocess_fn,

```

(continues on next page)

(continued from previous page)

```

        postprocess_fn=postprocess_fn,
        multiprocessing=False,
    )

    report = reports[0]

    report["throughput_avg"] = round(generated_token_count / (report["latency_ms_avg"] / 1000), 2)
    report["latency_per_token_ms_p50"] = round((report["latency_ms_p50"])/generated_token_count, 2)
    report["latency_per_token_ms_p99"] = round((report["latency_ms_p99"])/generated_token_count, 2)

    # display and save results
    npf.print_reports(reports, cols=["throughput_avg", "latency_per_token_ms_p50", "latency_per_token_ms_p99"])
    print(f"Results saved to: {npf.write_json(report)}")

benchmark()

```

- T5 inference tutorial [\[html\]](#) [\[notebook\]](#)
- Llama 3.2 1B inference example

*This document is relevant for:* Inf2, Trn1, Trn2

- [Training Tutorials](#)
- [Inference Tutorials](#)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### 3.3.6 Misc

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## NxD Core Release Notes (neuronx-distributed)

### Table of contents

- [NxD Core \[0.13.14393\]](#)
- [NxD Core \[0.12.12111\]](#)
- [NxD Core \[0.11.0\]](#)
- [NxD Core \[0.10.1\]](#)
- [NxD Core \[0.10.0\]](#)
- [NxD Core \[0.9.0\]](#)
- [NxD Core \[0.8.0\]](#)

- [NxD Core \[0.7.0\]](#)
- [NxD Core \[0.6.0\]](#)
- [NxD Core \[0.5.0\]](#)
- [NxD Core \[0.4.0\]](#)
- [NxD Core \[0.3.0\]](#)
- [NxD Core \[0.2.0\]](#)

This document lists the release notes for Neuronx-Distributed library.

### NxD Core [0.13.14393]

Date: 6/24/2025

#### New in this release

##### Inference:

- Add `--auto-cast=none` compiler arg by default in ModelBuilder to ensure model dtypes are preserved during compilation.
- Update ModelBuilder to cast model weights based on dtypes defined in module parameters.
- Add support for PyTorch 2.7. This release includes support for PyTorch 2.5, 2.6, and 2.7.
- Other minor fixes and improvements.

##### Training:

- Added support for transformers 4.48.0

### NxD Core [0.12.12111]

Date: 5/20/2025

#### New in this release

##### Inference:

- Improve the Model Builder API. Note: The Model Builder API is in beta.
  - Add Neuron Persistent Cache support to Model Builder. Now, Model Builder caches compiled model artifacts to reduce compilation time.
  - Improve the performance of weight sharding in Model Builder to support shard-on-load in NxD Inference.
  - Improve the performance of Model Builder trace when HLO debug mode is enabled.
- Add a Llama-3.2-1B reference inference sample using NxD Core.
- Remove the deprecated NxD inference examples. You can use the NxD Inference library to run inference with on Neuron using NxD.
- Other minor fixes and improvements.

**Training:**

- Context parallel support for sequence lengths up to 32k on TRN1 (beta feature)

**General:**

- Update the package version to include additional information.

**NxD Core [0.11.0]**

Date: 4/3/2025

**New in this release****Inference:**

- Improve the performance of weight sharding by up to 60-70%, depending on the model.
- You can now configure modules to skip during quantization with the `modules_to_not_convert` argument.
- Other minor fixes and improvements.

**Training:**

- Fixed issue with wikicorpus dataset download
- Updated model load for LoRA checkpoints

**Known Issues and Limitations**

- With PT2.5, some of the key workloads like Llama3-8B training may show reduced performance when using `-llm-training` compiler flag as compared to PT2.1.

In such a case, try removing `-llm-training` flag from `NEURON_CC_FLAGS` in the `run.sh` only if using Neuron Kernel Interface.

**NxD Core [0.10.1]**

Date: 1/14/2025

**New in this release****Inference:**

- Fix an issue with sequence parallel support for quantized models.

### NxD Core [0.10.0]

Date: 12/20/2024

#### New in this release

##### Training:

- Added support for HuggingFace Llama3 70B with Trn2 instances
- Added support for PyTorch 2.5
- Added DPO support for post-training model alignment
- Added fused QKV optimization in GQA models
- Support for Mixture-of-Experts with Tensor, Sequence, and Pipeline parallelism

#### Known Issues and Limitations

- With PT2.5, some of the key workloads like Llama3-8B training may show reduced performance when using `-llm-training` compiler flag as compared to PT2.1.

In such a case, try removing `-llm-training` flag from `NEURON_CC_FLAGS` in the `run.sh`

### NxD Core [0.9.0]

Date: 09/16/2024

#### New in this release

##### Training:

- Added LoRA adaptor support
- Added support for GPU compatible precision support using ZeRO-1

##### Inference:

- Added inference example for DBRX, and Mixtral models
- Improved inference performance with sequence length autobucketing
- Improved trace time for inference examples
- Reduced memory usage by sharing weights across prefill and decode traced models

## NxD Core [0.8.0]

Date: 07/03/2024

### New in this release

- Added support for Interleave pipeline parallel. At large cluster sizes, interleave pipeline schedule should help to reduce the pipeline bubble, thereby increasing training throughput.
- Added integration with flash attention kernel for longer sequence length training. See [Llama3 8K sequence-length training sample](#).
- Added support for naive speculative decoding, enabling assistance during the token generation process by predicting tokens with a draft model and verifying the predicted tokens with the original target model. Refer to the Neuronx Distributed inference developer guide for an example.
- Added integration with flash attention kernel for longer sequence length inference. See an end to end example of CodeLlama-13b model with 16K sequence length.
- Added support for scaled inference to run for Llama-2 70b or similar sized models

### Known Issues and Limitations

- Model checkpointing saves sharded checkpoints. Users will have to write a script to combine the shards
- Validation/Evaluation with interleaved pipeline feature is not supported.
- Due to weights not being able to be shared across context encoding and token generation trace, inference scale is tested for models up to size Llama-2-70b. For model configurations above this, there is a risk of OOM errors.
- Tracing Llama-2-70b sized models for inference and loading them to device can take close to two hours. This is due to duplicate sharding of weights for both context encoding and token generation traces.

## NxD Core [0.7.0]

Date: 04/01/2024

### New in this release

- Added support for Pipeline-parallelism training using PyTorch-lightning
- Added support for fine-tuning a model and running evaluation on the fine-tuned model using optimum-neuron
- Added support for auto-partitioning the pipeline parallel stages for training large models
- Added support for async checkpointing, optimizing the checkpoint saving time.
- Added support for auto-resume from a checkpoint, in case training job crashes.
- Added support for sequence length autobucketing in inference
- Added support for inference with bfloat16
- Improved performance for Llama-2-7b inference example.

## Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

## NxD Core [0.6.0]

Date: 12/21/2023

### New in this release

- Added support for Model/Optimizer wrapper that handles the parallelization in both model and optimizer.
- Added support for PyTorch-lightning. This allows users to train models using Tensor-parallelism and Data-parallelism.
- Added new checkpoint save/load APIs that handles the parallelization and dumps/loads the checkpoint.
- Added a new QKV module which has the ability to replicate the KV heads and produce the query, key and value states.
- Reduced the model initialization time when pipeline-parallel distributed strategy is used.
- Added support for limiting max parallel compilations in `parallel_model_trace`. This resolves many out of memory errors by reducing the host memory usage.
- Added example for Llama-2-7b inference. This is still early in development and is not well-optimized. The current recommendation is to use *transformers-neuronx* for optimal performance of llama inference.

## Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.
- Pipeline-parallelism is not supported as part of PyTorch-lightning integration.

## NxD Core [0.5.0]

Date: 10/26/2023

### New in this release

- Added support for pipeline-parallelism for distributed training.
- Added support for serialized checkpoint saving/loading, resulting in better checkpoint saving/loading time.
- Added support for mixed precision training using *torch.autocast*.
- Fixed an issue with Zero1 checkpoint saving/loading.



### Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

### NxD Core [0.4.0]

Date: 9/15/2023

#### New in this release

- Added API for padding attention heads when they are not divisible by tensor-parallel degree
- Added a constant threadpool for distributed inference
- Fixed a bug with padding\_idx in ParallelEmbedding layer
- Fixed an issue with checkpoint loading to take into account the stride parameter in tensor parallel layers

### Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

### NxD Core [0.3.0]

Date: 8/28/2023

#### New in this release

- Added Zero1 Optimizer support that works with tensor-parallelism
- Added support for sequence-parallel that works with tensor-parallelism
- Added IO aliasing feature in parallel\_trace api, which can allow marking certain tensors as state tensors
- Fixed hangs when tracing models using parallel\_trace for higher TP degree

### Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

### NxD Core [0.2.0]

Date: 7/19/2023

#### New in this release

- Added parallel cross entropy loss function.

#### Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

Date: 6/14/2023

#### New in this release

- Releasing the Neuron Distributed (`neuronx-distributed`) library for enabling large language model training/inference.
- Added support for tensor-parallelism training/inference.

#### Known Issues and Limitations

- Currently the model checkpointing saves a sharded checkpoint, and users have to write a script to combine the shards.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

- [NxD Core Release Notes \(`neuronx-distributed`\)](#)

*This document is relevant for: Inf2, Trn1, Trn2*

#### Setup

*Install PyTorch Neuron on Trn1* to create a pytorch environment. It is recommended to work out of python virtual env so as to avoid package installation issues.

You can install the `neuronx-distributed` package using the following command:

```
python -m pip install neuronx_distributed --extra-index-url https://pip.repos.neuron.  
↪ amazonaws.com
```

## App Notes

- *Tensor Parallelism Overview*
- *Pipeline Parallelism Overview*
- *Activation Memory Reduction*
- `context_parallelism_overview`

## API Reference Guide

- *Distributed Strategies APIs*
- *Training APIs*
- *Inference APIs*

## Developer Guide

- *Training Developer Guides*
- *Inference Developer Guide*

## Tutorials

- *Training Tutorials*
- *Inference Tutorials*

## Misc

- *NxD Core Release Notes (neuronx-distributed)*

*This document is relevant for: Inf2, Trn1, Trn2*



## ADDITIONAL ML LIBRARIES

### 4.1 Third-party libraries

#### 4.1.1 Third-party partner libraries

AWS Neuron integrates with multiple third-party partner products that allow you to run deep learning workloads on Amazon EC2 instances powered by AWS Trainium and AWS Inferentia chips. The following list gives an overview of the third-party libraries working with AWS Neuron.

##### Table of contents

- [\*Hugging Face Optimum Neuron\*](#)
- [\*PyTorch Lightning\*](#)
- [\*AXLearn\*](#)

#### Hugging Face Optimum Neuron

Optimum Neuron bridges Hugging Face Transformers and the AWS Neuron SDK, providing standard Hugging Face APIs for [AWS Trainium](#) and [AWS Inferentia](#). It offers solutions for both training and inference, including support for large-scale model training and deployment for AI workflows. Supporting Amazon SageMaker and pre-built Deep Learning Containers, Optimum Neuron simplifies the use of Trainium and Inferentia for machine learning. This integration allows developers to work with familiar Hugging Face interfaces while leveraging Trainium and Inferentia for their transformer-based projects.

[Optimum Neuron documentation](#)

#### PyTorch Lightning

PyTorch Lightning is a deep learning framework for professional AI researchers and machine learning engineers who need maximal flexibility without sacrificing performance at scale. Lightning organizes PyTorch code to remove boilerplate and unlock scalability.

[Get Started with Lightning](#)

Use PyTorch Lightning Trainer with *NxD*.

## **AXLearn**

AXLearn is an open-source JAX-based library used by AWS Neuron for training deep learning models on AWS Trainium. Integrates with JAX ecosystem and supports distributed training.

Check [AXLearn Github repository](#)

### **4.1.2 Additional third-party libraries**

#### **NeMo**

*NxD Training* offers a **NeMo**-compatible YAML interface for training PyTorch models on AWS Trainium chips. The library supports both Megatron-LM and HuggingFace model classes through its model hub. NxD Training leverages key NeMo components, including Experiment Manager for tracking ML experiments and data loaders for efficient data processing. This library simplifies the process of training deep learning models on AWS Trainium while providing compatibility with familiar NeMo YAML Interface.

*This document is relevant for:* Inf2, Trn1

## **4.2 Transformers NeuronX (transformers-neuronx)**

*This document is relevant for:* Inf2, Trn1

### **4.2.1 Transformers NeuronX Setup (transformers-neuronx)**

If you already have setup your environment to run PyTorch NeuronX, you just need to install Transformers NeuronX library using the following instruction.

```
pip install transformers-neuronx --extra-index-url=https://pip.repos.neuron.amazonaws.com
```

If you are starting from scratch, Neuron Multi Framework DLAMI is recommended as it comes pre-installed with Transformers NeuronX virtual environment. You can refer to the instructions to launch a Neuron instance using Multi Framework DLAMI

*This document is relevant for:* Inf2, Trn1

*This document is relevant for:* Inf2, Trn1

### **4.2.2 Transformers Neuron Developer Guide (transformers-neuronx)**

*This document is relevant for:* Inf2, Trn1

## Transformers NeuronX (transformers-neuronx) Developer Guide

Transformers NeuronX for Trn1 and Inf2 is a software package that enables PyTorch users to perform large language model (LLM) *performant inference* on second-generation Neuron hardware (See: NeuronCore-v2). The Neuron performance page lists expected inference performance for commonly used Large Language Models.

### Introduction

The [Transformers NeuronX repository](#) contains the source code of the AWS Neuron Transformers integration project. As it stands now, it mainly serves the purpose of running transformer decoder inference (autoregressive sampling) workflows on the Neuron platform.

Note: This project is **actively** in development. The Neuron team is still heavily modifying the Neuron optimized module classes. The functionality provided in this repository will not maintain long-term API stability until version `>= 1.0.0`. For applications willing to reuse code from this repository, we recommend treating the Neuron optimized module implementations as samples, and pin the version of the main library package `torch-neuronx` to avoid breaking interface changes as new features are developed.

### Checkpoint compatibility with HuggingFace Transformers

`transformers-neuronx` is checkpoint-compatible with HuggingFace Transformers. While the Neuron team reimplemented some HuggingFace Transformers models from scratch for the purpose of maximizing the execution efficiency of transformer decoders on Neuron, the implementations are done with maximizing compatibility in mind, meaning one can train transformer decoder models, say GPT2, using the standard HuggingFace Transformers library, and then construct an inference-optimized decoder model using `transformers-neuronx`'s `GPT2ForSampling` class. If training was done with other libraries such as MegatronLM, then it is still possible to convert the obtained checkpoint to the standard HuggingFace Transformers checkpoint format, and then move on to `transformers-neuronx`'s optimized decoder implementations.

### Neuron optimized transformer decoders implemented in XLA High Level Operations (HLO)

Due to the stateful nature of the autoregressive sampling computation, an efficient implementation of autoregressive sampling using the Neuron SDK requires rewriting the model forward function into a pure-function computation running on fixed-shape tensors. Furthermore, we want the pure-function computation be implemented in a compiled language so that the Neuron compiler can perform extensive code analysis and optimization. We chose XLA High Level Operations (HLO) as the compiled language for implementing Neuron optimized transformer decoder classes. The source code of these classes contains Python functions written in a syntax called “PyHLO”, name of a Neuron internal tool for writing/compiling the HLO language in Python. As an example, a “language model head” implemented in PyHLO may look like the following.

```
class LmHeadHlo:
    ...

    def lm_head(self, scribe):
        dtype = self.dtype
        hidden_size = self.hidden_size
        n_active_tokens = self.n_active_tokens
        batch_size = self.batch_size
        vocab_size = self.vocab_size
        hidden = dtype[hidden_size, n_active_tokens, batch_size].Parameter(parameter_
```

(continues on next page)

(continued from previous page)

```

↪number=0)
    weight = dtype[hidden_size, vocab_size].Parameter(parameter_number=1)
    rhs_size = n_active_tokens * batch_size
    hidden = dtype[hidden_size, rhs_size].Reshape(hidden)
    dot_dims = dict(lhs_contracting_dimensions=[0], rhs_contracting_dimensions=[0])
    logits = dtype[vocab_size, rhs_size].Dot(weight, hidden, dot_dimension_
↪numbers=dot_dims)
    return dtype[vocab_size, n_active_tokens, batch_size].Reshape(logits)

...

```

The `transformers_neuronx.compiler.compile_py_func` function can convert the Python `lm_head` function into `HloModuleProto`, a valid input format for the `neuronx-cc` compiler.

## Tensor-parallelism support

For transformer decoders used in large language models, tensor-parallelism is necessary as it provides a way to shard the models' large weight matrices onto multiple NeuronCores, and having NeuronCores working on the same matrix multiply operation collaboratively. `transformers-neuronx`'s tensor-parallelism support makes heavy use of collective operations such as all-reduce, which is supported natively by the Neuron runtime.

There are some principles for setting tensor-parallelism degree (number of NeuronCores participating in sharded matrix multiply operations) for Neuron-optimized transformer decoder models.

1. The number of attention heads needs to be divisible by the tensor-parallelism degree.
2. The total data size of model weights and key-value caches needs to be smaller than 16 GB times the tensor-parallelism degree.
3. Currently, the Neuron runtime supports tensor-parallelism degrees 1, 2, 8, and 32 on Trn1 and supports tensor-parallelism degrees 1, 2, 4, 8, and 24 on Inf2.

Some examples:

1. `facebook/opt-13b` has 40 attention heads, and when running at batch size 1 and float16 precision the model requires ~29 GB memory, therefore a `trn1.2xlarge` with 32 GB device memory is sufficient.
2. `facebook/opt-30b` has 56 attention heads, and at batch size 1 and float16 precision the model requires ~66 GB memory, therefore it can run on 8 NeuronCores on one `trn1.32xlarge` using 128 GB device memory.
3. `gpt2-xl` has 25 attention heads and requires ~4 GB memory at bfloat16 precision. It runs without tensor-parallelism only.

## Features

### Compile-time Configurations

Transformers Neuron models support a variety of compile-time configurations that can be used to tune model performance. All models support the following configurations:

- `batch_size`: The batch size to compile a model for. Once the batch size has been set, this is the only size that is supported at inference time. Neuron uses ahead-of-time compilation to achieve high performance which requires that the compiled artifact shapes must be known at compilation time.



- **n\_positions**: The maximum number of positions (or sequence length) to allow during generation. This parameter directly controls the width of the KV cache. This parameter should be set to the maximum expected sequence length for the end application.
- **tp\_degree**: This parameter controls the number of tensor parallel shards to split the model into. Each shard will execute on a separate NeuronCore. To minimize latency, it is recommended to set the tensor parallelism to be equal to the number of NeuronCores that are available on an instance.
- **amp**: This allows a models weights and compute to be cast to a different type. The options are; 'bf16', 'f16', or 'f32'. For models trained in float32, the 16-bit mixed precision options ('bf16', 'f16') generally provide sufficient accuracy while significantly improving performance.
- **context\_length\_estimate**: This parameter controls the maximum sequence length of the prompt/context handling compute graph. This parameter is not supported in GPTNeoXForSampling and GPTJForSampling.

```
from transformers_neuronx import NeuronAutoModelForCausalLM

model = NeuronAutoModelForCausalLM.from_pretrained(
    'gpt2',                                # Uses the GPT2 checkpoint from https://huggingface.co/
    ↪gpt2
    batch_size=1,                          # Allow inference with batch size 1 inputs
    n_positions=128,                       # Allow a maximum size of 128 prompt & output tokens
    tp_degree=2,                           # Shard the model weights & compute across 2 NeuronCores
    amp='f16',                             # Downcast the weights & compute to float16
    context_length_estimate=64,             # Build an optimized context encoding network for a
    ↪maximum prompt size of 64
)
model.to_neuron() # Load/compile the model
```

## Checkpoint support and automatic model selection

*New in release 2.18*

Transformers Neuron now supports a greater variety of checkpoints including older pytorch binary checkpoints and newer [safetensors](#) checkpoints. For improved load speed and reduced host memory consumption, it is recommended to always use safetensors by default. Both regular and sharded variants of checkpoints are supported. It is no longer recommended to use the `save_pretrained_split` function which was used in older Transformers Neuron examples.

In addition to supporting standard checkpoint formats, Transformers Neuron provides an `AutoModel` class `NeuronAutoModelForCausalLM` which can be used to load the correct model without explicitly importing the architecture-specific class.

```
from transformers_neuronx import NeuronAutoModelForCausalLM

# Loads: https://huggingface.co/bigscience/bloom-560m
bloom = NeuronAutoModelForCausalLM.from_pretrained('bigscience/bloom-560m')
bloom.to_neuron()

# Loads: https://huggingface.co/openlm-research/open_llama_3b_v2
llama = NeuronAutoModelForCausalLM.from_pretrained('openlm-research/open_llama_3b_v2')
llama.to_neuron()

# This is equivalent to the following:
from transformers_neuronx import BloomForSampling
model = BloomForSampling.from_pretrained('bigscience/bloom-560m')
```

(continues on next page)

(continued from previous page)

```
model.to_neuron()

from transformers_neuronx import LlamaForSampling
llama = LlamaForSampling.from_pretrained('openlm-research/open_llama_3b_v2')
llama.to_neuron()
```

**Note:** Advanced features of huggingface hub access are not supported. This includes private repositories which require access tokens and branches.

In order to support more advanced repository downloads, please download the model to a local directory and load it from there.

### Hugging Face generate() API support

Transformers Neuron models support the Hugging Face `generate()` API via the `HuggingFaceGenerationModelAdapter` adapter class. In the following example we demonstrate how to run sampling with temperature using the GPT2 model:

```
import torch
from transformers import AutoTokenizer, AutoConfig
from transformers_neuronx import GPT2ForSamplingWithContextBroadcasting,
↳ HuggingFaceGenerationModelAdapter

# Create and compile the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained('gpt2')
model.to_neuron()

# Use the `HuggingFaceGenerationModelAdapter` to access the generate API
config = AutoConfig.from_pretrained('gpt2')
model = HuggingFaceGenerationModelAdapter(config, model)

# Get a tokenizer and example input
tokenizer = AutoTokenizer.from_pretrained('gpt2')
tokenizer.pad_token_id = tokenizer.eos_token_id
tokenizer.padding_side = 'left'
text = "Hello, I'm a language model,"
encoded_input = tokenizer(text, return_tensors='pt', padding=True)

# Run inference using temperature
with torch.inference_mode():
    model.reset_generation()
    generated_sequence = model.generate(
        input_ids=encoded_input.input_ids,
        attention_mask=encoded_input.attention_mask,
        do_sample=True,
        max_length=256,
        temperature=0.7,
    )

print([tokenizer.decode(tok) for tok in generated_sequence])
```

Note: As the Hugging Face generation API can expand the input's batch dimension based on different generation configurations, we need to compile the neuron model with different compile batch\_size compared to the run time batch\_size (batch dimension of inputs to generation API). - if do\_sample=True, compile\_batch\_size = runtime\_batch\_size x num\_return\_sequences x beam\_size - otherwise, compile\_batch\_size = runtime\_batch\_size x num\_return\_sequences

## Neuron Persistent Cache

The Neuron Persistent Cache is now enabled for Transformers Neuron by default. Model artifacts which have been compiled once will be cached and reused on successive runs when possible. Model artifacts will only be reused when compiling with the same compiler version (neuronx-cc), model configurations, and compiler flags. It also includes other features (i.e. using an S3 bucket as the cache backend). For more detailed information, see the [Persistent cache documentation](#)

## int8 weight storage support

Transformers Neuron supports int8 weight storage for the GPT2 model class. int8 weight storage can be used to reduce memory bandwidth usage to improve model performance. int8 weight storage support for additional model classes will be added in an upcoming release. In the following example we demonstrate how to apply int8 weight storage to the GPT2 model via the QuantizationConfig and NeuronConfig configs:

```
import torch
from transformers import AutoTokenizer
from transformers_neuronx import GPT2ForSamplingWithContextBroadcasting, NeuronConfig,
    QuantizationConfig

# Set the weight storage config use int8 quantization and bf16 dequantization
neuron_config = NeuronConfig(
    quant=QuantizationConfig(quant_dtype='s8', dequant_dtype='bf16'),
)

# Create and compile the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained(
    'gpt2',
    amp='bf16', # NOTE: When using quantization, amp type must match dequant type
    neuron_config=neuron_config
)
model.to_neuron()

# Get a tokenizer and example input
tokenizer = AutoTokenizer.from_pretrained('gpt2')
text = "Hello, I'm a language model,"
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model.sample(encoded_input.input_ids, sequence_length=256,
    start_ids=None)
    print([tokenizer.decode(tok) for tok in generated_sequence])
```

## Parallel Input Prompt Context Encoding

Transformers Neuron supports parallel input prompt context encoding for the GPT2 model class. Parallel context encoding can be used to significantly reduce the latency of the input prompt context encoding before the autoregressive decoder token generation loop. Parallel context encoding support for additional model classes will be added in an upcoming release.

The `GPT2ForSamplingWithContextBroadcasting` class has a `context_length_estimate` variable that determines the number of input prompt tokens that will be processed in parallel. For optimal results, this should be set to a power of 2 that is closest to the most frequently seen input prompt length. In the following example we demonstrate how to apply parallel context encoding to the GPT2 model via the `GPT2ForSamplingWithContextBroadcasting` class. In this example, we set the `context_length_estimate` to be 128, which is the closest power of 2 the length of the input prompt (97 tokens).

```
import torch
from transformers import AutoTokenizer
from transformers_neuronx import GPT2ForSamplingWithContextBroadcasting

# Create and compile the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained(
    'gpt2',
    context_length_estimate=256 # Create an optimized network which handles prompts up
    ↳ to 256 tokens
)
model.to_neuron()

# Get a tokenizer and example input
tokenizer = AutoTokenizer.from_pretrained('gpt2')
text = "Hello, I'm a generative AI language model. Generative AI is a type of AI that
↳ can create new content and ideas, including conversations, stories, images, videos,
↳ and music. It is powered by large models that are pre-trained on vast amounts of data
↳ and commonly referred to as foundation models (FMs). With generative AI on AWS, you
↳ can reinvent your applications, create entirely new customer experiences, drive
↳ unprecedented levels of productivity, and transform your business. "
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model.sample(encoded_input.input_ids, sequence_length=256)
print([tokenizer.decode(tok) for tok in generated_sequence])
```

The `GPT2ForSamplingWithContextBroadcasting` class can also process an input prompt that has a different batch size from the batch size of the autoregressive decoder output. For example, an input prompt with batch size = 1 can be used to produce an output of batch size = 5 to generate multiple suggestions for the same input prompt. The input prompt batch size can be specified using the `prompt_batch_size` argument and the autoregressive decoder output batch size can be specified using the `batch_size` argument. In the following example we demonstrate how to apply parallel context encoding to the GPT2 model to generate 5 outputs for a single input.

```
import torch
from transformers import AutoTokenizer
from transformers_neuronx import GPT2ForSamplingWithContextBroadcasting

# Create and compile the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained(
```

(continues on next page)

(continued from previous page)

```

    'gpt2',
    prompt_batch_size=1, # This allows prompt and output batch to vary
    batch_size=5,
    context_length_estimate=256
)
model.to_neuron()

# Get a tokenizer and example input
tokenizer = AutoTokenizer.from_pretrained('gpt2')
text = "Hello, I'm a generative AI language model. Generative AI is a type of AI that
↳ can create new content and ideas, including conversations, stories, images, videos,
↳ and music. It is powered by large models that are pre-trained on vast amounts of data,
↳ and commonly referred to as foundation models (FMs). With generative AI on AWS, you
↳ can reinvent your applications, create entirely new customer experiences, drive
↳ unprecedented levels of productivity, and transform your business. "
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model.sample(encoded_input.input_ids, sequence_length=256)

for i, output in enumerate(generated_sequence):
    print('-' * 50)
    print(f'Batch {i} output:')
    print(tokenizer.decode(output))

```

## Serialization support

Transformers NeuronX supports model serialization (model saving and loading) for all models except the GPTJForSampling and GPTNeoXForSampling` model classes. In the following example we demonstrate how to save and load the compiled artifacts for the GPT2 model:

```

import torch
from transformers import AutoTokenizer
from transformers_neuronx import GPT2ForSamplingWithContextBroadcasting

# Create and compile the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained('gpt2')
model.to_neuron()

# Save the compiled Neuron model
model.save('gpt2-compiled-artifacts')

# Load the Neuron model
model = GPT2ForSamplingWithContextBroadcasting.from_pretrained('gpt2')
# Load the compiled Neuron artifacts
model.load('gpt2-compiled-artifacts')
# Since prior artifacts are loaded, this skips compilation
model.to_neuron()

# Get a tokenizer and example input

```

(continues on next page)

(continued from previous page)

```

tokenizer = AutoTokenizer.from_pretrained('gpt2')
text = "Hello, I'm a language model,"
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model.sample(encoded_input.input_ids, sequence_length=256,
    ↪ start_ids=None)
print([tokenizer.decode(tok) for tok in generated_sequence])

```

Transformers NeuronX also supports the serialization of presharded weights. This reduces future model load time by saving a transformed and sharded set of weights as a new safetensors checkpoint. When this checkpoint is loaded, sharding and transformations normally done by Transformers NeuronX will be skipped, reducing model load time significantly. The saving of presharded weights is only available when `on_device_embedding` is true. In the following example we demonstrate how to save and load presharded weights along with compiled artifacts on a Llama model:

```

from transformers_neuronx import LlamaForSampling
from transformers_neuronx import NeuronConfig
from transformers import AutoTokenizer

neuron_config = NeuronConfig(on_device_embedding=True)

# Create and compile the Neuron model
model_neuron = LlamaForSampling.from_pretrained('openlm-research/open_llama_3b', batch_
    ↪ size=1, tp_degree=8, n_positions=128, neuron_config=neuron_config)
model_neuron.to_neuron()

# save the presharded weights and compiled artifacts to a directory
model_neuron.save('llama-artifacts', sharded_weights=True)

del model_neuron

# use the presharded checkpoint to reduce model load time
model_neuron_presharded = LlamaForSampling.from_pretrained('llama-artifacts', batch_
    ↪ size=1, tp_degree=8, n_positions=128, neuron_config=neuron_config)

# load in the compiled artifacts to skip compilation
model_neuron_presharded.load('llama-artifacts')
model_neuron_presharded.to_neuron()

```

## CPU Compilation Support

Transformers NeuronX now supports compilation on CPU. CPU compilation is compatible with model serialization and presharding weights, and is available for all models except the `GPTJForSampling` and `GPTNeoXForSampling` model classes. To compile on CPU, the initial call to `to_neuron()` is replaced with `cpu_compile()`. In the following example we demonstrate how to compile on CPU for the LLaMA model:

```

from transformers_neuronx import LlamaForSampling
from transformers_neuronx import NeuronConfig
from transformers import AutoTokenizer

```

(continues on next page)

(continued from previous page)

```
neuron_config = NeuronConfig(on_device_embedding=True)

# Create and compile the model on CPU
model_neuron = LlamaForSampling.from_pretrained('openlm-research/open_llama_3b', batch_
↪size=1, tp_degree=8, n_positions=128, neuron_config=neuron_config)
model_neuron.cpu_compile() # instead of model_neuron.to_neuron()

# save the weights and compiled artifacts to a directory
model_neuron.save('llama-artifacts')
```

To use the saved artifacts generated by CPU compilation on a Neuron device:

```
from transformers_neuronx import LlamaForSampling
from transformers_neuronx import NeuronConfig
from transformers import AutoTokenizer

neuron_config = NeuronConfig(on_device_embedding=True)

# use the presharded checkpoint to reduce model load time
model_neuron_presharded = LlamaForSampling.from_pretrained('llama-artifacts', batch_
↪size=1, tp_degree=8, n_positions=128, neuron_config=neuron_config)

# load in the compiled artifacts to skip compilation
model_neuron_presharded.load('llama-artifacts')

# now, use CPU compiled artifacts to run the model
model_neuron_presharded.to_neuron()
```

## Compilation worker count support

Transformers-neuronx supports providing compilation worker count for all models. This setting controls how many workers will execute HLO graph compilation tasks in parallel. A lower setting reduces CPU memory utilization when compiling a model, but increases the compilation time. This setting is useful to prevent out of CPU memory errors when compiling large models. By default, the number of workers used is equal to the total HLO graphs required for compilation. Compilation worker count integrates with both CPU compilation flow using `cpu_compile()` and neuron device compilation flow using `to_neuron()`. To set the compilation worker count, use the `compilation_worker_count` argument in `NeuronConfig`. The following sample shows how to compile the graphs one by one.

```
neuron_config = NeuronConfig(compilation_worker_count=1)
```

## Grouped-query attention (GQA) support [Beta]

Transformers Neuron supports grouped-query attention (GQA) models for Llama and Mistral model classes. There are multiple sharding strategies for K/V cache, in order to satisfy different constraints.

- `GQA.SHARD_OVER_HEADS` distributes K/V caches along head dimension. This can be only used when K/V heads is multiple of tensor-parallelism degree. This is the default configuration.
- `GQA.SHARD_OVER_BATCH` distributes K/V caches along batch dimension. This can be only used when batch size is multiple of tensor-parallelism degree. This can be useful for large-batch inference.



- `GQA.REPLICATED_HEADS` replicates K/V heads. This can be used when neither batch size nor K/V heads can be divisible by tensor-parallelism degree. This can be useful for low-latency small-batch inference.
- `GQA.ALL_GATHER_HEADS` evenly splits the K/V heads across all NeuronCores. This is optimized for large-batch inference of GQA model without replication.

In the following example we demonstrate how to configure these distributed inference strategies and perform inference with the Mistral model:

```
import torch
from transformers import AutoTokenizer
from transformers_neuronx import MistralForSampling, GQA, NeuronConfig

# Set sharding strategy for GQA to be shard over heads
neuron_config = NeuronConfig(
    group_query_attention=GQA.SHARD_OVER_HEADS
)

# Create and compile the Neuron model
model_neuron = MistralForSampling.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2',
    ↪ amp='bf16', neuron_config=neuron_config)
model_neuron.to_neuron()

# Get a tokenizer and exaple input
tokenizer = AutoTokenizer.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2')
text = "[INST] What is your favourite condiment? [/INST]"
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model_neuron.sample(encoded_input.input_ids, sequence_
    ↪ length=256, start_ids=None)
print([tokenizer.decode(tok) for tok in generated_sequence])
```

## Repeated Ngram Filtering

Repeated Ngram Filtering reduces redundant ngram phrases within the generated text. It uses the same API as Hugging-Face API for `NoRepeatedNGram`. Set the parameter `no_repeat_ngram_size` to the size of ngram phrases to be filtered and pass it to the sampling function as in the example `model.sample(inputs_ids, no_repeat_ngram_size=3)`

## On-device sampling support [Beta]

Transformers-neuronx supports on-device sampling for all models except Mixtral models. The features can be enabled by setting `on_device_generation` in `NeuronConfig` to an instance of `GenerationConfig`.

In the following example, we demonstrate how to use on-device generation for a Llama model using `top_k`, `top_p`, `top_p_min_tokens` and `temperature`.



## Top-K on-device sampling support [Beta]

Transformers Neuron supports Top-K Sampling on-device for all models except Mixtral models. In the following example, we demonstrate how to use on-device Top-K for the Llama model via the GenerationConfig and NeuronConfig configs.

```
import torch
from transformers_neuronx import LlamaForSampling
from transformers_neuronx.config import NeuronConfig, GenerationConfig
from transformers import AutoTokenizer

neuron_config = NeuronConfig(
    on_device_generation=GenerationConfig(max_length=128, top_k=10, top_p=0.9, top_p_min_
    ↪tokens=1, temperature=0.9, do_sample=True)
)

# Create and compile the Neuron model
model_neuron = LlamaForSampling.from_pretrained('openlm-research/open_llama_3b', batch_
    ↪size=1, tp_degree=8, n_positions=128, neuron_config=neuron_config)
model_neuron.to_neuron()

# Get a tokenizer and exaple input
tokenizer = AutoTokenizer.from_pretrained('openlm-research/open_llama_3b')
text = "Hello, I'm a language model,"
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model_neuron.sample(encoded_input.input_ids, sequence_
    ↪length=128, top_k=10)
    print([tokenizer.decode(tok) for tok in generated_sequence])
```

By default, transformers-neuronx uses the same, fixed sampling parameters for all sequences across all invocations of the model when on-device generation is enabled. It is possible to provide new sampling parameters per model invocation by enabling the dynamic feature in the GenerationConfig. It is also possible to provide different sampling parameters for each sequence in the batch by using the per\_batch\_line feature. When using this feature, it is recommended to limit the number of tokens that are considered during sampling across all sequences by setting global\_top\_k to a reasonably low number e.g. 250 to prevent poor performance when computing top\_p tokens over a large vocabulary without any prior filtering. When using per\_batch\_line, top\_k, top\_p, top\_p\_min\_tokens and temperature accept lists with value per sequence in the batch.

In the following example, we demonstrate how to use the dynamic and per\_batch\_line features together.

```
import torch
from transformers_neuronx import LlamaForSampling
from transformers_neuronx.config import NeuronConfig, GenerationConfig
from transformers import AutoTokenizer

batch_size = 2
generation_config = GenerationConfig(
    max_length=128, dynamic=True, per_batch_line=True, do_sample=True,
    top_k=[1] * batch_size,
    top_p=[1.0] * batch_size,
    top_p_min_tokens=[1] * batch_size,
```

(continues on next page)

(continued from previous page)

```

        temperature=[1.0] * batch_size,
        global_top_k=256
    )

neuron_config = NeuronConfig(
    on_device_generation=generation_config
)

# Create and compile the Neuron model
model_neuron = LlamaForSampling.from_pretrained('openlm-research/open_llama_3b', batch_
↪size=2, tp_degree=8, n_positions=128, neuron_config=neuron_config)
model_neuron.to_neuron()

# Get a tokenizer and exaple input
tokenizer = AutoTokenizer.from_pretrained('openlm-research/open_llama_3b')
tokenizer.pad_token = tokenizer.eos_token
text = ["Hello, I'm a language model,", "Hello, I'm also a language model,"]
encoded_input = tokenizer(text, return_tensors='pt')

# Run inference
with torch.inference_mode():
    generated_sequence = model_neuron.sample(encoded_input.input_ids, sequence_
↪length=128)
    print([tokenizer.decode(tok) for tok in generated_sequence])

    # Use different settings for each sequence in the batch
    # Supported because we use `generation_config.per_batch_line = True`
    generation_config.top_k = [1, 20]
    generation_config.top_p = [1.0, 0.9]
    generation_config.top_p_min_tokens = [1, 1]
    generation_config.temperature = [1.0, 0.9]

    # Update the generation configuration dynamically
    # Supported because we use `generation_config.dynamic = True`
    model_neuron.update_generation_config(generation_config)

    generated_sequence = model_neuron.sample(encoded_input.input_ids, sequence_
↪length=128)
    print([tokenizer.decode(tok) for tok in generated_sequence])

```

## Running inference with multiple models

Multiple transformers-neuronx models can be loaded at the same time as long as the total number of consumed NeuronCores is less than or equal to the total number of NeuronCores on the instance. For example, three tp-degree=8 models can be loaded and run in parallel on an inf2.48xlarge which has 24 NeuronCores. The NEURON\_RT\_NUM\_CORES and NEURON\_RT\_VISIBLE\_CORES environment variables can be used to allocate the necessary number of NeuronCores to each process to run multiple transformers-neuronx models in parallel. See the [NeuronCore Allocation and Model Placement for Inference \(torch-neuronx\)](#) section for additional information about how to use these environment variables.

It is important to notice that when multiple models are used on a single instance, the number of threads should be reduced to avoid race condition on host side. Assume the neuron instance (i.e. trn1) has 192 CPU cores. If one of

the models keeps all CPU cores busy, there would be significant performance degradation in the rest of models. As a result, the number of threads for each model should be limited to part of available cores. To do this, OMP\_NUM\_THREADS environment variable can be set. For example, if there are 192 CPU cores available and four tp-degree=8 models are used, one can export OMP\_NUM\_THREADS=48 to avoid race condition.

## Streamer

LLMs generate tokens in auto-regressive loop. A model.sample call waits till the end of full sequence generation before returning the generated response. It is possible to output an output token as soon as it is generated. To do this, a streamer object can be used. Streamer is an object which has 2 methods: put and end. There are several predefined streamer in transformers library such as TextIteratorStreamer. The following example shows how to define a streamer and use it in transformers-neuronx:

```
import torch
from transformers import AutoTokenizer
from transformers_neuronx import MistralForSampling, GQA

import transformers
from time import time

# Create a custom streamer inherited from transformers.generation.streamers.BaseStreamer
class CustomStreamer(transformers.generation.streamers.BaseStreamer):
    def __init__(self) -> None:
        self.reset()

    def reset(self):
        self.token_latencies = []
        self.iter = 0
        self.now = time()

    def put(self, tokens):
        now = time()
        token_latency = now - self.now
        print(f"Iteration {self.iter:4d}: Latency [s] {token_latency:6.3f} -- Token
↪{tokens}")
        self.now = now
        self.iter += 1
        self.token_latencies.append(token_latency)

    def end(self):
        print("First 10 token latencies:", self.token_latencies[:10])

# Create and compile the Neuron model
model_neuron = MistralForSampling.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2',
↪amp='bf16')
model_neuron.to_neuron()

# Get a tokenizer and exaple input
tokenizer = AutoTokenizer.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2')
text = "[INST] What is your favourite condiment? [/INST]"
```

(continues on next page)

(continued from previous page)

```

encoded_input = tokenizer(text, return_tensors='pt')

streamer = CustomStreamer()
# Run inference
with torch.inference_mode():
    generated_sequence = model_neuron.sample(encoded_input.input_ids, sequence_
    ↳length=256, start_ids=None, streamer=streamer)

```

## Stopping Criteria

We can define custom stopping criteria to stop autoregressive loop. For example, if we want to limit autoregressive loop after 0.5s, we can define and use stopping criteria class as follows:

```

import torch
import transformers
from transformers import AutoModelForCausalLM, AutoTokenizer, TextIteratorStreamer
from transformers_neuronx import MistralForSampling, GQA, NeuronConfig
from transformers_neuronx.stopping_criteria import StoppingCriteria, StoppingCriteriaList

from time import time
from typing import List, Optional, Callable

class MaxTimeCriteria(StoppingCriteria):
    """
    This class can be used to stop generation whenever the full generation exceeds some
    ↳amount of time. By default, the
    time will start being counted when you initialize this function. You can override
    ↳this by passing an
    `initial_time`.

    Args:
        max_time (float):
            The maximum allowed time in seconds for the generation.
        initial_time (float, *optional*, defaults to `time()`):
            The start of the generation allowed time.
    """

    def __init__(self, max_time: float, initial_timestamp: Optional[float] = None):
        self.max_time = max_time
        self.initial_timestamp = time() if initial_timestamp is None else initial_
        ↳timestamp

    def __call__(self, input_ids: torch.LongTensor, scores: torch.FloatTensor, **kwargs)
    ↳-> bool:
        dt = time() - self.initial_timestamp
        end_condition = dt > self.max_time
        if end_condition:
            print("Stopping!")
        return end_condition

```

(continues on next page)

(continued from previous page)

```

# Create a streamer. This can be a custom streamer too inherited from transformers.
↳ generation.streamers.BaseStreamer
class CustomStreamer(transformers.generation.streamers.BaseStreamer):
    def __init__(self) -> None:
        self.reset()

    def reset(self):
        self.token_latencies = []
        self.iter = 0
        self.now = time()

    def put(self, tokens):
        now = time()
        token_latency = now - self.now
        print(f"Iteration {self.iter:4d}: Latency [s] {token_latency:6.3f} -- Token
↳ {tokens}")
        self.now = now
        self.iter += 1
        self.token_latencies.append(token_latency)

    def end(self):
        pass

# Create and compile the Neuron model
model_neuron = MistralForSampling.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2',
↳ amp='bf16')
model_neuron.to_neuron()

# Get a tokenizer and exaple input
tokenizer = AutoTokenizer.from_pretrained('mistralai/Mistral-7B-Instruct-v0.2')
text = "[INST] What is your favourite condiment? [/INST]"
encoded_input = tokenizer(text, return_tensors='pt')

# Add stopping criteria to stop after 0.5 seconds
stopping_criteria_list= StoppingCriteriaList([MaxTimeCriteria(0.5)])
streamer = CustomStreamer()

# Run inference
with torch.inference_mode():
    model_neuron.sample(input_ids=encoded_input.input_ids, sequence_length=256, stopping_
↳ criteria_list=stopping_criteria_list, streamer=streamer)

```

## Speculative sampling [Beta]

Transformers Neuron supports speculative sampling for the Llama and GPT2 model classes. In speculative sampling, we use a smaller draft model to speculate future tokens. These are then sent to the larger target model, which accepts or rejects these tokens. For more detailed information, see the original proposal by DeepMind titled [Accelerating Large Language Model Decoding with Speculative Sampling](#). Our implementation for speculative sampling is lossless. In addition to standalone draft models, we also support [Eagle draft models](#). Currently we only support Eagle v1.

In the following example, we demonstrate how to perform speculative sampling using the Llama model. In this example, we are performing multinomial sampling.

```
import torch
from transformers import LlamaTokenizer
from transformers_neuronx import NeuronAutoModelForCausalLM, NeuronConfig, \
    GenerationConfig
from transformers_neuronx.fused_speculation import FusedSpeculativeDecoder

# Specify path to draft and target
draft = '/home/ubuntu/Llama-2-7b-chat-hf'
target = '/home/ubuntu/Llama-2-70b-chat-hf'

# Specify generation parameters
gen_kwargs = {
    "top_k": 50,
    "top_p": 0.9,
    "do_sample": True,
    "temperature": 0.7,
}

# Load draft model
draft_neuron_model = NeuronAutoModelForCausalLM.from_pretrained(
    draft,
    n_positions=1024,
    batch_size=1,
    tp_degree=32,
    amp='bf16',
    neuron_config=NeuronConfig(
        padding_side="right",
        attention_layout=Layout.BSH,
        collectives_layout="BSH",
        on_device_embedding=True,
        on_device_generation=GenerationConfig(**gen_kwargs),
    ),
)
draft_neuron_model.to_neuron()

# Load target model
target_neuron_model = NeuronAutoModelForCausalLM.from_pretrained(
    target,
    n_positions=1024,
    batch_size=1,
    tp_degree=32,
    amp='bf16',
    neuron_config=NeuronConfig(
```

(continues on next page)

(continued from previous page)

```

        padding_side="right",
        attention_layout=Layout.BSH,
        collectives_layout="BSH",
        on_device_embedding=True,
        on_device_generation=GenerationConfig(**gen_kwargs),
    ),
)
target_neuron_model.to_neuron()

# Compile the speculative sampling model
# Here we set sepculation length to be 4
fsd = FusedSpeculativeDecoder(
    draft_neuron_model,
    target_neuron_model,
    4,
)
fsd.to_neuron()

# Initialize tokenizer and text prompt
tokenizer = LlamaTokenizer.from_pretrained(target)
prompt = "Hello, I'm a generative AI language model."
inputs = tokenizer(prompt, return_tensors="pt")

# Call speculative sampling on given input
response = fsd.sample(
    input_ids=inputs.input_ids,
    attention_mask=inputs.attention_mask,
    sequence_length=30,
)

# Decode the response
generated_text = tokenizer.decode(response[0])
print(f"\nDecoded tokens: {generated_text}")

```

The following sample shows how to enable EAGLE speculation. To get the EAGLE draft model to work, manually copy the LM head weights from the target model to the draft model. Additionally, you need to rename the keys in the draft model's state\_dict to match those in the target model.

```

import torch
from transformers import LlamaTokenizer
from transformers_neuronx import NeuronAutoModelForCausalLM, NeuronConfig, \
    GenerationConfig
from transformers_neuronx.fused_speculation import FusedSpeculativeDecoder

# Specify path to draft and target
# The Eagle draft model can be downloaded from Eagle website
draft = '/home/ubuntu/EAGLE-llama2-chat-70B'
target = '/home/ubuntu/Llama-2-70b-chat-hf'

# Specify generation parameters
gen_kwargs = {
    "top_k": 50,

```

(continues on next page)

(continued from previous page)

```

    "top_p": 0.9,
    "do_sample": True,
    "temperature": 0.7,
}

# Load draft model
draft_neuron_model = NeuronAutoModelForCausalLM.from_pretrained(
    draft,
    n_positions=1024,
    batch_size=1,
    tp_degree=32,
    amp='bf16',
    neuron_config=NeuronConfig(
        is_eagle_draft=True,
        has_pre_attention_norm=False,
        # Need the above two configs for Eagle
        padding_side="right",
        attention_layout=Layout.BSH,
        collectives_layout="BSH",
        on_device_embedding=True,
        on_device_generation=GenerationConfig(**gen_kwargs),
    ),
)
draft_neuron_model.to_neuron()

# Load target model
target_neuron_model = NeuronAutoModelForCausalLM.from_pretrained(
    target,
    n_positions=1024,
    batch_size=1,
    tp_degree=32,
    amp='bf16',
    neuron_config=NeuronConfig(
        is_eagle_target=True,
        # Need the above config for Eagle
        padding_side="right",
        attention_layout=Layout.BSH,
        collectives_layout="BSH",
        on_device_embedding=True,
        on_device_generation=GenerationConfig(**gen_kwargs),
    ),
)
target_neuron_model.to_neuron()

# Compile the speculative sampling model
# Here we set sepculation length to be 4
fsd = FusedSpeculativeDecoder(
    draft_neuron_model,
    target_neuron_model,
    4,
)
fsd.to_neuron()

```

(continues on next page)



(continued from previous page)

```
# The rest are the same
```

## QKV Weight Fusion

Concatenating a model's query, key and value weight matrices often achieves better performance because larger matrices allow for more efficient data movement and compute. QKV weight fusion can be enabled by setting `fuse_qkv=True` in the `NeuronConfig`:

```
neuron_config = NeuronConfig(fuse_qkv=True)
```

## Attention Layout

The intermediate tensor layouts in a model's attention layer can impact the compiler's optimization opportunities and thus can impact a model's performance. Using (batch, sequence, hidden) (or BSH) layout for attention often achieves better performance since it can enable better overlapping of compute with collectives and can reduce transposes. We intend to enable BSH attention by default in a future release. For now, BSH attention layout can be enabled by setting `attention_layout="BSH"` in the `NeuronConfig`:

```
neuron_config = NeuronConfig(attention_layout="BSH")
```

## Bucketing

LLM inference is a generate process that can produce variable length sequences. This poses a problem since the Neuron compiler produces executables which expect statically shaped inputs and outputs. To make LLM work with different shapes, `transformers_neuronx` generates buckets and applies padding wherever it is required.

There are at least two set of buckets for each LLM inference that can be set by user: 1) Context encoding (pre-fill) buckets and 2) output token generation buckets.

### Token generation buckets

In token generation, tokens are generated iteratively. At each token position, transformer need to attend to the previous tokens only. But in the naive implementation with static shapes, one may attend to all KV-cache (full sequence length). To solve this problem, we use token generation buckets. Token generation buckets determine the attention lengths. For instance, if the max sequence length is 1024 tokens and current token is at position 120, there is no need to attend to all 1024 tokens in the current step. We can use token generation buckets to attend to different portions of KV-cache. By default, token generation buckets which are powers of 2 starting from 128 tokens are used (i.e. 128, 256, 512, up to sequence length). In the example above, bucket 128 would be used for position 120 which would reduce the wasted compute significantly. User can change these buckets by setting a list for `n_positions` (see example below). Otherwise, if a number is given for `n_positions` (sequence length), instead of a list, then the powers of 2 buckets starting from 128 will be used. The last bucket would be `n_positions` (sequence length), even if it is not a power of 2.

### Context encoding buckets

The prompt tokens can be processed in parallel. As a result, we need to set the bucket sizes for different estimated length of input prompts. We can specify these context bucket sizes using the `context_length_estimate` argument. In general, it is better to have all the bucket to be multiples of 256 tokens. But adding too many buckets would increase device memory consumption and add extra latency for bucket switching. Usually, the powers of 2 starting from 128 tokens are used for context encoding buckets. If the total sequence length (`n_positions`) is beyond 2048 tokens, it is desirable to add extra buckets with multiple of 512 or 1024 tokens. It is not recommended to add buckets of multiples

of 256 tokens or smaller for context buckets beyond 2k to avoid bucket switching latency. At runtime, the smallest bucket which fits the input context will be used. By default, the context encoding buckets set to half of output-token buckets. Adding extra context buckets would reduce the wasted compute and improves performance. However, the extra executables would reduce memory space since executables require device memory space.

Notice that the default output token generation buckets work well for wide range of applications. However, ideal context encoding buckets depends on the specific use case. For instance, if all the requests have a context length of about 1500 +/- 500 tokens, adding more buckets closer to 1500 might help context encoding time. In this example, adding buckets of 1024, 1280, 1536, 1792, 2048 tokens (distance of 256 tokens) could help. Moreover, the largest context encoding bucket should be larger than the largest context length. Otherwise, the performance would degrade significantly.

To set context encoding and token generation buckets manually:

```
context_length_estimate = [1024, 1280, 1536, 1792, 2048]    # The best context estimate,
↳ depends on the use case
n_positions = [128, 256, 512, 1024, 2048, 3072]          # Usually default buckets,
↳ are appropriate

model = NeuronAutoModelForCausalLM.from_pretrained(
    'gpt2',
    batch_size=1,
    n_positions=n_positions,
    tp_degree=2,
    amp='f16',
    context_length_estimate=context_length_estimate,
)
```

## Multi-node inference support (TP/PP)

Prerequisite: <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/frameworks/torch/torch-neuronx/setup-trn1-multi-node-execution.html>

When models are too large to fit on single node, Transformers NeuronX multi-node inference (tensor parallel and pipeline parallel) can be used to shard model weights across multiple Neuron instances (only supported on Trn1 and Trn1n). Single node inference code can easily be extended to multi-node inference.

Note that Transformers Neuronx currently doesn't support multi-node Tensor Parallel and Pipeline Parallel at same time, when Pipeline Parallel is used, the Tensor Parallel has to be within a node (TP<=32 on Trn1/Trn1n).

In the below sections, we first outline the sample code for single node execution and then provide instructions to migrate the code to use multi-node tensor parallel or multi-node pipeline parallel. To start with, the code below is for single node script, running llama2-3b model with tensor parallel degree as 32.

```
import torch
from transformers import AutoTokenizer, AutoConfig
from transformers_neuronx import LlamaForSampling, HuggingFaceGenerationModelAdapter

# Create and compile the Neuron model
model = LlamaForSampling.from_pretrained("openlm-research/open_llama_3b", tp_degree=32)
model.to_neuron()

# Use the `HuggingFaceGenerationModelAdapter` to access the generate API
config = AutoConfig.from_pretrained("openlm-research/open_llama_3b")
model = HuggingFaceGenerationModelAdapter(config, model)
```

(continues on next page)

(continued from previous page)

```
# Get a tokenizer and example input
tokenizer = AutoTokenizer.from_pretrained("openlm-research/open_llama_3b")
tokenizer.pad_token_id = tokenizer.eos_token_id
tokenizer.padding_side = 'left'
text = "Hello, I'm a language model,"
encoded_input = tokenizer(text, return_tensors='pt', padding=True)

# Run inference using temperature
with torch.inference_mode():
    model.reset_generation()
    generated_sequence = model.generate(
        input_ids=encoded_input.input_ids,
        attention_mask=encoded_input.attention_mask,
        do_sample=True,
        max_length=256,
        temperature=0.7,
    )

print([tokenizer.decode(tok) for tok in generated_sequence])
```

command line:

```
python3 multi_node_dev_example.py
```

### Multi-Node Tensor Parallel

Compared to single node tensor parallel, multi-node tensor parallel shards the model weights in the same way but having more cores across nodes. In the meantime, it requires each node's `model.forward()` receives the exact same input, otherwise there would be unexpected behaviors (runtime failure, wrong output).

Configurations (environment variables to be configured on each node):

- `NEURON_RT_ROOT_COMM_ID`: the master node's <IP address>:<port>
- `NEURON_RANK_ID`: rank of the node, 0 means master node
- `NEURON_LOCAL_TP`: the local tensor parallel degree on each node

example:

Change the single node script to use `tp=64` (2 node). Set the `torch.manual_seed` to ensure the sampling loop running on each node will sample same token as next input.

Node 1 command line:

```
NEURON_RT_ROOT_COMM_ID=10.1.201.64:63423 NEURON_RANK_ID=0 NEURON_LOCAL_TP=32 python3 \
multi_node_dev_example.py
```

Node 2 command line (same as Node 1 but set `NEURON_RANK_ID` as 1):

```
NEURON_RT_ROOT_COMM_ID=10.1.201.64:63423 NEURON_RANK_ID=1 NEURON_LOCAL_TP=32 python3 \
multi_node_dev_example.py
```

You can also refer to [Tutorial](#) to run llama 3.1 405b multinode 16k tutorial with multi-node tensor parallel.

### Multi-Node Pipeline Parallel

While having the weight tensor sharded as tensor parallel, one can utilize pipeline parallel to partition the layers across different node, the intermediate tensor (hidden) will be transferred from one pipeline stage (nodes) to the next pipeline stage (nodes). The final output will be sent from last pipeline stage back to first pipeline stage.

Compared to multi-node tensor parallel, for non-zero rank, the `model.forward` in pipeline parallel will fallback to while loop and block on the input broadcasting from master.

Configurations (environment variables to be configured on each node):

- `NEURON_RT_ROOT_COMM_ID`: the master node's <IP address>:<port>
- `CPU_COMM_ID`: similar to `NEURON_RT_ROOT_COMM_ID`, but need to set with different port
- `NEURON_RANK_ID`: rank of the node, 0 means master node
- `NEURON_PP_STAGES`: number of pipeline stages (nodes)

example:

Keep the original single node script with `tp=32`.

Node 1 command line:

```
NEURON_PP_STAGES=2 CPU_COMM_ID=10.1.201.64:8989 NEURON_RT_ROOT_COMM_ID=10.1.201.64:63423_
↪NEURON_RANK_ID=0 python3 multi_node_dev_example.py
```

Node 2 command line (same as Node 1 but set `NEURON_RANK_ID` as 1):

```
NEURON_PP_STAGES=2 CPU_COMM_ID=10.1.201.64:8989 NEURON_RT_ROOT_COMM_ID=10.1.201.64:63423_
↪NEURON_RANK_ID=1 python3 multi_node_dev_example.py
```

## Long Sequence length support up to 128k

### Flash Attention

With the integration of FlashAttention kernel, developers can use longer sequence lengths for LLAMA models. The Flash Attention kernel is automatically used when the input sequence length is greater than 8k without any additional configuration. Refer to [Tutorial](#) for usage of 32k sequence length on a variation of LLAMA3-8B Model.

### Flash Decoding

Flash Decoding (FD) is a technique that significantly speeds up attention during inference, especially for long-context tasks in large language models (LLMs) with GQA.

With integration of FD, developers can achieve faster inference with larger sequence and batch size by reducing the KV cache replication. Refer to [Tutorial](#) on flash decoding usage for 128k sequence length sampling. Flash decoding can be enabled by setting the flag `shard_over_sequence=True` in *NeuronConfig*

```
neuron_config = NeuronConfig(shard_over_sequence=True)
```

Note that you can skip the first Allgather introduced by flash decoding at the cost of duplicate Q weights, this is only recommended for relatively small models (i.e. 3B, 8B) and large batch size.

```
neuron_config = NeuronConfig(shard_over_sequence=True, duplicate_q_weight_sos=True)
```

### Known limitations and FAQs

- Flash decoding is expected to have performance degradation (PTL) for smaller sequence and batch sizes. We recommend flash decoding when **batch-size x sequence length > 16k**

- Flash decoding support is not enabled for the following features
- Speculative Decoding
- Multi Head Attention (MHA) models

*This document is relevant for: Inf2, Trn1*

*This document is relevant for: Inf2, Trn1*

## Transformers NeuronX (transformers-neuronx) Developer Guide for Continuous Batching

Transformers NeuronX is integrated with vLLM to enable continuous batching for high-throughput LLM serving and inference. This guide aims to help users get started with continuous batching for Transformers NeuronX and vLLM by providing:

- *Transformers NeuronX* An overview of Transformers NeuronX.
- *Continuous Batching with Transformers NeuronX and vLLM* The continuous batching procedure implemented by Transformers NeuronX and vLLM.
- *Install vLLM and Get Started with Offline Inference* Installation and usage instructions for Transformers NeuronX and vLLM.
- *New Features in Neuron Release 2.21* A showcase of new features in Transformers NeuronX and vLLM.
- *Frequently Asked Questions*

## Transformers NeuronX (transformers-neuronx)

Transformers NeuronX for Trn1 and Inf2 is a software package that enables PyTorch users to perform large language model (LLM) *performant inference* on second-generation Neuron hardware (See: NeuronCore-v2). The Neuron performance page lists expected inference performance for commonly used Large Language Models.

## Continuous Batching with Transformers NeuronX and vLLM

Transformers NeuronX implements the following operational flow with vLLM for continuous batching support:

1. Context encode multiple prompts using virtual dynamic batching.
2. Decode all sequences simultaneously until a sequence generates an EOS token.
3. Evict the finished sequence and insert a new prompt encoding.
4. Resume the decoding process, repeating steps 2 and 3 until all sequences are decoded.

## Supported Model Architectures

Transformers NeuronX supports continuous batching for models compatible with the following Hugging Face classes:

- LlamaForCausalLM
- MistralForCausalLM

## Install vLLM and Get Started with Offline Inference

Neuron maintains a fork of vLLM (v0.6.2) that contains the necessary changes to support inference with Transformers NeuronX. Neuron is working with the vLLM community to upstream these changes to make them available in a future version.

### Install vLLM

First install `neuronx-cc` and the `transformers-neuronx` packages. Then install the vLLM fork from source:

```
git clone -b v0.6.x-neuron https://github.com/aws-neuron/upstreaming-to-vllm.git
cd upstreaming-to-vllm
pip install -r requirements-neuron.txt
VLLM_TARGET_DEVICE="neuron" && pip install -e .
```

---

**Note:** Please note the vLLM pip package from PyPI is not compatible with Neuron. To work with Neuron, install vLLM using the source as outlined above.

---

---

**Note:** The current supported version of Pytorch for Neuron installs `triton` version `2.1.0`. This is incompatible with `vllm >= 0.5.3`. You may see an error `cannot import name 'default_dump_dir'...` To work around this, run `pip install --upgrade triton==3.0.0` after installing the vLLM wheel.

---

If Neuron packages are detected correctly in the installation process, `vllm-0.1.dev2830+g22c56ee.neuron216` will be installed (The neuron version depends on the installed `neuronx-cc` version).

## Run Offline Batched Inference with Transformers NeuronX and vLLM

In the following example we demonstrate how to perform continuous batching with a Llama model.

---

**Note:** Since Llama models are gated, please accept the Llama Community License Agreement and request access to the model. Then use a Hugging Face user access token to download the model.

---

```
from vllm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]

# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)

# Create an LLM.
llm = LLM(
    model="meta-llama/Meta-Llama-3.1-8B-Instruct",
```

(continues on next page)

(continued from previous page)

```

max_num_seqs=8,
    # The max_model_len and block_size arguments are required to be same as max sequence_
    ↪length,
    # when targeting neuron device. Currently, this is a known limitation in continuous_
    ↪batching
    # support in transformers-neuronx.
    max_model_len=128,
    block_size=128,
    # The device can be automatically detected when AWS Neuron SDK is installed.
    # The device argument can be either unspecified for automated detection, or_
    ↪explicitly assigned.
    device="neuron",
    tensor_parallel_size=2)

# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

```

## Run the API Server

To run the OpenAI-compatible API server in vLLM, run either command below:

```

vllm serve meta-llama/Meta-Llama-3.1-8B-Instruct --tensor-parallel-size 32 --max-num-
    ↪seqs 4 --max-model-len 2048 --block-size 8

```

```

python3 -m vllm.entrypoints.openai.api_server --model meta-llama/Meta-Llama-3.1-8B-
    ↪Instruct --tensor-parallel-size 32 --max-num-seqs 4 --max-model-len 2048 --block-size 8

```

## New Features in Neuron Release 2.21

Neuron's vLLM integration with Transformers NeuronX is tested using a public fork of vLLM v0.6.2. New features and enhancements introduced in this fork will be described below. Neuron's intent is to upstream these features to vLLM as soon as possible after release. Prior to upstreaming, these features can be accessed in the AWS Neuron GitHub repository <https://github.com/aws-neuron/upstreaming-to-vllm/tree/v0.6.x-neuron>.

### Neuron Release 2.21 Features for the v0.6.2 vLLM Neuron Fork

- *Sequence bucketing* configuration for context encoding and token generation.
- *Granular NeuronConfig control* in vLLM entrypoints.
- Inference support for *speculative decoding*.
- Inference support for *EAGLE speculative decoding*.

### Neuron Release 2.20 Features

- Multi-node inference support for larger models. Example scripts are included in vLLM .

- Direct loading of Hugging Face-compatible checkpoints without creation of a `-split` directory.

## Sequence Bucketing

To configure buckets, set the following environment variables. Refer to the [developer guide](#) for details on how to configure the values. These environment variables need to be set before starting the vLLM server or instantiating the LLM object.

- `NEURON_CONTEXT_LENGTH_BUCKETS`: Bucket sizes for context encoding.
- `NEURON_TOKEN_GEN_BUCKETS`: Bucket sizes for token generation.

For example: `export NEURON_CONTEXT_LENGTH_BUCKETS="128,512,1024"`

## NeuronConfig Override

The default `NeuronConfig` in vLLM uses the latest optimizations from the Neuron SDK. However, you can override the default values or add a new configuration from the [developer guide](#) by setting the `override_neuron_config` parameter while creating the LLM object.

```
llm = LLM(  
    model="meta-llama/Meta-Llama-3.1-8B-Instruct",  
    max_num_seqs=8,  
    max_model_len=128,  
    block_size=128  
    device="neuron",  
    tensor_parallel_size=32,  
    #Override or update the NeuronConfig  
    override_neuron_config={"shard_over_sequence":True})
```

While standing up the API server, set the `override-neuron-config` argument. For example:

```
python3 -m vllm.entrypoints.openai.api_server --model meta-llama/Meta-Llama-3.1-8B-  
→ Instruct --tensor-parallel-size 32 --max-num-seqs 4 --max-model-len 2048 --block-size_  
→ 8 --override-neuron-config {"shard_over_sequence":"True"}
```

## Quantization

To use [int8 weight storage](#), set the environment variable `NEURON_QUANT_DTYPE` to `s8`.

## Speculative Decoding

Speculative decoding is a token generation optimization technique that uses a small draft model to generate `K` tokens autoregressively and a larger target model to determine which draft tokens to accept, all in a combined forward pass. For more information on speculative decoding, please see [\[Leviathan, 2023\]](#) and [\[Chen et al., 2023\]](#).

Speculative decoding is now available for inference with Transformers NeuronX and vLLM:

```
from vllm import LLM, SamplingParams  
  
# Sample prompts.
```

(continues on next page)



(continued from previous page)

```

prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)

# Create an LLM.
llm = LLM(
    model="meta-llama/Meta-Llama-3.1-70B-Instruct",
    speculative_model="meta-llama/Llama-3.2-1B-Instruct",
    # The max_model_len, speculative_max_model_len, and block_size arguments are
    ↪required to be same as max sequence length,
    # when targeting neuron device. Currently, this is a known limitation in continuous
    ↪batching
    # support in transformers-neuronx.
    max_model_len=128,
    block_size=128,
    speculative_max_model_len=128,
    dtype="bfloat16",
    max_num_seqs=4,
    num_speculative_tokens=4,
    # The device can be automatically detected when AWS Neuron SDK is installed.
    # The device argument can be either unspecified for automated detection, or
    ↪explicitly assigned.
    device="neuron",
    tensor_parallel_size=32,
    use_v2_block_manager=True,
)

outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

```

**Note:** Please ensure that the selected target and draft model are from the same model family. For example, if the target model is an instruction-tuned Llama model, the draft model must also be a lower-capacity instruction-tuned Llama model.

## EAGLE Speculative Decoding

Extrapolation Algorithm for Greater Language-model Efficiency (EAGLE) extends the speculative decoding technique described above by:

- Utilizing a specially trained EAGLE draft model that predicts feature outputs through an Autoregression Head and next token outputs through an LM Head.
- Reducing sampling uncertainty by using the next autoregressively sampled token and a current feature map as draft model inputs.

For more information on EAGLE, please see [\[Li et al., 2024\]](#)

EAGLE speculative decoding can be applied without changes to the speculative decoding code sample above. Transformers NeuronX and vLLM will recognize a draft model as an EAGLE draft when `is_eagle: True` is set in the model's Hugging Face `config.json` file.

## Frequently Asked Questions

### Is PagedAttention supported in the vLLM integration?

No, PagedAttention is not currently supported. It will be supported in a future Neuron release.

*This document is relevant for:* Inf2, Trn1

- [Transformers NeuronX \(transformers-neuronx\) Developer Guide](#)

*This document is relevant for:* Inf2, Trn1

*This document is relevant for:* Inf2, Trn1

## 4.2.3 Transformers NeuronX Tutorials

- [Hugging Face meta-llama/Llama-2-13b autoregressive sampling on Inf2 & Trn1](#)
- [Hugging Face facebook/opt-13b autoregressive sampling on Inf2 & Trn1](#)
- [Hugging Face facebook/opt-30b autoregressive sampling on Inf2 & Trn1](#)
- [Hugging Face facebook/opt-66b autoregressive sampling on Inf2](#)

*This document is relevant for:* Inf2, Trn1

*This document is relevant for:* Inf2, Trn1

## 4.2.4 Misc (transformers-neuronx)

*This document is relevant for:* Inf2, Trn1, Trn2

## Transformers Neuron (transformers-neuronx) release notes

### Table of Contents

- *Model classes status*
- *Model features*
- *Release [0.13.380.0]*
- *Release [0.13.322.0]*
- *Release [0.12.313]*
- *Release [0.11.351.0]*
- *Release [0.10.0.332]*
- *Release [0.10.0.21]*
- *Release [0.9.474]*
- *Release [0.8.268]*
- *Release [0.7.84]*
- *Release [0.6.106]*
- *Release [0.5.58]*
- *Release [0.4.0]*
- *Release [0.3.0]*

Transformers Neuron for Trn1/Inf2 is a software package that enables PyTorch users to perform large language model (LLM) inference on second-generation Neuron hardware (See: NeuronCore-v2).

### Model classes status

- **BLOOM**: [Beta]
- **GPT2**: [Beta]
- **GPT-J**: [Beta]
- **GPT-Neox**: [Beta]
- **Llama**: [Beta]
- **Llama 2**: [Beta]
- **Mistral**: [Beta]

## Model features

| Model     | Flexible Tensor Parallelism | Prompt Estimate Support | Serialization Support |
|-----------|-----------------------------|-------------------------|-----------------------|
| BLOOM     | Yes                         | Yes                     | Yes                   |
| GPT2      | Yes                         | Partial                 | Yes                   |
| GPT-J     | No                          | No                      | No                    |
| GPT-NeoX  | No                          | No                      | No                    |
| Llama     | Yes                         | Yes                     | Yes                   |
| Llama 2   | Yes                         | Yes                     | Yes                   |
| Llama 3.1 | Yes                         | Yes                     | Yes                   |
| Mistral   | Yes                         | Yes                     | Yes                   |

## Release [0.13.380.0]

Date: 01/14/2025

### What's new in this release

- The transformers dependency has been pinned to `transformers<4.48`

## Release [0.13.322.0]

Date: 12/20/2024

## Summary

### What's new in this release

- Flash decoding support for speculative decoding
- Enabled on-device generation support in speculative decoding flows
- Added support for EAGLE speculative decoding support with greedy and lossless sampling
- Support for CPU compilation and sharded model saving

## Performance Improvements

- Performance optimized MLP and QKV kernels added for llama models with support for sequence parallel norm
- Added support to control concurrent compilation workers
- Added option to skip AllGather using duplicate Q weights during shard over sequence

## Resolved Issues

- Fixed padding issues when requested batch size is smaller than neff compiled size
- Fixed sequence parallel norm issue when executor is used with speculative decoding flows

## Known Issues and Limitations

- GPT-NeoX is sensitive to `fp16` and customers are advised to use only `amp="f32"` for GPT-NeoX.
- Using `cache_layout=constants.LAYOUT_BSH` in `NeuronConfig` has known limitations with compilation. Customers are advised to use `constants.LAYOUT_SBH` instead.

## Release [0.12.313]

Date: 09/16/2024

## Summary

### What's new in this release

- Support for model serialization (save and load) of all models except the `GPTJForSampling` and `GPTNeoXForSampling` model classes, which reduces future model load time by saving a transformed and sharded set of weights as a new safetensors checkpoint.
- Support for on device sampling (Top P) with Continuous batching
- Support for Scaled RoPE for LLAMA 3.1 models
- Support for multi-node inference for LLAMA 3.1 405B model for specific sequence lengths
- Support for FlashDecoding (using `shard_over_sequence`) for supporting long context lengths upto 128k [Tutorial](#)

## Resolved Issues

- Fixes to handle `seq_ids` consistently across vLLM versions
- Fixes for KV head full replication logic errors

## Known Issues and Limitations

- GPT-NeoX is sensitive to `fp16` and customers are advised to use only `amp="f32"` for GPT-NeoX.
- Using `cache_layout=constants.LAYOUT_BSH` in `NeuronConfig` has known limitations with compilation. Customers are advised to use `constants.LAYOUT_SBH` instead.

**Release [0.11.351.0]**

Date: 07/03/2024

**Summary****What's new in this release**

- Support for compiler optimized flash attention kernel to support context lengths of 16k/32k for Llama models
- Streamer support enabled for BLOOM, GPTJ, GPT2, GPT-NeoX and LLAMA models
- Support for on device generation for TopK in Mixtral models
- Continuous batching support for Mistral v0.2
- Minor API improvements with type annotations for NeuronConfig, deprecation warnings for old arguments, and exposing top-level configurations
- Performance improvements such as an optimized logit ordering for continuous batching in Llama models, optimized QKV padding for certain GQA models, faster implementation of cumsum operation to improve TopP performance

**Resolved Issues**

- Removed `start_ids=None` from `generate()`
- Mistral decoding issue that occurs during multiple sampling runs
- Mistralv0.1 sliding window error
- Off-by-one error in window context encoding
- Better error messaging

**Known Issues and Limitations**

- `on_device_generation=GenerationConfig(do_sample=True)` has some known failures for Llama models. Customers are advised not to use `on_device_generation` in such cases.
- GPT-NeoX is sensitive to `fp16` and customers are advised to use only `amp="f32"` for GPT-NeoX.
- Using `cache_layout=constants.LAYOUT_BSH` in `NeuronConfig` has known limitations with compilation. Customers are advised to use `constants.LAYOUT_SBH` instead.

**Release [0.10.0.332]**

Date: 04/10/2024

## Summary

### What's new in this release

- [Beta] Added support for continuous batching and a reference integration with vLLM (Llama models only)

### Known Issues and Limitations

- There is a known compiler issue for inference of some configurations of Llama-2 70B that can cause accuracy degradation. Customers are advised to use the `--enable-mixed-precision-accumulation` compiler flag if Llama-2 70B accuracy issues occur.
- There is a known compiler issue for inference of some configurations of Llama-2 13B that can cause accuracy degradation. Customers are advised to use the `--enable-saturate-infinity` `--enable-mixed-precision-accumulation` compiler flags if Llama-2 13B accuracy issues occur.
- There is a known compiler issue for inference of some configurations of GPT-2 that can cause accuracy degradation. Customers are advised to use the `--enable-saturate-infinity` `--enable-mixed-precision-accumulation` compiler flags if GPT-2 accuracy issues occur.
- GPT-NeoX is sensitive to `fp16` and customers are advised to use only `amp="f32"` for GPT-NeoX.
- Using `cache_layout=constants.LAYOUT_BSH` in `NeuronConfig` has known limitations with compilation. Customers are advised to use `constants.LAYOUT_SBH` instead.

### Release [0.10.0.21]

Date: 04/01/2024

## Summary

### What's new in this release

- Added support for on device log-softmax and on device sampling for TopK
- Added support for on device embedding for all models.
- Added support for Speculative Decoding
- [Beta] Added support for Mixtral-8x7b MoE
- [Beta] Added support for mistralai/Mistral-7B-Instruct-v0.2 with no sliding window
- Added faster checkpoint loading support for both sharded and whole checkpoints
- Added the ability to download checkpoints directly from huggingface hub repositories
- Added `NeuronAutoModelForCausalLM` class which automatically loads architecture-specific classes
- Added a warmup to all kernels to avoid unexpected initialization latency spikes

## Resolved Issues

- Users no longer need a copy of the original checkpoint and can use safetensor checkpoints for optimal speed.

## Known Issues and Limitations

- There is a known compiler issue for inference of some configurations of Llama-2 70B that can cause accuracy degradation. Customers are advised to use the `--enable-mixed-precision-accumulation` compiler flag if Llama-2 70B accuracy issues occur.
- There is a known compiler issue for inference of some configurations of Llama-2 13B that can cause accuracy degradation. Customers are advised to use the `--enable-saturate-infinity` `--enable-mixed-precision-accumulation` compiler flags if Llama-2 13B accuracy issues occur.
- There is a known compiler issue for inference of some configurations of GPT-2 that can cause accuracy degradation. Customers are advised to use the `--enable-saturate-infinity` `--enable-mixed-precision-accumulation` compiler flags if GPT-2 accuracy issues occur.
- GPT-NeoX is sensitive to `fp16` and customers are advised to use only `amp="f32"` for GPT-NeoX.

## Release [0.9.474]

Date: 12/21/2023

## Summary

### What's new in this release

- [Llama] [Beta] Added support for Llama-2 70B.
- [Mistral] [Beta] Added support for Mistral 7B.
- [Beta] Added support for PyTorch 2.1.
- [Beta] Added support for Grouped Query Attention (GQA).
- [Beta] Added support for `safetensors` serialization.
- [Llama] [Beta] Added support for early stopping in the `sample_llama` function.
- [GPT2] [Beta] Added sparse attention support.
- [Stable] Added support for `BatchNorm`.
- Use the `--auto-cast=none` compiler flag by default for all models. This flag improves accuracy for `float32` operations.



## Resolved Issues

- Resolved an issue in `top_p` in the `sample_llama` function so that it now selects the same number of tokens that the Hugging Face `top_p` implementation selects.

## Known Issues and Limitations

- There is a known compiler issue for inference of some configurations of Llama-2 70B that can cause accuracy degradation. Customers are advised to use the `--enable-mixed-precision-accumulation` compiler flag if Llama-2 70B accuracy issues occur.
- There are known compiler issues impacting inference accuracy of certain model configurations of Llama-2-13b when `amp = fp16` is used. If this issue is observed, `amp=fp32` should be used as a work around. This issue will be addressed in future Neuron releases.

## Release [0.8.268]

Date: 10/26/2023

## Summary

### What's new in this release

- [Llama] [Beta] Added support for `int8` quantization for Llama.
- [BLOOM] [Beta] Added multi bucket context encoding support for BLOOM.
- [Beta] Added model Serialization for all supported models (except GPT-J and GPT-NeoX).
- [Beta] Added the ability to return output logit scores during sampling.
- [Stable] Added support for SOLU activation and GroupNorm.

## Resolved Issues

- [GPT2] Fixed an issue in `GPT2ForSamplingWithContextBroadcasting` where the input prompt would get truncated if it was longer than the `context_length_estimate`.

## Known Issues and Limitations

## Release [0.7.84]

Date: 09/15/2023

## Summary

### What's new in this release

- Use the `--model-type=transformer` compiler flag by default for all models. This flag improves performance and compilation time for all models. This flag replaces the `--model-type=transformer-inference` flag, which is now deprecated.

### Resolved Issues

- Fixed an issue where the `HuggingFaceGenerationModelAdapter` class falls back to serial context encoding for models that have parallel context encoding (`GPT2ForSamplingWithContextBroadcasting`, `LlamaForSampling`, etc.)
- [GPT2 / OPT] Fixed an issue in the parallel context encoding network where incorrect results could be generated due to incorrect masking logic.

### Known Issues and Limitations

- Some configurations of Llama and Llama-2 inference models fail compilation with the error `IndirectLoad/Save requires contiguous indirect access per partition`. This is fixed in the compiler version 2.10.0.35 (Neuron SDK 2.14.1).
- Some configurations of Llama and Llama-2 inference model fail compilation with the error `Too many instructions after unroll for function sg0000`. To mitigate this, please try with `-O1` compiler option (or `--optlevel 1`) by adding `os.environ["NEURON_CC_FLAGS"] = "-O1"` to your script or set in the environment. A complete fix will be coming in the future release which will not require this option. Note: Using `-O1` in the Llama-2 13B tutorial results in about 50% increase in latency compared to Neuron SDK 2.13.2. If this is not acceptable, please use compiler version from Neuron SDK 2.13.2.

### Release [0.6.106]

Date: 08/28/2023

## Summary

### What's new in this release

- Added support for Llama 2 (excluding grouped/multi-query versions, such as Llama 2 70B) [Beta]
- Improved the performance of BLOOM and Llama models [Beta]
- Reduced execution latency of token generation in tensor parallel models by improving thread synchronization. (supported in Llama only)
- Added an optimized vector implementation of RoPE positional embedding. (supported in Llama only)
- Added support for faster context encoding on sequences of varying lengths. This is implemented by allowing multiple buckets for parallel context encoding. During inference the best fit bucket is chosen. (supported in Llama/GPT-2 only)
- Added the Neuron Persistent Cache for compilation to automatically load pre-compiled model artifacts. (supported by all models)

- Improved compilation time by compiling models used for different sequence length buckets in parallel. (not supported in GPT-NeoX/GPT-J)

### Resolved Issues

- [Llama] Fixed an issue in the parallel context encoding network where incorrect results could be generated if the context length is shorter than the context length estimate
- [GPT2 / OPT] Fixed an issue in the parallel context encoding network where incorrect results could be generated

### Known Issues and Limitations

- The `HuggingFaceGenerationModelAdapter` class currently falls back to serial context encoding for models that have parallel context encoding (`GPT2ForSamplingWithContextBroadcasting`, `LlamaForSampling`, etc. )
- Beam search can introduce memory issues for large models
- There can be accuracy issues for the GPT-J model for certain use-cases

### Release [0.5.58]

Date: 7/21/2023

### Summary

#### What's new in this release

- Added support for GPT-NeoX models [Beta].
- Added support for BLOOM models [Beta].
- Added support for Llama models [Alpha].
- Added support for more flexible tensor-parallel configurations to GPT2, OPT, and BLOOM. The attention heads doesn't need to be evenly divisible by *tp\_degree* anymore. (Note: The *tp\_degree* still needs to satisfy the runtime topologies constraint for collective communication (i.e Allreduce). For more details on supported topologies, see: [Tensor-parallelism-support](https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-features/collective-communication.html) and <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/arch/neuron-features/collective-communication.html>.)
- Added multi-query / multi-group attention support for GPT2.

### Resolved Issues

- Fixed NaN issues for GPT2 model.
- Fixed OPT/GPT-NeoX gibberish output.
- Resolved an issue where NaN values could be produced when the `context_length` argument was used in GPT2/OPT.

### Known Issues and Limitations

- Missing cache reorder support for beam search.
- For more info, please see [features-support](#).

### Release [0.4.0]

Date: 6/14/2023

#### Summary

##### What's new in this release

- Added int8 weight storage for [GPT2](#) models.
- Improved prompt context encoding performance for [GPT2](#) models.
- Improved collective communications performance for tp-degrees 4, 8, and 24 on Inf2.
- Improved collective communications performance for tp-degrees 8 and 32 on Trn1.
- Support for the `--model-type=transformer-inference` compiler flag for optimized decoder-only LLM inference.

#### Resolved Issues

##### Incorrect GPT-J linear layer sharding

Added padding to the [GPT-J](#) linear layer to correctly handle odd vocabulary sizes.

##### Incorrect output with HuggingFace `beam_search()`

Issues where the HuggingFace `generate()` method produces incorrect results when `beam_search()` is used have been resolved.

### Release [0.3.0]

Date: 05/01/2023

#### Summary

##### What's new in this release

- Added `transformers-neuronx` artifacts to PyPI repository.
- Added support for the HuggingFace `generate()`.
- Added model serialization support for GPT2 models, including model saving, loading, and weight swapping.
- Added support for caching compiled artifacts.

- Improved performance by removing unnecessary KV-cache tensor resetting.
- Improved prompt context encoding performance (OPT, GPT2).

## Resolved Issues

### Incorrect GPT-J `amp_callback` import

Fixed the GPT-J demo to import the correct `amp_callback` function.

## Known Issues and Limitations

### Incorrect output with HuggingFace `beam_search()`

When the HuggingFace `generate()` method is configured to use `beam_search()`, this can produce incorrect results for certain configurations. It is recommended to use other generation methods such as `sample()` or `greedy_search()`. This will be fixed in a future Neuron release.

*This document is relevant for: Inf2, Trn1, Trn2*

- *Transformers Neuron (transformers-neuronx) release notes*

*This document is relevant for: Inf2, Trn1*

## Setup (transformers-neuronx)

If you already have setup your environment to run PyTorch NeuronX, you just need to install Transformers NeuronX library using the following instruction.

```
pip install transformers-neuronx --extra-index-url=https://pip.repos.neuron.amazonaws.com
```

If you are starting from scratch, Neuron Multi Framework DLAMI is recommended as it comes pre-installed with Transformers NeuronX virtual environment. You can refer to the instructions to launch a Neuron instance using Multi Framework DLAMI

## Developer Guide (transformers-neuronx)

- *Transformers NeuronX (transformers-neuronx) Developer Guide*

## Tutorials (transformers-neuronx)

- Hugging Face meta-llama/Llama-2-13b autoregressive sampling on Inf2 & Trn1
- Hugging Face facebook/opt-13b autoregressive sampling on Inf2 & Trn1
- Hugging Face facebook/opt-30b autoregressive sampling on Inf2 & Trn1
- Hugging Face facebook/opt-66b autoregressive sampling on Inf2

**Misc (transformers-neuronx)**

- *Transformers Neuron (transformers-neuronx) release notes*

*This document is relevant for:* Inf2, Trn1

*This document is relevant for:* Inf2, Trn1, Trn2

## 4.3 AWS Neuron Reference for NeMo Megatron

AWS Neuron Reference for NeMo Megatron is a library that includes modified versions of the open-source packages [NeMo](#) and [Apex](#) that have been adapted for use with AWS Neuron and AWS EC2 Trn1 instances. The library supports Tensor Parallel, Pipeline parallel and Data Parallel configurations for distributed training of large language models like GPT-3 175B. The APIs have been optimized for XLA based computation and high performance communication over Trainium instances. The library uses various techniques to improve memory utilization such as sequence parallelism which reduces activation memory footprint, selective or full activation checkpointing which allows larger model configurations to fit. SPMD optimizations are also used whenever possible to reduce the number of graphs obtained.

**Setup (neuronx-nemo-megatron)**

The library can be installed from [neuronx-nemo-megatron github repo](#)

**Tutorials (neuronx-nemo-megatron)**

- Launch a GPT-3 pretraining job using neuronx-nemo-megatron
- Launch a Llama 2 pretraining job using neuronx-nemo-megatron

### 4.3.1 Important Tips for Training with Neuron NeMo Megatron

**Do Not Create the Attention Mask**

If you are using your own data pipeline, do not create an attention mask for each record. Neuron NeMo Megatron is optimized to create an attention mask on Neuron Cores directly before use. Creating an attention mask per sample consumes excess CPU memory and often causes out of memory errors on CPU.

*This document is relevant for:* Inf2, Trn1, Trn2

## DEVELOPER FLOWS

### 5.1 Neuron DLAMI User Guide

#### Table of Contents

- *Neuron DLAMI Overview*
- *Neuron Multi Framework DLAMI*
  - *Multi Framework DLAMIs supported*
  - *Virtual Environments pre-installed*
- *Neuron Single Framework DLAMI*
  - *Single Framework DLAMIs supported*
  - *Virtual Environments pre-installed*
- *Neuron Base DLAMI*
  - *Base DLAMIs supported*
- *Using SSM parameters to find DLAMI id and trigger Cloud Automation flows*
  - *Finding specific DLAMI image id with the latest neuron release*
- *Other Resources*

### 5.1.1 Neuron DLAMI Overview

Neuron DLAMIs are an easy way to get started on Neuron SDK as they come pre-installed with Neuron SDK. Neuron currently supports 3 types of DLAMIs, multi-framework DLAMIs, single framework DLAMIs and base DLAMIs to easily get started on single Neuron instance. Below sections describe the supported Neuron DLAMIs, corresponding virtual environments and easy way to retrieve the DLAMI id using SSM parameters.

### 5.1.2 Neuron Multi Framework DLAMI

Neuron Deep Learning AMI (DLAMI) is a multi-framework DLAMI that supports multiple Neuron framework/libraries. Each DLAMI is pre-installed with Neuron drivers and support all Neuron instance types. Each virtual environment that corresponds to a specific Neuron framework/library comes pre-installed with all the Neuron libraries including Neuron compiler and Neuron runtime needed for you to easily get started.

**Note:** Tensorflow-neuron 2.10 (inf1) released in SDK v2.20.2 is not compatible with the latest runtime in v2.21 SDK. Code that compiles will face runtime errors with the latest SDK 2.21.1 version.

Neuron team is aware of this issue and we will ship a single-framework AMI for TF 2.10 inf1 in a future release.

You can use multi-framework DLAMIs from Neuron SDK v2.20.0 for inf1 workloads to avoid this issue. For example:

Deep Learning AMI Neuron (Ubuntu 22.04/AL2023) 20241027

Ubuntu22: ami-017ff4652165fd617

AL2023: ami-06fdb253ce8a32239

```
aws ec2 run-instances --image-id <ami-id>
```

Alternatively, you can use the latest Neuron DLAMIs on Ubuntu and run this command as a work-around:

```
sudo apt-get remove -y aws-neuronx-dkms aws-neuronx-collectives aws-neuronx-runtime-lib_
↪aws-neuronx-tools
sudo apt-get install aws-neuronx-dkms=2.18.* -y
sudo apt-get install aws-neuronx-collectives=2.22.* -y
sudo apt-get install aws-neuronx-runtime-lib=2.22.* -y
sudo apt-get install aws-neuronx-tools=2.19.* -y
```

<https://github.com/aws-neuron/aws-neuron-sdk/issues/1071> for more information on the issue.

### Multi Framework DLAMIs supported

| Operating System  | Neuron Instances Supported    | DLAMI Name                                   |
|-------------------|-------------------------------|--|
| Ubuntu 22.04      | Inf1, Inf2, Trn1, Trn1n, Trn2 | Deep Learning AMI Neuron (Ubuntu 22.04)      |
| Amazon Linux 2023 | Inf1, Inf2, Trn1, Trn1n, Trn2 | Deep Learning AMI Neuron (Amazon Linux 2023) |



## Virtual Environments pre-installed

| Neuron Framework/Libraries supported     | Virtual Environment                             |
|--|---|
| PyTorch 2.7 Torch NeuronX, NxD Core      | /opt/aws_neuronx_venv_pytorch_2_7               |
| PyTorch 2.7 NxD Training, Torch NeuronX  | /opt/aws_neuronx_venv_pytorch_2_7_nxd_training  |
| PyTorch 2.7 NxD Inference, Torch NeuronX | /opt/aws_neuronx_venv_pytorch_2_7_nxd_inference |
| Transformers NeuronX (PyTorch 2.7)       | /opt/aws_neuronx_venv_pytorch_2_7_transformers  |
| JAX 0.6 NeuronX                          | /opt/aws_neuronx_venv_jax_0_6                   |
| Tensorflow 2.10 NeuronX                  | /opt/aws_neuronx_venv_tensorflow_2_10           |
| Tensorflow 2.10 Neuron (Inf1)            | /opt/aws_neuron_venv_tensorflow_2_10_inf1       |
| PyTorch 1.13 Neuron (Inf1)               | /opt/aws_neuron_venv_pytorch_1_13_inf1          |

Within the PyTorch 2.7 NxD Training virtual environment, we have included a setup script that installs required dependencies for the package. To run this script, activate the virtual environment and run `setup_nxd_t.sh` and this will run [the setup steps here](#).

You can easily get started with the multi-framework DLAMI through AWS console by following this setup guide. If you are looking to use the Neuron DLAMI in your cloud automation flows, Neuron also supports [SSM parameters](#) to easily retrieve the latest DLAMI id.

### 5.1.3 Neuron Single Framework DLAMI

Neuron supports single framework DLAMIs that correspond to a single framework version (ex:- TensorFlow 2.10). Each DLAMI is pre-installed with Neuron drivers and supports all Neuron instance types. Each virtual environment corresponding to a specific Neuron framework/library comes pre-installed with all the relevant Neuron libraries including Neuron compiler and Neuron run-time.

#### Single Framework DLAMIs supported

| Framework              | Operating System  | Neuron Instances Supported | DLAMI Name   |
|------------------------|-------------------|----------------------------|--|
| PyTorch 2.7            | Ubuntu 22.04      | Inf2, Trn1, Trn1n, Trn2    | Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04)          |
| PyTorch 2.7            | Amazon Linux 2023 | Inf2, Trn1, Trn1n, Trn2    | Deep Learning AMI Neuron PyTorch 2.7 (Amazon Linux 2023)     |
| JAX 0.6                | Ubuntu 22.04      | Inf2, Trn1, Trn1n, Trn2    | Deep Learning AMI Neuron JAX 0.6 (Ubuntu 22.04)              |
| JAX 0.6                | Amazon Linux 2023 | Inf2, Trn1, Trn1n, Trn2    | Deep Learning AMI Neuron JAX 0.6 (Amazon Linux 2023)         |
| Tensorflow 2.10        | Ubuntu 22.04      | Inf2, Trn1, Trn1n, Trn2    | Deep Learning AMI Neuron TensorFlow 2.10 (Ubuntu 22.04)      |
| Tensorflow 2.10 (Inf1) | Ubuntu 22.04      | Inf1                       | Deep Learning AMI Neuron TensorFlow 2.10 Inf1 (Ubuntu 22.04) |
| PyTorch 1.13 (Inf1)    | Ubuntu 22.04      | Inf1                       | Deep Learning AMI Neuron PyTorch 1.13 Inf1 (Ubuntu 22.04)    |

## Virtual Environments pre-installed

| DLAMI Name   | Neuron Libraries supported               | Virtual Environment                             |
|--|--|---|
| Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04, Amazon Linux 2023) | PyTorch 2.7 Torch NeuronX, NxD Core      | /opt/aws_neuronx_venv_pytorch_2_7               |
| Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04, Amazon Linux 2023) | PyTorch 2.7 NxD Training, Torch NeuronX  | /opt/aws_neuronx_venv_pytorch_2_7_nxd_training  |
| Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04, Amazon Linux 2023) | PyTorch 2.7 NxD Inference, Torch NeuronX | /opt/aws_neuronx_venv_pytorch_2_7_nxd_inference |
| Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04, Amazon Linux 2023) | Transformers NeuronX PyTorch 2.7         | /opt/aws_neuronx_venv_pytorch_2_7_transformers  |
| Deep Learning AMI Neuron JAX 0.6 (Ubuntu 22.04, Amazon Linux 2023)     | JAX NeuronX 0.6                          | /opt/aws_neuronx_venv_jax_0_6                   |
| Deep Learning AMI Neuron PyTorch 1.13 (Ubuntu 22.04)                   | Pytorch Neuron (Inf1)                    | /opt/aws_neuron_venv_pytorch_1_13_inf1          |
| Deep Learning AMI Neuron TensorFlow 2.10 (Ubuntu 22.04)                | Tensorflow Neuronx                       | /opt/aws_neuronx_venv_tensorflow_2_10           |
| Deep Learning AMI Neuron TensorFlow 2.10 (Ubuntu 22.04)                | Tensorflow Neuron (Inf1)                 | /opt/aws_neuron_venv_tensorflow_2_10_inf1       |

You can easily get started with the single framework DLAMI through AWS console by following one of the corresponding setup guides . If you are looking to use the Neuron DLAMI in your cloud automation flows , Neuron also supports *SSM parameters* to easily retrieve the latest DLAMI id.

### 5.1.4 Neuron Base DLAMI

Neuron Base DLAMIs comes pre-installed with Neuron driver, EFA, and Neuron tools. Base DLAMIs might be relevant if you are extending the DLAMI for containerized applications.

#### Base DLAMIs supported

| Operating System  | Neuron Instances Supported    | DLAMI Name  |
|-------------------|-------------------------------|---|
| Amazon Linux 2023 | Inf1, Inf2, Trn1n, Trn1, Trn2 | Deep Learning Base Neuron AMI (Amazon Linux 2023) |
| Ubuntu 22.04      | Inf1, Inf2, Trn1n, Trn1, Trn2 | Deep Learning Base Neuron AMI (Ubuntu 22.04)      |

### 5.1.5 Using SSM parameters to find DLAMI id and trigger Cloud Automation flows

Neuron DLAMIs support AWS SSM parameters to easily find the Neuron DLAMI id. Currently we only support finding the latest DLAMI id that corresponds to latest Neuron SDK release with SSM parameter support. In the future releases, we will add support for finding the DLAMI id using SSM parameters for a specific Neuron release.

## Finding specific DLAMI image id with the latest neuron release

You can find the DLAMI that supports latest Neuron SDK by using the SSM get-parameter.

```
aws ssm get-parameter \
--region us-east-1 \
--name <dlami-ssm-parameter-prefix>/latest/image_id \
--query "Parameter.Value" \
--output text
```

The SSM parameter prefix for each DLAMI can be seen below

### SSM Parameter Prefix

| AMI Name  | SSM parameter Prefix  |
|---|---|
| Deep Learning AMI Neuron (Ubuntu 22.04)                   | /aws/service/neuron/dlami/multi-framework/ubuntu-22.04      |
| Deep Learning AMI Neuron (Amazon Linux 2023)              | /aws/service/neuron/dlami/multi-framework/amazon-linux-2023 |
| Deep Learning AMI Neuron PyTorch 2.7 (Ubuntu 22.04)       | /aws/service/neuron/dlami/pytorch-2.7/ubuntu-22.04          |
| Deep Learning AMI Neuron PyTorch 2.7 (Amazon Linux 2023)  | /aws/service/neuron/dlami/pytorch-2.7/amazon-linux-2023     |
| Deep Learning AMI Neuron JAX 0.6 (Ubuntu 22.04)           | /aws/service/neuron/dlami/jax-0.6/ubuntu-22.04              |
| Deep Learning AMI Neuron JAX 0.6 (Amazon Linux 2023)      | /aws/service/neuron/dlami/jax-0.6/amazon-linux-2023         |
| Deep Learning AMI Neuron PyTorch 1.13 Inf1 (Ubuntu 22.04) | /aws/service/neuron/dlami/pytorch-1.13-inf1/ubuntu-22.04    |
| Deep Learning AMI Neuron TensorFlow 2.10 (Ubuntu 22.04)   | /aws/service/neuron/dlami/tensorflow-2.10/ubuntu-22.04      |
| Deep Learning Base Neuron AMI (Amazon Linux 2023)         | /aws/service/neuron/dlami/base/amazon-linux-2023            |
| Deep Learning Base Neuron AMI (Ubuntu 22.04)              | /aws/service/neuron/dlami/base/ubuntu-22.04                 |

For example to find the latest DLAMI id for Multi-Framework DLAMI (Ubuntu 22) you can use the following

```
aws ssm get-parameter \
--region us-east-1 \
--name /aws/service/neuron/dlami/multi-framework/ubuntu-22.04/latest/image_id \
--query "Parameter.Value" \
--output text
```

You can find all available parameters supported in Neuron DLAMIs via CLI

```
aws ssm get-parameters-by-path \
--region us-east-1 \
--path /aws/service/neuron \
--recursive
```

You can also view the SSM parameters supported in Neuron through AWS parameter store by selecting the “Neuron” service.

## Use SSM Parameter to launch instance directly via CLI

You can use CLI to find the latest DLAMI id and also launch the instance simultaneously. Below code snippet shows an example of launching inf2 instance using multi-framework DLAMI

```
aws ec2 run-instances \
--region us-east-1 \
--image-id resolve:ssm:/aws/service/neuron/dlami/tensorflow-2.10/ubuntu-22.04/latest/
↪image_id \
--count 1 \
--instance-type inf2.48xlarge \
--key-name <my-key-pair> \
--security-groups <my-security-group>
```

## Use SSM alias in EC2 launch templates

SSM Parameters can also be used directly in launch templates. So, you can update your Auto Scaling groups to use new AMI IDs without needing to create new launch templates or new versions of launch templates each time an AMI ID changes. Ref: <https://docs.aws.amazon.com/autoscaling/ec2/userguide/using-systems-manager-parameters.html>

### 5.1.6 Other Resources

<https://docs.aws.amazon.com/dlami/latest/devguide/what-is-dlami.html>

<https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html>

<https://docs.aws.amazon.com/systems-manager/latest/userguide/systems-manager-parameter-store.html>

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.2 Neuron Containers

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 5.2.1 Getting started with Neuron DLC using Docker

#### Training

##### Launch Trn1 Instance

- Please follow the instructions at [launch an Amazon EC2 Instance](#) to Launch an instance, when choosing the instance type at the EC2 console. Please make sure to select the correct instance type.
- To get more information about instances sizes and pricing see: [Trn1 web page](#), [Inf2 web page](#), [Inf1 web page](#)
- Select your Amazon Machine Image (AMI) of choice, please note that Neuron supports Amazon Linux 2 AMI(HVM) - Kernel 5.10.
- When launching a Trn1, please adjust your primary EBS volume size to a minimum of 512GB.
- After launching the instance, follow the instructions in [Connect to your instance](#) to connect to the instance

**Note:** If you are facing a connectivity issue during the model loading process on a Trn1 instance with Ubuntu, that could probably be because of Ubuntu limitations with multiple interfaces. To solve this problem, please follow the steps mentioned [here](#).

Users are highly encouraged to use DLAMI to launch the instances, since DLAMIs come with the required fix.

## Install Drivers

```
# Configure Linux for Neuron repository updates

sudo tee /etc/yum.repos.d/neuron.repo > /dev/null <<EOF
[neuron]
name=Neuron YUM Repository
baseurl=https://yum.repos.neuron.amazonaws.com
enabled=1
metadata_expire=0
EOF
sudo rpm --import https://yum.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-NEURON.
↪PUB

# Update OS packages
sudo yum update -y

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Remove preinstalled packages and Install Neuron Driver and Runtime
sudo yum remove aws-neuron-dkms -y
sudo yum remove aws-neuronx-dkms -y
sudo yum install aws-neuronx-dkms-2.* -y

# Install EFA Driver(only required for multi-instance training)
curl -O https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz
wget https://efa-installer.amazonaws.com/aws-efa-installer.key && gpg --import aws-efa-
↪installer.key
cat aws-efa-installer.key | gpg --fingerprint
wget https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz.sig && gpg --
↪verify ./aws-efa-installer-latest.tar.gz.sig
tar -xvf aws-efa-installer-latest.tar.gz
cd aws-efa-installer && sudo bash efa_installer.sh --yes
cd
sudo rm -rf aws-efa-installer-latest.tar.gz aws-efa-installer
```

## Install Docker

```
sudo yum install -y docker.io
sudo usermod -aG docker $USER
```

Logout and log back in to refresh membership.

## Verify Docker

```
docker run hello-world
```

Expected result:

```
Hello from Docker!
This message shows that your installation appears to be working correctly.

To generate this message, Docker took the following steps:
1. The Docker client contacted the Docker daemon.
2. The Docker daemon pulled the "hello-world" image from the Docker Hub.
   (amd64)
3. The Docker daemon created a new container from that image which runs the
   executable that produces the output you are currently reading.
4. The Docker daemon streamed that output to the Docker client, which sent it
   to your terminal.

To try something more ambitious, you can run an Ubuntu container with:
$ docker run -it ubuntu bash

Share images, automate workflows, and more with a free Docker ID:
https://hub.docker.com/

For more examples and ideas, visit:
https://docs.docker.com/get-started/
```

## Verify Neuron Component

Once the environment is setup, a container can be started with `--device=/dev/neuron#` to specify desired set of Inferentia/Trainium devices to be exposed to the container. To find out the available neuron devices on your instance, use the command `ls /dev/neuron*`.

When running `neuron-ls` inside a container, you will only see the set of exposed Trainiums. For example:

```
docker run --device=/dev/neuron0 neuron-test neuron-ls
```

Would produce the following output in `trn1.32xlarge`:

```
+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | PCI   |
| DEVICE | CORES  | MEMORY | BDF   |
+-----+-----+-----+-----+
| 0      | 2      | 32 GB  | 10:1c.0 |
+-----+-----+-----+-----+
```

## Build and Run Docker Image

- [how-to-build-neuron-container](#)

## Run Tutorial

[tutorial-training](#)

## Inference

## Install Drivers

```
# Configure Linux for Neuron repository updates
sudo tee /etc/yum.repos.d/neuron.repo > /dev/null <<EOF
[neuron]
name=Neuron YUM Repository
baseurl=https://yum.repos.neuron.amazonaws.com
enabled=1
metadata_expire=0
EOF
sudo rpm --import https://yum.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-NEURON.
↪PUB

# Update OS packages
sudo yum update -y

#####
↪#####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade to
↪latest Neuron driver
#####
↪#####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum install aws-neuron-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#         Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####
```

## Install Docker

```
sudo yum install -y docker.io
sudo usermod -aG docker $USER
```

Logout and log back in to refresh membership.

## Verify Docker

```
docker run hello-world
```

Expected result:

```
Hello from Docker!
This message shows that your installation appears to be working correctly.

To generate this message, Docker took the following steps:
1. The Docker client contacted the Docker daemon.
2. The Docker daemon pulled the "hello-world" image from the Docker Hub.
   (amd64)
3. The Docker daemon created a new container from that image which runs the
   executable that produces the output you are currently reading.
4. The Docker daemon streamed that output to the Docker client, which sent it
   to your terminal.

To try something more ambitious, you can run an Ubuntu container with:
$ docker run -it ubuntu bash

Share images, automate workflows, and more with a free Docker ID:
https://hub.docker.com/

For more examples and ideas, visit:
https://docs.docker.com/get-started/
```

## Verify Neuron Component

Once the environment is setup, a container can be started with `--device=/dev/neuron#` to specify desired set of Inferentia/Trainium devices to be exposed to the container. To find out the available neuron devices on your instance, use the command `ls /dev/neuron*`.

When running `neuron-ls` inside a container, you will only see the set of exposed Inferentias. For example:

```
docker run --device=/dev/neuron0 neuron-test neuron-ls
```

Would produce the following output in `inf1.xlarge`:

| PCI BDF      | LOGICAL ID | NEURON CORES | MEMORY CHANNEL 0 | MEMORY CHANNEL 1 | EAST | WEST |
|--------------|------------|--------------|------------------|------------------|------|------|
| 0000:00:1f.0 | 0          | 4            | 4096 MB          | 4096 MB          | 0    | 0    |



## Run Tutorial

tutorial-infer

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.2.2 Neuron Deep Learning Containers

### Table of Contents

- *Overview*
- *Inference Containers*
- *Training Containers*
- *Getting started with Neuron DLC using Docker*
- *Using containers on AWS services*
  - Amazon EKS
  - Amazon ECS
  - Amazon SageMaker
  - AWS Batch
- *Customizing Neuron Deep Learning Containers*

### Overview

AWS Deep Learning Containers (DLCs) provide a set of Docker images that are pre-installed with deep learning frameworks. The containers are optimized for performance and available in Amazon Elastic Container Registry (Amazon ECR). DLCs make it straightforward to deploy custom ML environments in a containerized manner, while taking advantage of the portability and reproducibility benefits of containers.

AWS Neuron DLCs are a set of Docker images for training and serving models on AWS Trainium and Inferentia instances using AWS Neuron SDK. The sections below list all of the AWS Neuron DLCs, as well as the AWS DLCs that come pre-installed with the Neuron SDK.

## Inference Containers

| DLC Name   | DLC Link(s)   | Tutorial(s)   |
|--|---|---|
| Neuron Inference Containers                                    | <a href="#">Neuron PyTorch Inference Containers</a><br><a href="#">Neuronx PyTorch Inference Containers</a>                       | <a href="#">tutorial-infer</a><br><a href="#">torchserve-neuron</a> |
| Large Model Inference (LMI)/Deep Java Library (DJL) Containers | <a href="#">LMI Containers</a>  |   |
| HuggingFace Inference Containers                               | <a href="#">HuggingFace Text Generation Inference (TGI) Containers</a><br><a href="#">HuggingFace Neuron Inference Containers</a> |   |
| Triton Inference Containers                                    | <a href="#">NVIDIA Triton Inference Containers</a>  |   |

## Training Containers

| DLC Name                        | DLC Link(s)  | Tutorial(s)                       |
|---------------------------------|--|-----------------------------------|
| Neuron Training Containers      | <a href="#">Neuronx PyTorch Training Containers</a>    | <a href="#">tutorial-training</a> |
| HuggingFace Training Containers | <a href="#">HuggingFace Neuron Training Containers</a> |                                   |

## Getting started with Neuron DLC using Docker

*Getting started with Neuron DLC using Docker*

## Using containers on AWS services

### Amazon EKS

### Amazon ECS

### Amazon SageMaker

### AWS Batch

## Customizing Neuron Deep Learning Containers

Deep Learning Containers can be customized to fit your specific project needs. To read more, visit [Customize Neuron DLC](#).

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### 5.2.3 Customize Neuron DLC

#### Table of Contents

- [Description](#)
- [Method 1: Using DLC as a Base Image](#)
- [Method 2: Modifying Published Dockerfiles](#)

#### Description

This guide covers how to customize and extend the Neuron Deep Learning Container (DLC) to fit your specific project needs. You can customize the DLC either by using the DLC as a base image in your Dockerfile or by modifying published Dockerfiles on GitHub.

#### Method 1: Using DLC as a Base Image

1. Create a New Dockerfile. In your Dockerfile, specify the Neuron DLC as your base image using the FROM directive.
2. Complete the Dockerfile. You can add additional packages, change the base environment, or any other modifications that suit your project. [AWS Batch Training](#) is a good example which needs customize Neuron DLC by using it as the base image. From its [Dockerfile](#), we can find the customized container copies llama\_batch\_training.sh to the container and runs it.
3. Navigate to the directory containing your Dockerfile and build your custom container.

#### Method 2: Modifying Published Dockerfiles

1. Visit the [Neuron DLC Github repo](#) and locate the Dockerfile for the container you wish to customize.
2. Modify the Dockerfile as needed. You can add additional packages, change the base environment, or any other modifications that suit your project. For example, if you do not need to use Neuron tools in your scenario and want to make the container smaller, you can remove aws-neuronx-tools at this [line](#).
3. Navigate to the directory containing your Dockerfile and build your custom container.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 5.2.4 Neuron Plugins for Containerized Environments

This section summarizes various neuron infrastructure artifacts for containerized environments.

- **Neuron Node Problem Detector** - This plugin enhances resiliency by detecting and remediating errors. For detailed instructions on running this plugin in EKS environment, please refer to [EKS Setup For Neuron](#). To leverage this plugin on ECS, please refer to [Neuron Problem Detector And Recovery](#)
- **Neuron Device Plugin** - The Neuron device plugin manages Neuron hardware resources in a Kubernetes environment. It integrates with the Kubernetes device plugin framework to advertise and manage Neuron resources, making them available for use by Pods. For more information on using Neuron with Kubernetes, please refer to [EKS Setup For Neuron](#)

- Neuron Scheduler Extension - Neuron scheduler extension is a Kubernetes artifact which helps with optimal allocation of neuron cores. Installing scheduler extension is optional if a workload pod consumes all neuron resources on a node. For more information on using Neuron with Kubernetes, please refer to [EKS Setup For Neuron](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.2.5 Neuron Containers FAQ

### Table of Contents

- [Where can I find DLC images](#)
- [What is OCI Neuron Hook and do we need that](#)
- [What container runtimes are supported](#)
- [How to expose Neuron Devices to Container](#)
- [How to expose Neuron Cores to Container](#)
- [Can Neuron Devices be shared by different Containers running in the same Host](#)
- [Can Neuron Cores be shared by different Containers running in the same Host](#)
- [When would you use Neuron K8 Scheduler Extension](#)
- [How to add EFA devices to the container](#)
- [Can distributed training jobs be run without EFA devices in container](#)

### Where can I find DLC images

- The Inference/Training DLC images can be found [here](#).
- In the [DLC release page](#) do a search for neuron to get the ECR repo location of specific neuron DLC release.

### What is OCI Neuron Hook and do we need that

Neuron devices are exposed to the containers using the `--device` option in the docker run command. Docker runtime (runc) does not yet support the ALL option to expose all neuron devices to the container.

With OCI neuron hook support is added to expose ALL devices to container using an environment variable, “AWS\_NEURON\_VISIBLE\_DEVICES=ALL”. For more details please refer [oci neuron hook](#)

In Kubernetes, if we are using the device plugin version 1.7 & below, then the oci neuron hook is needed. If using device plugin version `>= 1.8` then oci neuron hook is not needed

## What container runtimes are supported

Neuron containers have been tested to work with docker, containerd, cri-o runtimes without any changes. If the oci neuron hook is used then they need to be enabled in the runtime config. For more details please refer [oci neuron hook](#)

## How to expose Neuron Devices to Container

Neuron Device: Represents the number of Inferentia/Trainium chips in the instance. Refer [Container Devices](#) for more details

## How to expose Neuron Cores to Container

Neuron Core: Represents the number of Neuron Cores in the instance. Refer [Container Cores](#) for more details. Each Inferentia1 device has 4 Neuron Cores and each Inferentia2 and Trainium1 device has 2 Neuron Cores. When the devices are exposed to the containers all the cores in the device are available for use in the container. Please refer [NeuronX Runtime Configuration](#) to see how the environment variables NEURON\_RT\_VISIBLE\_CORES and NEURON\_RT\_NUM\_CORES can be used to assign core to containers

## Can Neuron Devices be shared by different Containers running in the same Host

Yes, except in Kubernetes environment where the devices cannot be shared

## Can Neuron Cores be shared by different Containers running in the same Host

No

## When would you use Neuron K8 Scheduler Extension

The neuron cores/devices that are exposed to the container needs to be contiguous. The kubernetes device plugin does not guarantee the devices to be contiguous. The K8 Neuron Scheduler Extension takes care of assigning contiguous devices to the containers.

## How to add EFA devices to the container

The EFA devices are exposed to the container using the `--device` option

```
--device /dev/infiniband/uverbs0
```

In a Kubernetes environment, the EFA device plugin is used to detect and advertise the available EFA interfaces. The EFA device plugin can be installed using the [Helm chart provided by Amazon EKS](#)

```
helm repo add eks https://aws.github.io/eks-charts
helm install aws-efa-k8s-device-plugin --namespace kube-system eks/aws-efa-k8s-device-
↪plugin
```

Once the plugin is deployed, applications can use the resource type `vpc.amazonaws.com/efa` in a pod request spec

```
resources:
  limits:
    vpc.amazonaws.com/efa: 4
```

### **Can distributed training jobs be run without EFA devices in container**

No. For distributed training jobs on Trainium, all EFA interfaces provided by trn1.32xlarge need to be attached to the container

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

In this section, you'll find resources to help you use containers for accelerating your deep learning models on Inferentia and Trainium instances.

## **5.2.6 Getting started with Neuron DLC using Docker**

AWS Neuron Deep Learning Containers (DLCs) are a set of Docker images for training and serving models on AWS Trainium and Inferentia instances using AWS Neuron SDK. To build a Neuron container using Docker, please refer to *Getting started with Neuron DLC using Docker*.

## **5.2.7 Neuron Deep Learning Containers**

In most cases, it is recommended to use a preconfigured [Deep Learning Container \(DLC\)](#) from AWS. Each DLC is pre-configured to have all of the Neuron components installed and is specific to the chosen ML Framework. For more details on Neuron Deep Learning Containers, please refer to *Neuron Deep Learning Containers*.

## **5.2.8 Customize Neuron DLC**

Neuron DLC can be customized as needed. To learn more about how to customize the Neuron Deep Learning Container (DLC) to fit your specific project needs, please refer to *Customize Neuron DLC*.

## **5.2.9 Neuron Plugins for Containerized Environments**

Neuron provides plugins for better observability and fault tolerance. For more information on the plugins, please refer to *Neuron Plugins for Containerized Environments*.

## **5.2.10 Neuron Containers FAQ**

For frequently asked questions and troubleshooting, please refer to *Neuron Containers FAQ*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## **5.3 AWS Workload Orchestration**

Neuron can be used in a wide selection of development flows. Each flow has its own starting point and requirements which are required to enable deep learning acceleration with AWS Neuron.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 5.3.1 Amazon EKS

This document is relevant for: Inf1, Inf2, Trn1, Trn2

#### Using Neuron with Amazon EKS

##### Table of Contents

- *EKS Setup For Neuron*
- *Prerequisites*
- *Neuron Device Plugin*
  - *Deploy Neuron Device Plugin*
- *Neuron Scheduler Extension*
  - *Container Device Allocation On Different Instance Types*
  - *Deploy Neuron Scheduler Extension*
- *Neuron Node Problem Detector Plugin*
  - *Permissions for Neuron Node Problem Detector Plugin*
  - *Deploy Neuron Node Problem Detector And Recovery*
- *Neuron Monitor Daemonset*
  - *Deploy Neuron Monitor Daemonset*
- *Neuron Helm Chart*

#### EKS Setup For Neuron

Customers that use Kubernetes can conveniently integrate Inf1/Trn1 instances into their workflows. This section will go through steps for setting up EKS cluster for Neuron.

#### Prerequisites

Please refer to [EKS instructions](#) to create a cluster. Once the cluster is ACTIVE, please add nodes to the cluster. We recommend using node template for neuron nodes. Following example demonstrates how to add neuron nodes using node template. The example adds managed nodes using `eksctl` tool. For more details, please refer to [EKS User Guide](#).

As first step, please create a script to capture the parameters for the node template:

```
#!/bin/bash

CLUSTER_NAME=$1
CLUSTER_SG=$(eksctl get cluster $CLUSTER_NAME -o json|jq -r "[0].ResourcesVpcConfig.
↳ClusterSecurityGroupId")
VPC_ID=$(eksctl get cluster $CLUSTER_NAME -o json|jq -r "[0].ResourcesVpcConfig.VpcId")

cat <<EOF > cfn_params.json
```

(continues on next page)

(continued from previous page)

```
[
  {
    "ParameterKey": "ClusterName",
    "ParameterValue": "$CLUSTER_NAME"
  },
  {
    "ParameterKey": "ClusterControlPlaneSecurityGroup",
    "ParameterValue": "$CLUSTER_SG"
  },
  {
    "ParameterKey": "VpcId",
    "ParameterValue": "$VPC_ID"
  }
]
EOF
```

These parameters include the name of the cluster, the security group the nodes can use to connect to the control plane and the vpcid. Next, get the node group template from tutorial below -

```
wget https://raw.githubusercontent.com/aws-neuron/aws-neuron-eks-samples/master/dp_bert_
hf_pretrain/cfn/eks_trn1_ng_stack.yaml
```

This template file has a few important config settings -

- It places the node in a placement group. This optimizes the network speed between the nodes.
- The template installs the EFA driver. Please note that the libfabric version should match between the AMI and the workload containers.
- It uses the [EKS optimized accelerated AMI](#) which has the necessary neuron components installed. The template uses AMI for Kubernetes version 1.25. Please update to appropriate version.
- The template adds trn1.32xlarge nodes to the cluster. Please update to the desired instance type.
- Trn2 instance types use a default LNC (Logical NeuronCore Configuration) setting of 2, if you want to change it to 1, update the UserData section of the launch template to a new LNC setting as shown below, and deploy the new/updated version of launch template.

```
====BOUNDARY==
Content-Type: text/x-shellscript; charset="us-ascii"

#!/bin/bash
set -ex
config_dir=/opt/aws/neuron
config_file=${config_dir}/logical_nc_config
[ -d "$config_dir" ] || mkdir -p "$config_dir"
[ -f "$config_file" ] || touch "$config_file"
if ! grep -q "^NEURON_LOGICAL_NC_CONFIG=1$" "$config_file" 2>/dev/null; then
  printf "NEURON_LOGICAL_NC_CONFIG=1" >> "$config_file"
fi
====BOUNDARY====
```

Finally, run the following command to create cloud formation stack:



```
aws cloudformation create-stack \
--stack-name eks-trn1-ng-stack \
--template-body file://eks_trn1_ng_stack.yaml \
--parameters file://cfn_params.json \
--capabilities CAPABILITY_IAM
```

The above command will create a stack named eks-trn1-ng-stack, which will be visible in cloudformation. Please wait for that stack creation to complete before proceeding to next step.

Now we are ready to add the nodes. The example will demonstrate creating node groups using eksctl tool.

Please run following command to determine the AZs:

```
aws ec2 describe-availability-zones \
--region $REGION_CODE \
--query "AvailabilityZones[]" \
--filters "Name=zone-id,Values=$1" \
--query "AvailabilityZones[].ZoneName" \
--output text
```

Next, create a script named create\_ng\_yaml.sh to generate node group yaml. The arguments to the script include the region, AZs, cluster name and name of the cloudformation stack created earlier (eks-trn1-ng-stack in case of this example):

```
#!/bin/bash

REGION_CODE=$1
EKSAZ1=$2
EKSAZ2=$3
CLUSTER_NAME=$4
STACKNAME=$5

LT_ID_TRN1=$(aws cloudformation describe-stacks --stack-name $STACKNAME \
--query "Stacks[0].Outputs[?OutputKey=='LaunchTemplateIdTrn1'].OutputValue" \
--output text)

cat <<EOF > trn1_nodegroup.yaml
apiVersion: eksctl.io/v1alpha5
kind: ClusterConfig

metadata:
  name: $CLUSTER_NAME
  region: $REGION_CODE
  version: "1.28"

iam:
  withOIDC: true

availabilityZones: ["$EKSAZ1","$EKSAZ2"]

managedNodeGroups:
- name: trn1-32x1-ng1
  launchTemplate:
    id: $LT_ID_TRN1
```

(continues on next page)

(continued from previous page)

```

minSize: 1
desiredCapacity: 1
maxSize: 1
availabilityZones: ["$EKSAZ1"]
privateNetworking: true
efaEnabled: true
EOF

```

Run the above script. It should produce a yaml similar to -

```

apiVersion: eksctl.io/v1alpha5
kind: ClusterConfig

metadata:
  name: nemo2
  region: us-west-2
  version: "1.25"

iam:
  withOIDC: true

availabilityZones: ["us-west-2d", "us-west-2c"]

managedNodeGroups:
- name: trn1-32xl-ng1
  launchTemplate:
    id: lt-093c222b35ea89009
  minSize: 1
  desiredCapacity: 1
  maxSize: 1
  availabilityZones: ["us-west-2d"]
  privateNetworking: true
  efaEnabled: true

```

The example shows kubernetes version 1.25. Please update the version as needed. This yaml can now be used with eksctl.

```
eksctl create nodegroup -f trn1_nodegroup.yaml
```

This will add the nodes to the cluster. Please wait for the nodes to be 'Ready'. This can be verified using the get node command.

If you are running a distributed training or inference job, you will need EFA resources. Please install the EFA device plugin using instructions at [EFA device plugin repository](#).

Next, we will install the Neuron Device Plugin.

## Neuron Device Plugin

Neuron device plugin exposes Neuron cores & devices to kubernetes as a resource. `aws.amazon.com/neuroncore` and `aws.amazon.com/neuron` are the resources that the neuron device plugin registers with the kubernetes. `aws.amazon.com/neuroncore` is used for allocating neuron cores to the container. `aws.amazon.com/neuron` is used for allocating neuron devices to the container. When resource name 'neuron' is used, all the cores belonging to the device will be allocated to container.

## Deploy Neuron Device Plugin

- Make sure prerequisite are satisfied
- Apply the Neuron device plugin as a daemonset on the cluster with the following command

```
helm upgrade --install neuron-helm-chart oci://public.ecr.aws/neuron/neuron-
helm-chart \
--set "npd.enabled=false"
```

- Verify that neuron device plugin is running

```
kubectl get ds neuron-device-plugin -n kube-system
```

Expected result (with 2 nodes in cluster):

| NAME                   | DESIRED | CURRENT | READY | UP-TO-DATE | AVAILABLE |
|------------------------|---------|---------|-------|------------|-----------|
| ↳NODE SELECTOR    AGE  |         |         |       |            |           |
| neuron-device-plugin   | 2       | 2       | 2     | 2          | 2         |
| ↳<none>            18h |         |         |       |            |           |

- Verify that the node has allocatable neuron cores and devices with the following command

```
kubectl get nodes "-o=custom-columns=NAME:.metadata.name,NeuronCore:.status.
allocatable.aws.amazon.com/neuroncore"
```

Expected result:

| NAME  | NeuronCore |
|---|------------|
| ip-192-168-65-41.us-west-2.compute.internal | 32         |
| ip-192-168-87-81.us-west-2.compute.internal | 32         |

```
kubectl get nodes "-o=custom-columns=NAME:.metadata.name,NeuronDevice:.
status.allocatable.aws.amazon.com/neuron"
```

Expected result:

| NAME  | NeuronDevice |
|---|--------------|
| ip-192-168-65-41.us-west-2.compute.internal | 16           |
| ip-192-168-87-81.us-west-2.compute.internal | 16           |

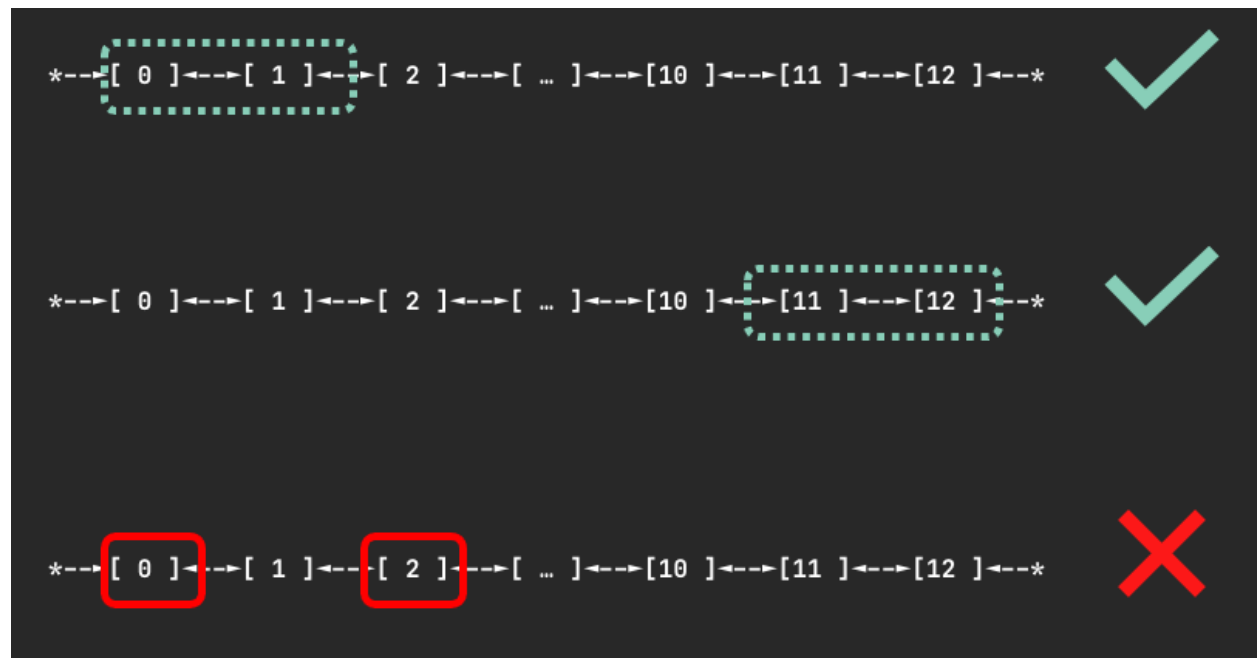
## Neuron Scheduler Extension

The Neuron scheduler extension is required for scheduling pods that require more than one Neuron core or device resource. For a graphical depiction of how the Neuron scheduler extension works, see `k8s-neuron-scheduler-flow`. The Neuron scheduler extension finds sets of directly connected devices with minimal communication latency when scheduling containers. On Inf1 and Inf2 instance types where Neuron devices are connected through a ring topology, the scheduler finds sets of contiguous devices. For example, for a container requesting 3 Neuron devices the scheduler might assign Neuron devices 0,1,2 to the container if they are available but never devices 0,2,4 because those devices are not directly connected. On Trn1.32xlarge and Trn1n.32xlarge instance types where devices are connected through a 2D torus topology, the Neuron scheduler enforces additional constraints that containers request 1, 4, 8, or all 16 devices. If your container requires a different number of devices, such as 2 or 5, we recommend that you use an Inf2 instance instead of Trn1 to benefit from more advanced topology.

## Container Device Allocation On Different Instance Types

The Neuron scheduler extension applies different rules when finding devices to allocate to a container on Inf1 and Inf2 instances than on Trn1. These rules ensure that when users request a specific number of resources, Neuron delivers *consistent* and *high* performance regardless of which cores and devices are assigned to the container.

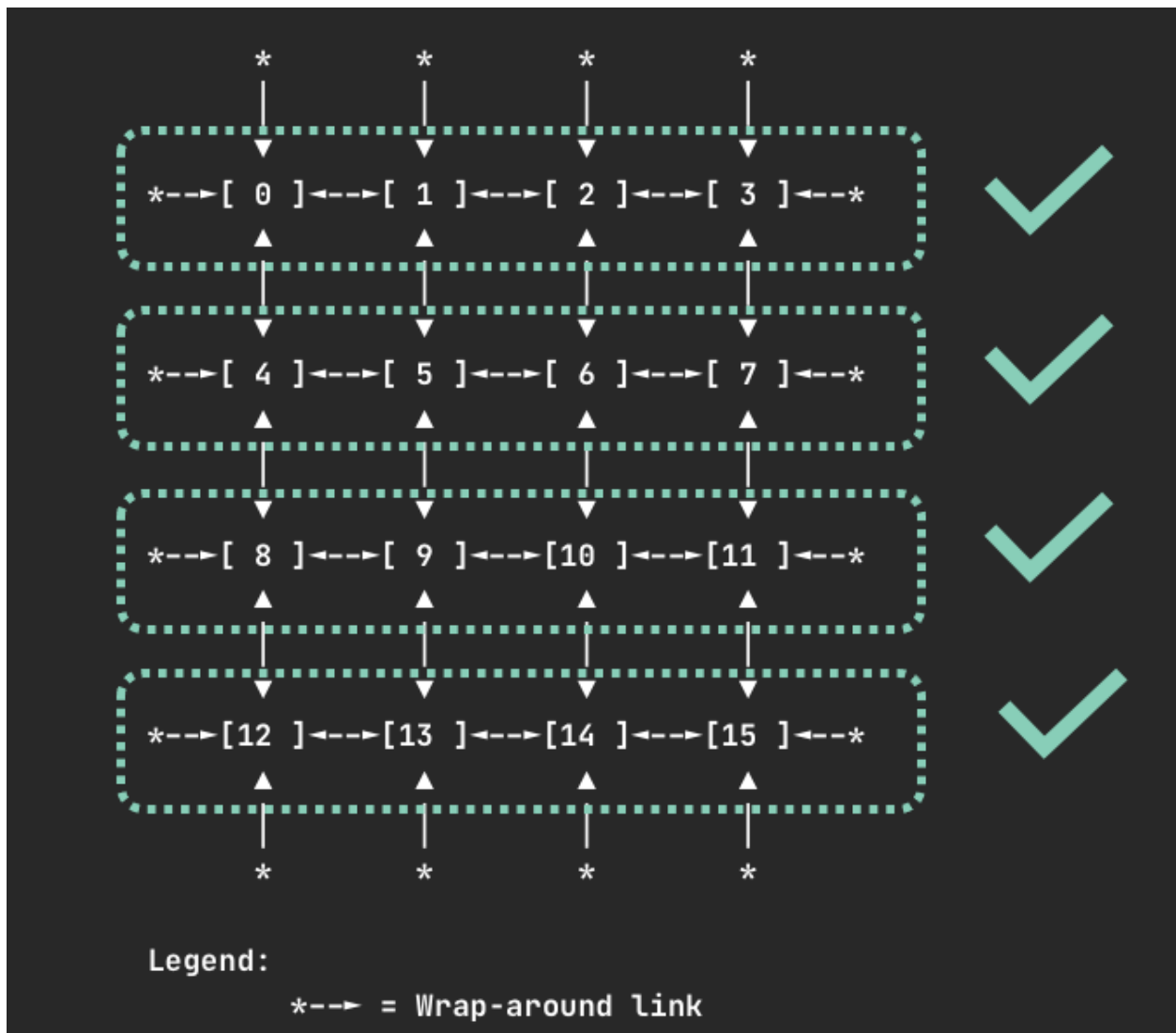
On Inf1 and Inf2 Neuron devices are connected through a ring topology. There are no restrictions on the number of devices requested as long as it is fewer than the number of devices on a node. When the user requests N devices, the scheduler finds a node where N contiguous devices are available. It will never allocate non-contiguous devices to the same container. The figure below shows examples of device sets on an Inf2.48xlarge node which could be assigned to a container given a request for 2 devices.



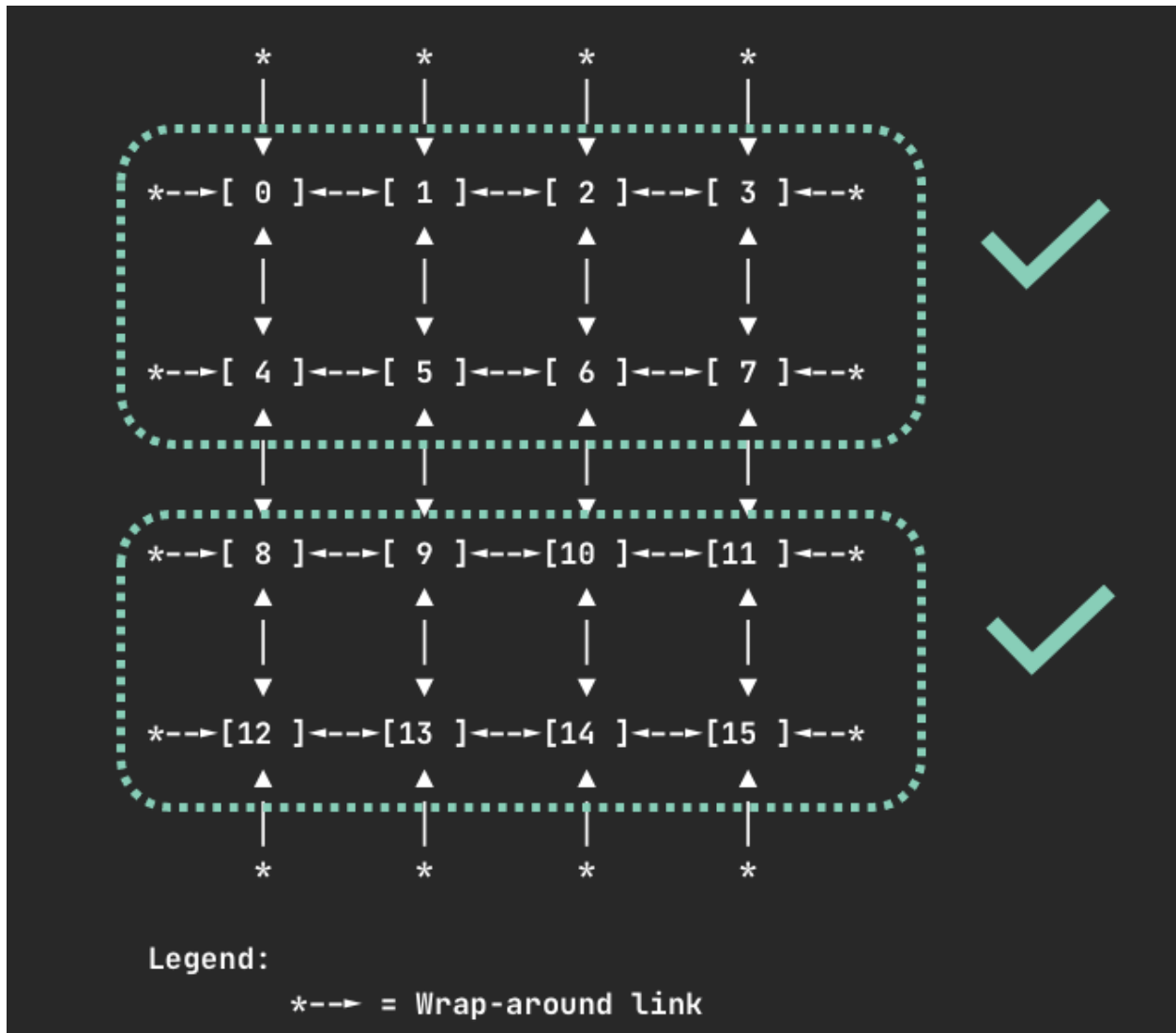
Devices on Trn1.32xlarge and Trn1n.32xlarge nodes are connected via a 2D torus topology. On Trn1 nodes containers can request 1, 4, 8, or all 16 devices. In the case you request an invalid number of devices, such as 7, your pod will not be scheduled and you will receive a warning:

Instance type `trn1.32xlarge` does not support requests for device: 7. Please request a different number of devices.

When requesting 4 devices, your container will be allocated one of the following sets of devices if they are available.



When requesting 8 devices, your container will be allocated one of the following sets of devices if they are available.



For all instance types, requesting one or all Neuron cores or devices is valid.

## Deploy Neuron Scheduler Extension

### Multiple Scheduler Approach

In cluster environments where there is no access to default scheduler, the neuron scheduler extension can be used with another scheduler. A new scheduler is added (along with the default scheduler) and then the pod's that needs to run the neuron workload use this new scheduler. Neuron scheduler extension is added to this new scheduler. EKS natively does not yet support the neuron scheduler extension and so in the EKS environment this is the only way to add the neuron scheduler extension.

- Make sure Neuron device plugin is running
- Install the neuron-scheduler-extension

```
helm upgrade --install neuron-helm-chart oci://public.ecr.aws/neuron/neuron-
helm-chart \
```

(continues on next page)

(continued from previous page)

```
--set "scheduler.enabled=true" \
--set "npd.enabled=false"
```

- Check there are no errors in the my-scheduler pod logs and the k8s-neuron-scheduler pod is bound to a node

```
kubectl logs -n kube-system my-scheduler-79bd4cb788-hq2sq
```

```
I1012 15:30:21.629611      1 scheduler.go:604] "Successfully bound pod to
↳ node" pod="kube-system/k8s-neuron-scheduler-5d9d9d7988-xcpqm" node="ip-
↳ 192-168-2-25.ec2.internal" evaluatedNodes=1 feasibleNodes=1
```

- When running new pod's that need to use the neuron scheduler extension, make sure it uses the my-scheduler as the scheduler. Sample pod spec is below

```
apiVersion: v1
kind: Pod
metadata:
  name: <POD_NAME>
spec:
  restartPolicy: Never
  schedulerName: my-scheduler
  containers:
    - name: <POD_NAME>
      command: ["<COMMAND>"]
      image: <IMAGE_NAME>
      resources:
        limits:
          cpu: "4"
          memory: 4Gi
          aws.amazon.com/neuroncore: 9
        requests:
          cpu: "1"
          memory: 1Gi
```

- Once the neuron workload pod is run, make sure logs in the k8s neuron scheduler has successful filter/bind request

```
kubectl logs -n kube-system k8s-neuron-scheduler-5d9d9d7988-xcpqm
```

```
2022/10/12 15:41:16 POD nrt-test-5038 fits in Node:ip-192-168-2-25.ec2.
↳ internal
2022/10/12 15:41:16 Filtered nodes: [ip-192-168-2-25.ec2.internal]
2022/10/12 15:41:16 Failed nodes: map[]
2022/10/12 15:41:16 Finished Processing Filter Request...
```

```
2022/10/12 15:41:16 Executing Bind Request!
2022/10/12 15:41:16 Determine if the pod %v is NeuronDevice podnrt-test-5038
2022/10/12 15:41:16 Updating POD Annotation with alloc devices!
2022/10/12 15:41:16 Return aws.amazon.com/neuroncore
2022/10/12 15:41:16 neuronDevUsageMap for resource:aws.amazon.com/
↳ neuroncore in node: ip-192-168-2-25.ec2.internal is [false false false
↳ false false false false false false false false false false]
```

(continues on next page)

(continued from previous page)

```

↳false]
2022/10/12 15:41:16 Allocated ids for POD nrt-test-5038 are: 0,1,2,3,4,5,6,
↳7,8
2022/10/12 15:41:16 Try to bind pod nrt-test-5038 in default namespace to
↳node ip-192-168-2-25.ec2.internal with &Binding{ObjectMeta:{nrt-test-5038
↳8da590b1-30bc-4335-b7e7-fe574f4f5538 0 0001-01-01 00:00:00 +0000 UTC
↳<nil> <nil> map[] map[] [] [] []},Target:ObjectReference{Kind:Node,
↳Namespace:,Name:ip-192-168-2-25.ec2.internal,UID:,APIVersion:,
↳ResourceVersion:,FieldPath:,},}
2022/10/12 15:41:16 Updating the DevUsageMap since the bind is successful!
2022/10/12 15:41:16 Return aws.amazon.com/neuroncore
2022/10/12 15:41:16 neuronDevUsageMap for resource:aws.amazon.com/
↳neuroncore in node: ip-192-168-2-25.ec2.internal is [false false false
↳false false false false false false false false false
↳false]
2022/10/12 15:41:16 neuronDevUsageMap for resource:aws.amazon.com/
↳neurondevice in node: ip-192-168-2-25.ec2.internal is [false false false
↳false]
2022/10/12 15:41:16 Allocated devices list 0,1,2,3,4,5,6,7,8 for resource
↳aws.amazon.com/neuroncore
2022/10/12 15:41:16 Allocated devices list [0] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [0] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [0] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [0] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [1] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [1] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [1] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [1] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Allocated devices list [2] for other resource aws.
↳amazon.com/neurondevice
2022/10/12 15:41:16 Return aws.amazon.com/neuroncore
2022/10/12 15:41:16 Successfully updated the DevUsageMap [true true true
↳true true true true true true false false false false false]
↳and otherDevUsageMap [true true true false] after alloc for node ip-192-
↳168-2-25.ec2.internal
2022/10/12 15:41:16 Finished executing Bind Request...

```



## Default Scheduler Approach

- Make sure Neuron device plugin is running
- Enable the kube-scheduler with option to use configMap for scheduler policy. In your cluster.yml Please update the spec section with the following

```
spec:
  kubeScheduler:
    usePolicyConfigMap: true
```

- Launch the cluster

```
kops create -f cluster.yml
kops create secret --name neuron-test-1.k8s.local sshpublickey admin -i ~/.
ssh/id_rsa.pub
kops update cluster --name neuron-test-1.k8s.local --yes
```

- Install the neuron-scheduler-extension [Registers neuron-scheduler-extension with kube-scheduler]

```
helm upgrade --install neuron-helm-chart oci://public.ecr.aws/neuron/neuron-
helm-chart \
  --set "scheduler.enabled=true" \
  --set "scheduler.customScheduler.enabled=false" \
  --set "scheduler.defaultScheduler.enabled=true" \
  --set "npd.enabled=false"
```

## Neuron Node Problem Detector Plugin

The Neuron Problem Detector Plugin facilitates error detection and recovery by continuously monitoring the health of Neuron devices across all Kubernetes nodes. It publishes CloudWatch metrics for node errors and can optionally trigger automatic recovery of affected nodes. Please follow the instructions below to enable the necessary permissions for the plugin.

### Permissions for Neuron Node Problem Detector Plugin

Neuron node problem detection and recovery is authorized via IAM roles for service accounts. For more information, see [IAM roles for service accounts](#) in the Amazon EKS User Guide. This documentation shows how to configure an IAM role for service accounts using the command line tool eksctl. Follow the instructions below to configure IAM authorization for service accounts:

- Install the eksctl CLI using instructions listed at <https://eksctl.io/installation/>.
- Create a policy as shown below:

Policy template

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Action": [
        "autoscaling:SetInstanceHealth",
```

(continues on next page)

(continued from previous page)

```

        "autoscaling:DescribeAutoScalingInstances"
    ],
    "Effect": "Allow",
    "Resource": <arn of the Auto Scaling group corresponding to the_
↪Neuron nodes for the cluster>
    },
    {
        "Action": [
            "ec2:DescribeInstances"
        ],
        "Effect": "Allow",
        "Resource": "*",
        "Condition": {
            "ForAllValues:StringEquals": {
                "ec2:ResourceTag/aws:autoscaling:groupName": <name of_
↪the Auto Scaling group corresponding to the Neuron nodes for the cluster>
            }
        }
    },
    {
        "Action": [
            "cloudwatch:PutMetricData"
        ],
        "Effect": "Allow",
        "Resource": "*",
        "Condition": {
            "StringEquals": {
                "cloudwatch:Namespace": "NeuronHealthCheck"
            }
        }
    }
}
]
}

```

To create the policy, the AWS CLI can be used as shown below, where npd-policy-trimmed.json is the JSON policy constructed from the template above.

```

aws iam create-policy \
  --policy-name NeuronProblemDetectorPolicy \
  --policy-document file://npd-policy-trimmed.json

```

- Create a namespace for the Neuron Node Problem Detector and its service account:

```
kubectl create ns neuron-healthcheck-system
```

- Associate the authorization with the service account using the following script:

```

#!/bin/bash
CLUSTER_NAME=<eks cluster name>
REGION_CODE=$(aws configure get region)
POLICY_ARN=<policy arn for NeuronProblemDetectorPolicy>

eksctl create iamserviceaccount \

```

(continues on next page)

(continued from previous page)

```
--name node-problem-detector \
--namespace neuron-healthcheck-system \
--cluster $CLUSTER_NAME \
--attach-policy-arn $POLICY_ARN \
--approve \
--role-name neuron-problem-detector-role-$CLUSTER_NAME \
--region $REGION_CODE \
--override-existing-serviceaccounts
```

- Verify that the service account is annotated correctly. An example is shown below:

```
kubectl describe sa node-problem-detector -n neuron-healthcheck-system
Name:                node-problem-detector
Namespace:           neuron-healthcheck-system
Labels:              app.kubernetes.io/managed-by=eksctl
Annotations:         eks.amazonaws.com/role-arn: arn:aws:iam::111111111111:
                    ↪role/neuron-problem-detector-role-cluster1
Image pull secrets:  <none>
Mountable secrets:   <none>
Tokens:              <none>
Events:              <none>
```

- To cleanup, deletion of the service account can be done using the following command:

```
#!/bin/bash
CLUSTER_NAME=<eks cluster name>
REGION_CODE=$(aws configure get region)

eksctl delete iamserviceaccount \
--name node-problem-detector \
--namespace neuron-healthcheck-system \
--cluster $CLUSTER_NAME \
--approve \
--region $REGION_CODE \
```

## Deploy Neuron Node Problem Detector And Recovery

Neuron node problem detector and recovery artifact checks the health of Neuron devices on each Kubernetes node. After detecting an unrecoverable Neuron error, it triggers a node replacement. In order to get started with Neuron node problem detector and recovery, make sure that the following requirements are satisfied:

- The Neuron node problem detector and recovery requires Neuron driver 2.15+, and it requires the runtime to be at SDK 2.18 or later.
- Make sure prerequisites are satisfied. This includes prerequisites for getting started with Kubernetes containers and prerequisites for the Neuron node problem detector and recovery.
- Install the Neuron node problem detector and recovery as a DaemonSet on the cluster with the following command:

---

**Note:** The installation pulls the container image from the upstream repository for node problem detector `registry.k8s.io/node-problem-detector`.

---

```
helm upgrade --install neuron-helm-chart oci://public.ecr.aws/neuron/neuron-
↪helm-chart
```

- By default, the Neuron node problem detector and recovery has monitor only mode enabled. To enable the recovery functionality:

```
helm upgrade --install neuron-helm-chart oci://public.ecr.aws/neuron/neuron-
↪helm-chart \
  --set "npd.nodeRecovery.enabled=true"
```

- Verify that the Neuron device plugin is running:

```
kubect1 get pod -n neuron-healthcheck-system
```

Expected result (with 4 nodes in cluster):

| NAME                        | READY | STATUS  | RESTARTS | AGE |
|-----------------------------|-------|---------|----------|-----|
| node-problem-detector-7qcrj | 1/1   | Running | 0        | 59s |
| node-problem-detector-j45t5 | 1/1   | Running | 0        | 59s |
| node-problem-detector-mr2cl | 1/1   | Running | 0        | 59s |
| node-problem-detector-vpjtk | 1/1   | Running | 0        | 59s |

- When any unrecoverable error occurs, Neuron node problem detector and recovery publishes a metric under the CloudWatch namespace NeuronHealthCheck. It also reflects in NodeCondition and can be seen with `kubect1 describe node`.

## Neuron Monitor Daemonset

Neuron monitor is primary observability tool for neuron devices. For details of neuron monitor, please refer to the [neuron monitor guide](#). This tutorial describes deploying neuron monitor as a daemonset on the kubernetes cluster.

### Deploy Neuron Monitor Daemonset

- Download the neuron monitor yaml file. `k8s-neuron-monitor-daemonset.yml`
- Apply the Neuron monitor yaml to create a daemonset on the cluster with the following command

```
kubect1 apply -f k8s-neuron-monitor.yml
```

- Verify that neuron monitor daemonset is running

```
kubect1 get ds neuron-monitor --namespace neuron-monitor
```

Expected result (with 2 nodes in cluster):

| NAME           | DESIRED | CURRENT | READY | UP-TO-DATE |  |
|----------------|---------|---------|-------|------------|--|
| ↪AVAILABLE     |         |         |       |            |  |
| node-selector  | AGE     |         |       |            |  |
| neuron-monitor | 2       | 2       | 2     | 2          |  |
| ↪2             | <none>  | 27h     |       |            |  |

- Get the neuron-monitor pod names

```
kubectl get pods
```

Expected result

| NAME                 | READY | STATUS  | RESTARTS | AGE |
|----------------------|-------|---------|----------|-----|
| neuron-monitor-slsxf | 1/1   | Running | 0        | 17m |
| neuron-monitor-wc4f5 | 1/1   | Running | 0        | 17m |

- Verify the prometheus endpoint is available

```
kubectl exec neuron-monitor-wc4f5 -- wget -q --output-document - http://127.0.0.1:8000
```

Expected result

```
# HELP python_gc_objects_collected_total Objects collected during gc
# TYPE python_gc_objects_collected_total counter
python_gc_objects_collected_total{generation="0"} 362.0
python_gc_objects_collected_total{generation="1"} 0.0
python_gc_objects_collected_total{generation="2"} 0.0
# HELP python_gc_objects_uncollectable_total Uncollectable objects found during GC
# TYPE python_gc_objects_uncollectable_total counter
```

## Neuron Helm Chart

To simplify the Kubernetes container deployment process, the Neuron Helm Chart has been provided with the following containers:

- Neuron Device Plugin
- Neuron Scheduler Extension
- Neuron Node Problem Detector and Recovery

For information on how to setup the containers on a Kubernetes cluster using the Neuron Helm Chart, please refer to <https://github.com/aws-neuron/neuron-helm-charts/>.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

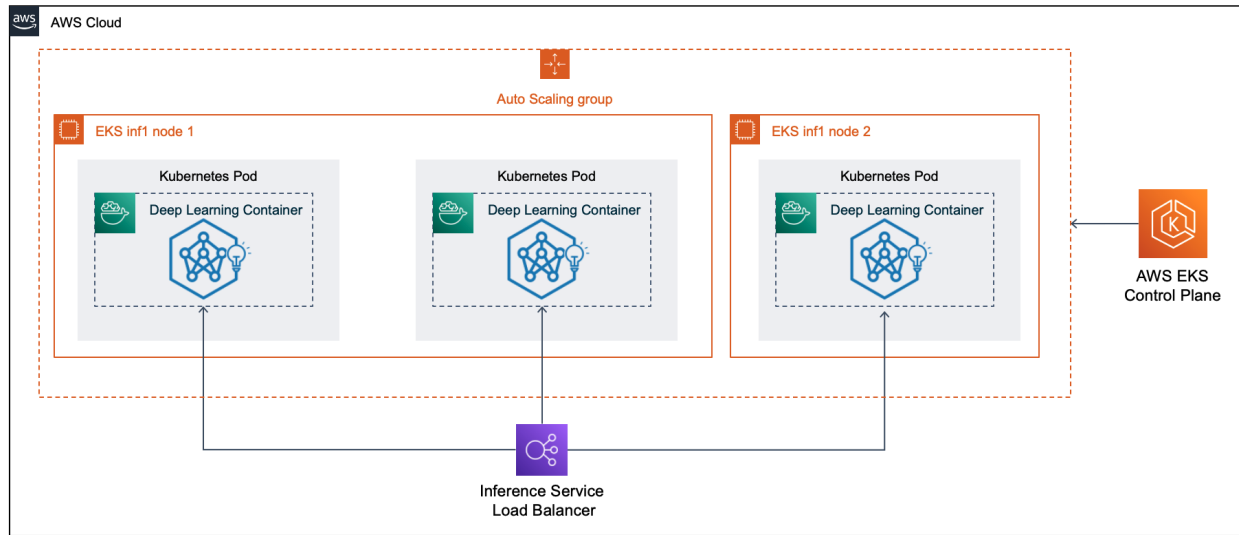
*This document is relevant for:* Inf1

## Deploy Neuron Container on Elastic Kubernetes Service (EKS) for Inference

### Table of Contents

- *Description*
- *Setup Environment*
- *Inference Example*

## Description



You can use the Neuron version of the [AWS Deep Learning Containers](#) to run inference on Amazon Elastic Kubernetes Service (EKS). In this developer flow, you set up an EKS cluster with Inf1 instances, create a Kubernetes manifest for your inference service and deploy it to your cluster. This developer flow assumes:

1. The model has already been compiled through [Compilation with Framework API on EC2 instance](#) or through [Compilation with Sagemaker Neo](#).
2. You already set up your container to retrieve it from storage.

## Setup Environment

Please add inferentia nodes using instructions at [EKS Setup For Neuron](#).

Using the YAML deployment manifest shown in the [EKS documentation for inferentia](#), replace the *image* in the *containers* specification with the one you built using [how-to-build-neuron-container](#).

---

**Note:** Before deploying the yaml to your EKS cluster, make sure to push the image to ECR. Refer to [Pushing a Docker image](#) for more information.

---

## Inference Example

Please refer to [example-deploy-rn50-as-k8s-service](#) run a simple inference example. Note that the container image referenced in the YAML manifest is created using [how-to-build-neuron-container](#).

*This document is relevant for:* Inf1

*This document is relevant for:* Trn1, Trn2

## Deploy a simple mlp training script as a Kubernetes job

This tutorial uses mlp train as a teaching example on how to deploy an training application using Kubernetes on the Trn1 instances. For more advanced example, please refer to [Tutorial: Launch a Multi-Node PyTorch Neuron Training Job on Trainium Using TorchX and EKS](#)

### Prerequisite:

- *EKS Setup For Neuron*: to setup k8s support on your cluster.
- Trn1 instances as worker nodes with attached roles allowing:
  - ECR read access policy to retrieve container images from ECR:  
**arn:aws:iam::aws:policy/AmazonEC2ContainerRegistryReadOnly**
- Have a container image that is build using tutorial-training

### Deploy a mlp training image

1. Create a file named *mlp\_train.yaml* with the contents below.

---

**Note:** In the image: add the appropriate location of the image

---

```
apiVersion: v1
kind: Pod
metadata:
  name: trn1-mlp
spec:
  restartPolicy: Never
  schedulerName: default-scheduler
  hostNetwork: true
  nodeSelector:
    beta.kubernetes.io/instance-type: trn1.32xlarge
    beta.kubernetes.io/instance-type: trn1.2xlarge
  containers:
    - name: trn1-mlp
      command: ["/usr/local/bin/python3"]
      args: ["/opt/ml/mlp_train.py"]
      image: 647554078242.dkr.ecr.us-east-1.amazonaws.com/sunda-pt:k8s_mlp_0907
      imagePullPolicy: IfNotPresent
      env:
        - name: NEURON_RT_LOG_LEVEL
          value: "INFO"
      resources:
        limits:
          aws.amazon.com/neuron: 2
        requests:
          aws.amazon.com/neuron: 2
```

2. Deploy the pod.

```
kubectl apply -f mlp_train.yaml
```

3. Check the logs to make sure training completed

```
kubectl logs <pod name>
```

Your log should have the following

```
Final loss is 0.1977
-----End Training -----
```

*This document is relevant for: Trn1, Trn2*

In this section, you'll find resources to help you use Neuron with EKS cluster, deploying inference and training workloads on Inferentia and Trainium EKS clusters.

## EKS Setup

This guide covers setting up the Neuron device plugin, scheduler extension, node problem detector, and monitoring plugins. These components enable efficient resource utilization, monitoring, and resilience when using Inferentia and Trainium instances for inference and training workloads on Kubernetes clusters. To get started with using AWS Neuron and setting up the required plugins on an EKS cluster, please refer to [EKS Setup For Neuron](#).

## Running Inference workload

This guide walks you through the end-to-end process of building and running a Docker container with your model and deploying it on an EKS cluster with Inferentia instances. For running machine learning inference workloads on Amazon EKS using AWS Deep Learning Containers, please refer to [Deploy Neuron Container on Elastic Kubernetes Service \(EKS\) for Inference](#).

## Running Training workload

This guide walks you through the end-to-end process of building and running a Docker container with your model and deploying it on an EKS cluster with Trainium instances. For running machine learning training workloads on Amazon EKS using AWS Deep Learning Containers, please refer to [Deploy a simple mlp training script as a Kubernetes job](#).

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.3.2 Amazon ECS

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*



## Neuron Problem Detector And Recovery

Neuron node problem detector and recovery artifact checks the health of Neuron devices on each ECS instance. After detecting an unrecoverable Neuron error, it triggers an instance replacement. In order to get started with Neuron node problem detector and recovery, make sure that the following requirements are satisfied:

- The Neuron node problem detector and recovery requires Neuron driver 2.15+, and it requires the runtime to be at SDK 2.18 or later.

## Creating a Task Definition

### Configuration

The task definition includes two containers:

- **npd-container**: This container is responsible for enabling Problem detection functionality in the ECS cluster.
- **recovery-container**: This container handles recovery operations in case of failures detected by Neuron Problem Detector.

The **recovery-container** has an environment variable called `ENABLE_RECOVERY` that controls whether recovery is enabled or disabled. Set the value to `true` to enable recovery, or `false` to disable it.

Follow these steps to create a task definition for NPD and recovery:

1. Go to the [ECS console](#) and select **Task Definitions** in the navigation pane.
2. Click **Create new Task Definition** and choose **Create new Task Definition with JSON**.
3. Paste the task definition JSON provided, replacing the placeholders with your account-specific values.

```
{
  "family": "neuron-npd-and-recovery",
  "containerDefinitions": [
    {
      "name": "npd",
      "image": "registry.k8s.io/node-problem-detector/node-problem-
↪detector:v0.8.19",
      "cpu": 0,
      "portMappings": [
        {
          "name": "npd-80-tcp",
          "containerPort": 80,
          "hostPort": 80,
          "protocol": "tcp",
          "appProtocol": "http"
        }
      ],
      "essential": true,
      "entryPoint": [
        "/bin/sh",
        "-c"
      ],
      "command": [
        "echo '{\"plugin\": \"kmsg\", \"logPath\": \"/dev/kmsg\", \"
↪lookback\": \"5m\", \"bufferSize\": 10, \"source\": \"kernel-monitor\", \"
```

(continues on next page)

(continued from previous page)

```

↪ "conditions\":[{"type\":"NeuronHealth\","\reason\":"NeuronHasNoError\","
↪ \message\":"Neuronhasnoerror\"}],\rules\":[{"type\":"permanent\","
↪ "condition\":"NeuronHealth\","\reason\":"NeuronHasError_SRAM_
↪ UNCORRECTABLE_ERROR\","\pattern\":"*.NEURON_HW_ERR=SRAM_UNCORRECTABLE_
↪ ERROR.*"}, {"type\":"permanent\","\condition\":"NeuronHealth\","\
↪ reason\":"NeuronHasError_NC_UNCORRECTABLE_ERROR\","\pattern\":"*.
↪ *NEURON_HW_ERR=NC_UNCORRECTABLE_ERROR.*"}, {"type\":"permanent\","\
↪ "condition\":"NeuronHealth\","\reason\":"NeuronHasError_HBM_
↪ UNCORRECTABLE_ERROR\","\pattern\":"*.NEURON_HW_ERR=HBM_UNCORRECTABLE_
↪ ERROR.*"}, {"type\":"permanent\","\condition\":"NeuronHealth\","\
↪ reason\":"NeuronHasError_DMA_ERROR\","\pattern\":"*.NEURON_HW_ERR=DMA_
↪ ERROR.*"}]]}' > /config/kernel-monitor.json && /node-problem-detector --
↪ v=2 --logtostderr --enable-k8s-exporter=false --config.system-log-
↪ monitor=/config/kernel-monitor.json"
    ],
    "environment": [],
    "mountPoints": [],
    "volumesFrom": [],
    "linuxParameters": {
      "devices": [
        {
          "hostPath": "/dev/kmsg",
          "containerPath": "/dev/kmsg",
          "permissions": [
            "read",
            "write"
          ]
        }
      ]
    },
    "privileged": true,
    "logConfiguration": {
      "logDriver": "awslogs",
      "options": {
        "awslogs-group": "/ecs/npd",
        "awslogs-create-group": "true",
        "awslogs-region": "us-west-2",
        "awslogs-stream-prefix": "ecs"
      },
      "secretOptions": []
    },
    "systemControls": []
  },
  {
    "name": "recovery",
    "image": "public.ecr.aws/neuron/neuron-node-recovery:1.3.0",
    "cpu": 0,
    "portMappings": [],
    "essential": true,
    "entryPoint": [
      "/bin/sh",
      "-c"
    ]
  }
}

```

(continues on next page)

(continued from previous page)

```

    ],
    "command": [
        "python scripts/check-health.py"
    ],
    "environment": [
        {
            "name": "ENABLE_RECOVERY",
            "value": "false"
        }
    ],
    "mountPoints": [],
    "volumesFrom": [],
    "readonlyRootFilesystem": true,
    "logConfiguration": {
        "logDriver": "awslogs",
        "options": {
            "awslogs-create-group": "true",
            "awslogs-group": "/ecs/recovery",
            "awslogs-region": "us-west-2",
            "awslogs-stream-prefix": "ecs"
        }
    },
    "systemControls": []
  },
  "executionRoleArn": "arn:aws:iam::012345678910:role/ecsTaskExecutionRole",
  "taskRoleArn": "arn:aws:iam::012345678910:role/ecsTaskExecutionRole",
  "networkMode": "awsvpc",
  "requiresCompatibilities": [
    "EC2"
  ],
  "cpu": "1024",
  "memory": "3072",
  "runtimePlatform": {
    "cpuArchitecture": "X86_64",
    "operatingSystemFamily": "LINUX"
  }
}

```

4. Review the task definition and click **Create**.

For more details on task definitions, refer to the [AWS documentation](#).

## Deploying the Service

After creating the task definition, follow these steps to deploy the service:

1. In the ECS console, select the task definition and click **Deploy** → **Create Service**.
2. Select your ECS cluster, set the launch type to **EC2**, and the service type to **Daemon**.
3. Click **Create** to deploy the service.

For more details on deploying services, refer to the [AWS documentation](#).

## Permissions

Ensure the ECS task execution role and task role have permissions to:

- Publish metrics to CloudWatch
- Read and set health status of EC2 instances in the Auto Scaling group

Refer to the [AWS documentation on IAM roles for ECS tasks](#) for more information.

When any unrecoverable error occurs, Neuron node problem detector and recovery publishes a metric under the CloudWatch namespace NeuronHealthCheck. It also reflects in NodeCondition and can be seen with `kubectl describe node`.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

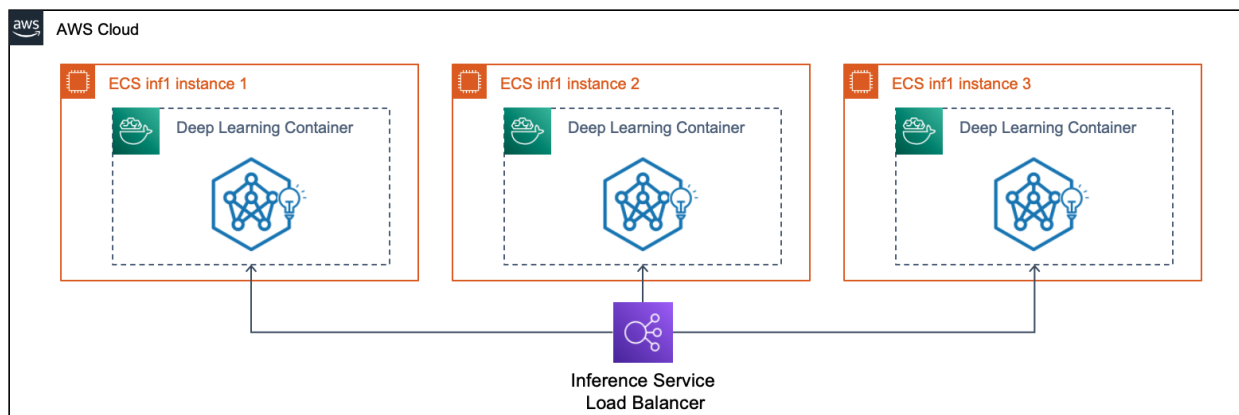
*This document is relevant for: Inf1*

## Deploy Neuron Container on Elastic Container Service (ECS) for Inference

### Table of Contents

- *Description*
- *Setup Environment*

## Description



You can use the Neuron version of the [AWS Deep Learning Containers](#) to run inference on Amazon Elastic Container Service (ECS). In this developer flow, you set up an ECS cluster with inf1/inf2 instances, create a task description for your inference service and deploy it to your cluster. This developer flow assumes:

1. The model has already been compiled through [Compilation with Framework API on EC2 instance](#) or through Compilation with Sagemaker Neo.
2. You already set up your container to retrieve it from storage.

## Setup Environment

1. **Set up an Amazon ECS cluster:**

Follow the instructions on [Setting up Amazon ECS for Deep Learning Containers](#)

2. **Define an Inference Task:**

Use the instruction on the [DLC Inference on ECS Tutorial](#) to define a task and create a service for the appropriate framework.

When creating tasks for inferentia instances on ECS, be aware of the considerations and requirements listed in [Working with inference workloads on Amazon ECS](#).

3. Use the container image created using how-to-build-neuron-container as the `image` in your task definition.

---

**Note:** Before deploying your task definition to your ECS cluster, make sure to push the image to ECR. Refer to [Pushing a Docker image](#) for more information.

---

*This document is relevant for: Inf1*

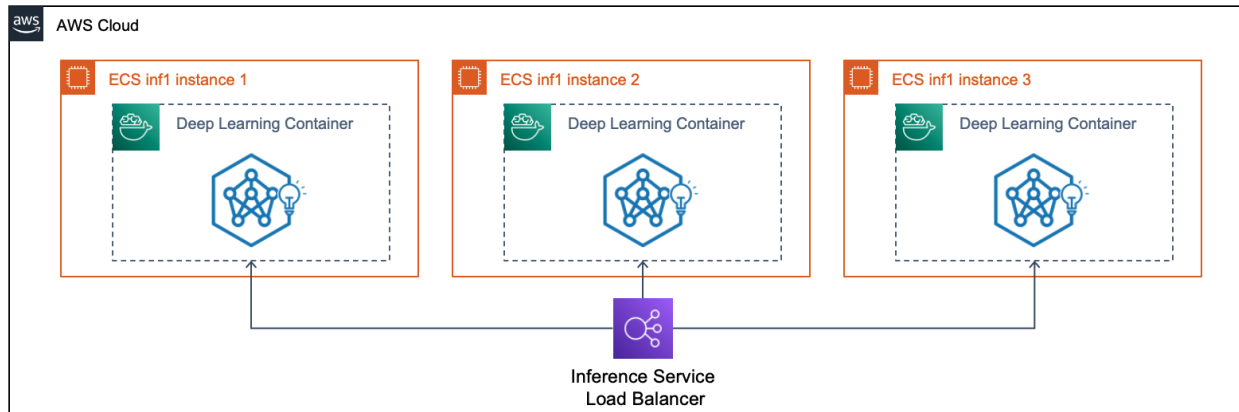
*This document is relevant for: Trn1, Trn2*

## Deploy Neuron Container on Elastic Container Service (ECS) for Training

### Table of Contents

- [Description](#)
- [Setup Environment](#)

## Description



You can use the Neuron version of the [AWS Deep Learning Containers](#) to run training on Amazon Elastic Container Service (ECS). In this developer flow, you set up an ECS cluster with trn1 instances, create a task description for your training container and deploy it to your cluster. This developer flow assumes:

1. The model has already been compiled through [Compilation with Framework API on EC2 instance](#) or through [Compilation with SageMaker Neo](#).
2. You already set up your container to retrieve it from storage.

## Setup Environment

### 1. Set up an Amazon ECS cluster:

Follow the instructions on [Setting up Amazon ECS for Deep Learning Containers](#)

### 2. Define a Training Task:

Use the instruction on the [DLC Training on ECS Tutorial](#) to define a task and create a service for the appropriate framework.

When creating tasks for trn1 instances on ECS, be aware of the considerations and requirements listed in [Working with training workloads on Amazon ECS](#).

3. Use the container image created using [how-to-build-neuron-container](#) as the `image` in your task definition.

---

**Note:** Before deploying your task definition to your ECS cluster, make sure to push the image to ECR. Refer to [Pushing a Docker image](#) for more information.

---

*This document is relevant for: Trn1, Trn2*

In this section, you'll find resources to help you use Neuron with ECS cluster, deploying inference and training workloads on Inferentia and Trainium ECS clusters.

## Using Neuron Node Problem Detector Plugin with ECS

Neuron node problem detector and recovery plugin enhances resiliency by detecting and remediating errors. To get started with using Neuron node problem detector plugin and recovery plugin on an ECS cluster, please refer to [Neuron Problem Detector And Recovery](#).

## Running Inference workload

This guide walks you through the end-to-end process of building and running a Docker container with your model and deploying it on an ECS cluster with Inferentia instances. For running machine learning inference workloads on Amazon ECS using AWS Deep Learning Containers, please refer to [Deploy Neuron Container on Elastic Container Service \(ECS\) for Inference](#).

## Running Training workload

This guide walks you through the end-to-end process of building and running a Docker container with your model and deploying it on an ECS cluster with Trainium instances. For running machine learning training workloads on Amazon ECS using AWS Deep Learning Containers, please refer to [Deploy Neuron Container on Elastic Container Service \(ECS\) for Training](#).

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.3.3 AWS ParallelCluster

*This document is relevant for: Trn1, Trn2*

### Parallel Cluster Flows- Training

*This document is relevant for: Inf2, Trn1, Trn2*

## Train your model on ParallelCluster

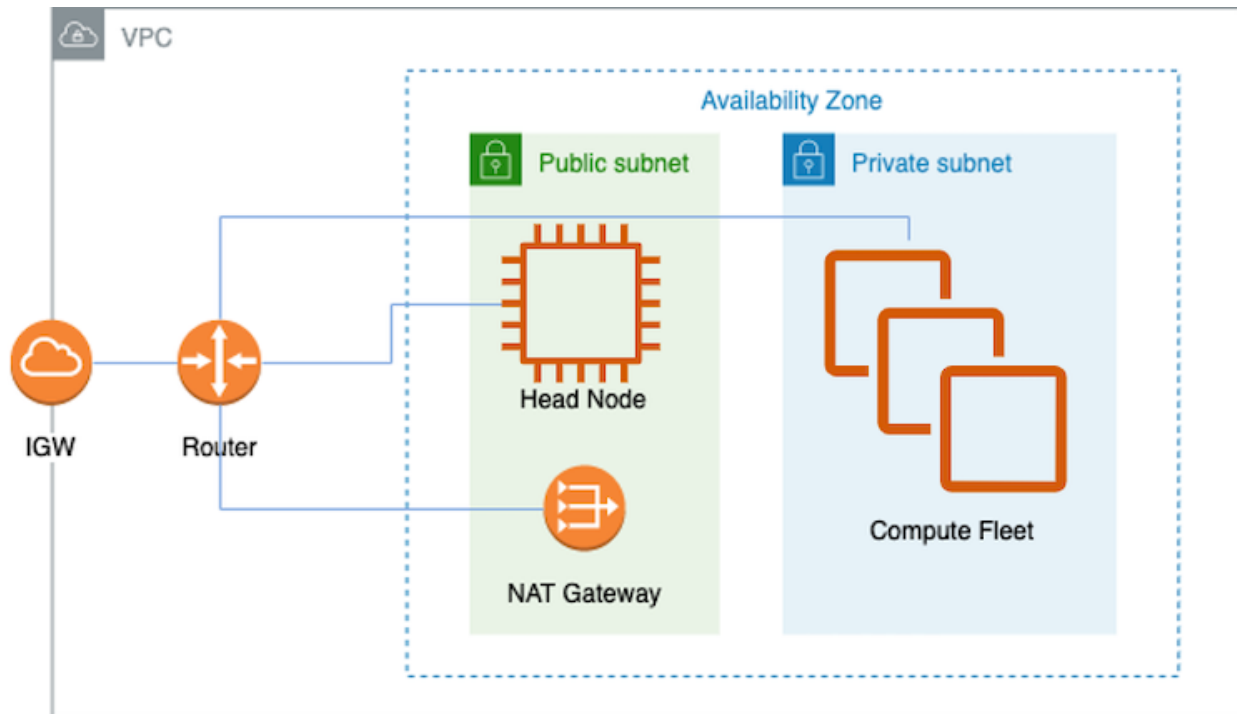
### Table of Contents

- [Description](#)
- [Setup environment](#)

## Description

This document explains how to use AWS ParallelCluster to build HPC compute environment that uses Trn1 compute nodes to run your distributed ML training job. Once the nodes are launched, we will run a training task to confirm that the nodes are working, and use slurm commands to check the job status. In this tutorial, we will use `AWS pcluster` command to run a yaml file in order to generate the cluster. As an example, we are going to launch multiple Trn1.32x1 nodes in our cluster.

We are going to set up our ParallelCluster infrastructure as below:



As shown in the figure above, inside a VPC, there are two subnets, a public and a private ones. Head Node resides in the public subnet, while the compute fleet (in this case, trn1 instances) are in the private subnet. A Network Address Translation (NAT) gateway is also needed in order for nodes in the private subnet to connect to clients outside the VPC. In the next section, we are going to describe how to set up all the necessary infrastructure for trn1 ParallelCluster.

## Setup environment

1. Install prerequisite infrastructure:

Follow [these setup](#) instructions to install VPC and all the necessary components for ParallelCluster.

2. Install AWS ParallelCluster in a virtual environment (recommended)

Follow <https://docs.aws.amazon.com/parallelcluster/latest/ug/install-v3-virtual-environment.html>

3. Create and launch ParallelCluster

Follow [these creating cluster](#) instructions to launch ParallelCluster in the VPC.

1. Launch training job

Follow [these running training](#) instructions to submit a model training script as a slurm job.

*This document is relevant for: Inf2, Trn1, Trn2*



- *Train your model on ParallelCluster*

*This document is relevant for:* Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### 5.3.4 AWS Batch

*This document is relevant for:* Inf2, Trn1, Trn2

#### Train your model on AWS Batch

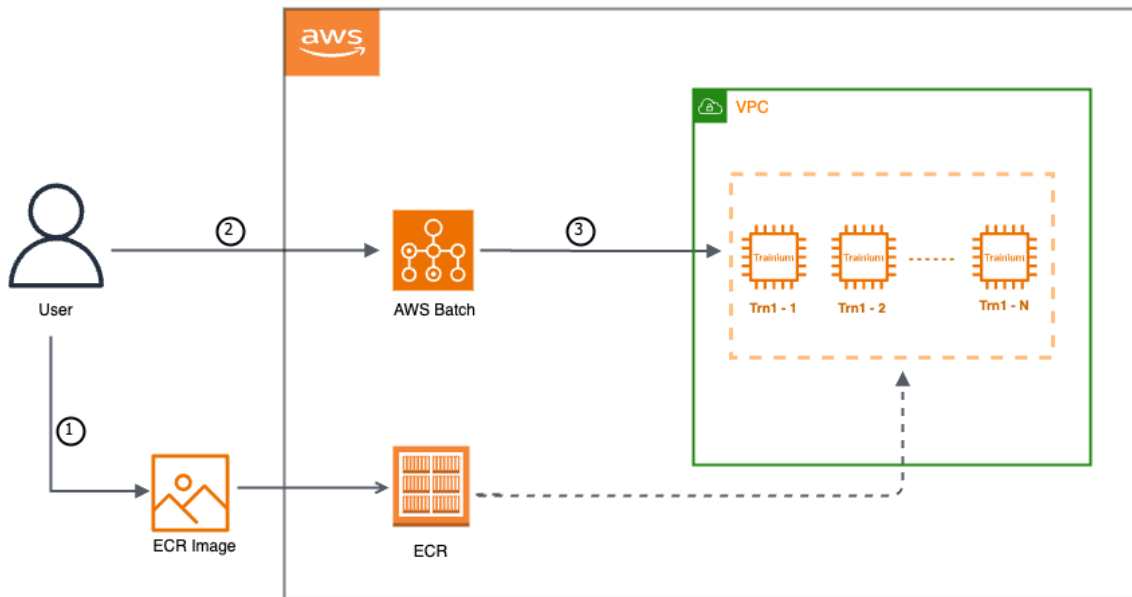
##### Table of Contents

- *Description*
- *How does AWS Batch work with Trainium*

#### Description

AWS Batch provides a scalable and cost-effective solution for running batch computing workloads in the AWS Cloud. Integrating Trainium with AWS Batch provides an efficient and cost-effective way of training deep learning models at scale. Once you configure your training job, AWS Batch effectively manages the orchestration, execution, and dynamic scaling of compute resources for your extensive machine learning workloads. To learn more about AWS Batch, see [here](#).

## How does AWS Batch work with Trainium



As depicted in the illustration above, our workflow begins by building a Docker container image for Trainium and pushing it to Amazon Elastic Container Registry (ECR). Following this, we configure our AWS Batch environment with the required capabilities, and subsequently submit the training job.

Please follow the below mentioned steps to run your training jobs on AWS Batch with Trainium.

### 1. Before you begin, please ensure that you have the following prerequisites completed:

- AWS VPC with at least one Subnet and EFA Enabled Security Group (learn more about EFA-enabled security group [here](#)). Please make sure subnet needs to be private, and the VPC needs to have a NAT gateway to allow internet connectivity for the private subnet.
- AWS ECR repository
- AWS CLI installed and configured with permissions for the above mentioned AWS resources
- Docker
- jq

### 2. Setup to start working with AWS Batch

Connect to your EC2 instance(x86\_64-based Linux instance) and clone the aws-neuron-samples repo. Once done, navigate to aws batch scripts directory.

```
cd ~/
git clone https://github.com/aws-neuron/aws-neuron-samples.git
cd ~/aws-neuron-samples/torch-neuronx/training/aws-batch/all-reduce
```

### 3. Configure resource requirements

Update the build\_configs\_and\_setup.sh with your environment variables. Once done, execute the bash script using the command ./build\_configs\_and\_setup.sh.

### 4. Build the required docker image and publish it to ECR

Run `./build_docker_image.sh` to build a Neuron Deep-Learning Container image using the latest Neuron packages and push this image to ECR.

#### 5. Prepare the AWS infrastructure required to submit the batch job

Run `./create_resources.sh` to create all AWS Batch resources needed for your training workload. Below is the brief description of various AWS Batch components this script will create for you -

- **Placement Group** enables you to influence the placement of your EC2 (Elastic Compute Cloud) instances within the AWS infrastructure.
- **Launch Template** allows you to define a set of instance configuration parameters, including the Amazon Machine Image (AMI), instance type, key pair, security groups, and other settings, in a template format.
- **Compute Environment** helps you to specify configuration that specifies the type of compute resources you want to use for your batch jobs. It includes details such as the EC2 instance types, the minimum and maximum number of instances, the VPC configuration, and other settings related to the compute environment.
- **Job Definition** is a blueprint that specifies how a batch job should be run. It encapsulates information about the job, such as the Docker image to be used, the command to execute within the container, the CPU and memory requirements, job dependencies, and other settings.
- **Job Queue** acts as a queueing mechanism for managing and scheduling the execution of batch computing workloads. By using job queues, AWS Batch provides a scalable and efficient way to process batch workloads, managing the allocation of resources and ensuring optimal use of compute capacity.

#### 6. Submit the job to AWS-Batch

Run `./submit_job.sh` to submit a basic all-reduce job in the provisioned AWS Batch environment

#### 7. Monitor the AWS-Batch job

You can use Amazon CloudWatch Logs to monitor, store, and view all your logs from AWS Batch job. To learn more about it, please see [here](#).

---

**Note:**

- You could run a full model training job using this setup. For example, [this sample](#) runs the Llama2-7B tutorial on AWS Batch using the same setup.
- You can further tailor your `Dockerfile` to include any additional dependencies specific to your needs.
- You have the option to leverage `trn1n.32xlarge` instances as an alternative to `trn1.32xlarge`. To make this transition, you only need to make adjustments to the `launch template` and `job definition` in order to accommodate the use of 16 EFA (Elastic Fabric Adapter) devices, whereas the current setup for `trn1` employs 8 EFA devices. Please check out [this document](#) to start with `trn1n.32xlarge` for multi-node execution.

---

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 5.4 Amazon SageMaker

Amazon SageMaker is a fully managed machine learning (ML) platform that streamlines the end-to-end ML workflow at scale. AWS Neuron integrates with Amazon SageMaker to provide optimized performance for ML workloads on AWS Inferentia and AWS Trainium chips.

### Table of contents

- [SageMaker JumpStart](#)
- [SageMaker HyperPod](#)
- [SageMaker Training](#)
- [SageMaker Inference](#)

### 5.4.1 SageMaker JumpStart

Use [Amazon SageMaker JumpStart](#) to train and deploy models using Neuron. SageMaker JumpStart is an ML hub that accelerates model selection and deployment. It provides support for fine-tuning and deploying popular models such as Meta's Llama family of models. Users can customize pre-trained models with their data and easily deploy them.

### 5.4.2 SageMaker HyperPod

Use [Amazon SageMaker HyperPod](#) to streamline ML infrastructure setup and optimization with AWS Neuron. SageMaker HyperPod leverages pre-configured distributed training libraries to split workloads across numerous AI accelerators, enhancing model performance. HyperPod ensures uninterrupted training through automatic checkpointing, fault detection, and recovery.

### 5.4.3 SageMaker Training

[Amazon SageMaker Model Training](#) reduces the time and cost to train and tune ML models at scale without the need to manage infrastructure.

### 5.4.4 SageMaker Inference

With [Amazon SageMaker](#), you can start getting predictions, or inferences, from your trained ML models. SageMaker provides a broad selection of ML infrastructure and model deployment options to help meet all your ML inference needs.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 5.5 Third-party solutions

AWS Neuron integrates with multiple third-party partner solutions that allow you to run deep learning workloads on Amazon EC2 instances powered by AWS Trainium and AWS Inferentia chips. The following list gives an overview of third-party solutions that work with AWS Neuron.

### Table of contents

- [Ray](#)
- [Domino](#)

### 5.5.1 Ray

Ray, by Anyscale, is the open source AI Compute Engine at the center of the world's most powerful AI Platforms. It precisely orchestrates infrastructure for any distributed AI workload like data processing, model training, and serving on any accelerator at any scale. Ray simplifies the complexity of distributed computing, improves efficiency, lower costs, and accelerates developer productivity.

[Ray Train documentation](#)

### 5.5.2 Domino

Domino is an open enterprise platform for data science, machine learning, and AI research. It works with an expansive list of industry leading tools and technologies to enrich data science research, development, and deployment processes. Domino works with a wide range of data sources, languages, IDEs, tools, libraries, and publication targets.

[Domino documentation](#)

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 5.6 Setup Guide

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### 5.6.1 Amazon EC2

*This document is relevant for:* Inf1

## EC2 Flows - Inference

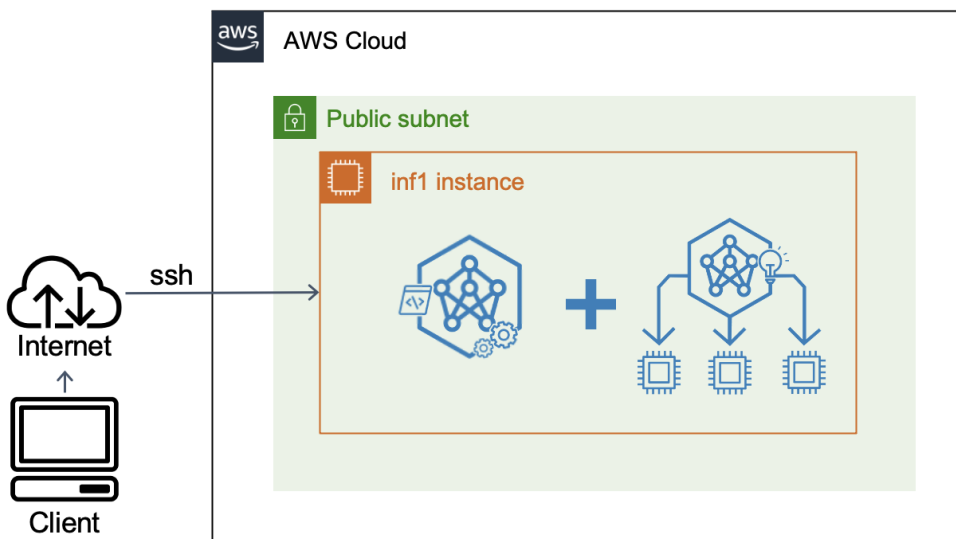
This document is relevant for: `Inf1`

### Compile with Framework API and Deploy on EC2 Inf1

#### Table of Contents

- *Description*
- *Setup Environment*
  - 1. *Launch an Inf1 Instance*
  - 2. *Set up a development environment*
    - \* *Enable PyTorch-Neuron*
    - \* *Enable TensorFlow-Neuron*
    - \* *Enable Apache MXNet*
  - 3. *Set up Jupyter notebook*

#### Description



You can use a single inf1 instance as a development environment to compile and deploy Neuron models. In this developer flow, you provision an EC2 inf1 instance using a Deep Learning AMI (DLAMI) and execute the two steps of the development flow in the same instance. The DLAMI comes pre-packaged with the Neuron frameworks, compiler, and required runtimes to complete the flow. Development happens through Jupyter Notebooks or using a secure shell (ssh) connection in terminal. Follow the steps below to setup your environment.

**Note:** Model compilation can be executed on a non-inf1 instance for later deployment. Follow the same EC2 Developer Flow Setup using other instance families and leverage [Amazon Simple Storage Service \(S3\)](#) to share the

compiled models between different instances.

---

## Setup Environment

### 1. Launch an Inf1 Instance

- Please follow the instructions at [launch an Amazon EC2 Instance](#) to Launch an Inf1 instance, when choosing the instance type at the EC2 console. Please make sure to select the correct instance type. To get more information about Inf1 instances sizes and pricing see [Inf1 web page](#).
- When choosing an Amazon Machine Image (AMI) make sure to select [Deep Learning AMI with Conda Options](#). Please note that Neuron Conda environments are supported only in Ubuntu 18 DLAMI and Amazon Linux2 DLAMI, Neuron Conda environments are not supported in Amazon Linux DLAMI.
- After launching the instance, follow the instructions in [Connect to your instance](#) to connect to the instance

---

**Note:** You can also launch the instance from AWS CLI, please see [AWS CLI commands to launch inf1 instances](#).

---

### 2. Set up a development environment

#### Enable PyTorch-Neuron

---

##### Important:

For successful installation or update to next releases (Neuron 1.20.0 and newer):

- Uninstall `aws-neuron-dkms` by running: `sudo apt remove aws-neuron-dkms` or `sudo yum remove aws-neuron-dkms`
  - Install or upgrade to latest Neuron driver (`aws-neuron-dkms`) by following the “Setup Guide” instructions.
- 

#### PyTorch 1.9.1

##### Ubuntu DLAMI

---

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

---

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↪ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↪ Neuron documentation
```

(continues on next page)

(continued from previous page)

```
#####
↪#####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↪to latest Neuron driver
#####
↪#####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↪include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↪Neuron documentation

#####
↪#####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↪to latest Neuron driver
#####
↪#####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
```

(continues on next page)



(continued from previous page)

```
# Warning: If Linux kernel is updated as a result of OS package update
#           Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate
```

## PyTorch 1.8.1

### Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#           Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate aws_neuron_pytorch_p36
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####

# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#           Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate aws_neuron_pytorch_p36
```

## Enable TensorFlow-Neuron

### Important:

For successful installation or update to next releases (Neuron 1.20.0 and newer):

- Uninstall aws-neuron-dkms by running: `sudo apt remove aws-neuron-dkms` or `sudo yum remove aws-neuron-dkms`
- Install or upgrade to latest Neuron driver (aws-neuron-dkms) by following the “Setup Guide” instructions.

## TensorFlow 2.5.1

### Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
Traceback (most recent call last):
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 1001, in <module>
    setup_cmd += nr_setup.instructions(framework=framework,action=action,framework_
↪version=args.framework_version,os=args.os,ami=args.ami,mode=args.mode)
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 977, in instructions
    setup_cmd=hlpr_instructions(self,self.neuron_version)
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 787, in hlpr_instructions
    fal_supported_rtd=nr_setup.fal_supported_runtime[fw][fw_ver]['neuron-rtd']
KeyError: '2.10.1'
```

### Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
Traceback (most recent call last):
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 1001, in <module>
    setup_cmd += nr_setup.instructions(framework=framework,action=action,framework_
↪version=args.framework_version,os=args.os,ami=args.ami,mode=args.mode)
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 977, in instructions
    setup_cmd=hlpr_instructions(self,self.neuron_version)
  File "/home/docs/checkouts/readthedocs.org/user_builds/awsdocs-neuron/checkouts/latest/
↪src/helperscripts/neuronsetuphelper.py", line 787, in hlpr_instructions
    fal_supported_rtd=nr_setup.fal_supported_runtime[fw][fw_ver]['neuron-rtd']
KeyError: '2.10.1'
```

## TensorFlow 1.15.5

### Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate TensorFlow
source activate aws_neuron_tensorflow_p36
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
```

(continues on next page)

(continued from previous page)

```

sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#         Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate TensorFlow
source activate aws_neuron_tensorflow_p36

```

## Enable Apache MXNet

### Important:

For successful installation or update to next releases (Neuron 1.20.0 and newer):

- Uninstall aws-neuron-dkms by running: `sudo apt remove aws-neuron-dkms` or `sudo yum remove aws-neuron-dkms`
- Install or upgrade to latest Neuron driver (aws-neuron-dkms) by following the “Setup Guide” instructions.

## MXNet 1.8.0

### Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```

# Note: There is no DLAMI Conda environment for this framework version
#       Framework will be installed/updated inside a Python environment

# Update OS packages
sudo apt-get update -y

#####
↪ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade ↪
↪ to latest Neuron driver
#####
↪ #####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

```

(continues on next page)

(continued from previous page)

```
# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Install Neuron Tools
sudo apt-get install aws-neuronx-tools -y

export PATH=/opt/aws/neuron/bin:$PATH

# Activate MXNet
source activate aws_neuron_mxnet_p36

# Set Pip repository to point to the Neuron repository
pip config set global.extra-index-url https://pip.repos.neuron.amazonaws.com

#Install Neuron MXNet
wget https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws_mx-1.8.0.2-py2.py3-none-
↳manylinux2014_x86_64.whl
pip install aws_mx-1.8.0.2-py2.py3-none-manylinux2014_x86_64.whl
pip install mx_neuron neuron-cc
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Note: There is no DLAMI Conda environment for this framework version
#       Framework will be installed/updated inside a Python environment

# Update OS packages
sudo yum update -y

#####
↳#####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳to latest Neuron driver
#####
↳#####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum versionlock delete aws-neuronx-dkms
```

(continues on next page)

(continued from previous page)

```

sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Install Neuron Tools
sudo yum install aws-neuronx-tools -y

export PATH=/opt/aws/neuron/bin:$PATH

# Activate MXNet
source activate aws_neuron_mxnet_p36

# Set Pip repository to point to the Neuron repository
pip config set global.extra-index-url https://pip.repos.neuron.amazonaws.com

#Install Neuron MXNet
wget https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws_mx-1.8.0.2-py2.py3-none-
↳manylinux2014_x86_64.whl
pip install aws_mx-1.8.0.2-py2.py3-none-manylinux2014_x86_64.whl
pip install mx_neuron neuron-cc

```

## MXNet 1.5.1

## Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```

# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳Neuron documentation

#####
↳#####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳to latest Neuron driver
#####
↳#####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

```

(continues on next page)

(continued from previous page)

```
#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate MXNet
source activate aws_neuron_mxnet_p36
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate MXNet
source activate aws_neuron_mxnet_p36
```



### 3. Set up Jupyter notebook

To develop from a Jupyter notebook see [setup-jupyter-notebook-steps-troubleshooting](#)

You can also run a Jupyter notebook as a script, first enable the ML framework Conda or Python environment of your choice and see [running-jupyter-notebook-as-script](#) for instructions.

*This document is relevant for: Inf1*

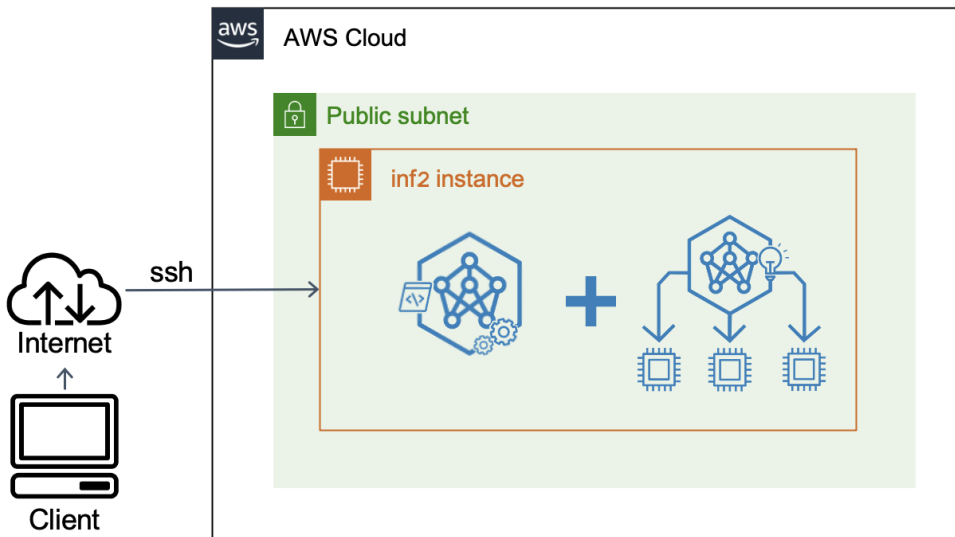
*This document is relevant for: Inf1*

### Compile with Framework API and Deploy on EC2 Inf2

#### Table of Contents

- *Description*
- *Setup Environment*
  - 1. *Launch an Inf2 Instance*
  - 2. *Set up a development environment*
    - \* *Enable PyTorch-Neuron*
  - 3. *Set up Jupyter notebook*

#### Description



You can use a single inf2 instance as a development environment to compile and deploy Neuron models. In this developer flow, you provision an EC2 inf2 instance using a Deep Learning AMI (DLAMI) and execute the two steps of the development flow in the same instance. The DLAMI comes pre-packaged with the Neuron frameworks, compiler, and required runtimes to complete the flow. Development happens through Jupyter Notebooks or using a secure shell (ssh) connection in terminal. Follow the steps below to setup your environment.

---

**Note:** Model compilation can be executed on a non-inf2 instance for later deployment. Follow the same EC2 Developer Flow Setup using other instance families and leverage [Amazon Simple Storage Service \(S3\)](#) to share the compiled models between different instances.

---

## Setup Environment

### 1. Launch an Inf2 Instance

- Please follow the instructions at [launch an Amazon EC2 Instance](#) to launch an Inf2 instance, when choosing the instance type at the EC2 console. Please make sure to select the correct instance type. To get more information about Inf2 instances sizes and pricing see [Inf2 web page](#).
- When choosing an Amazon Machine Image (AMI) make sure to select [Deep Learning AMI with Conda Options](#). Please note that Neuron Conda environments are supported only in Ubuntu 18 DLAMI and Amazon Linux2 DLAMI, Neuron Conda environments are not supported in Amazon Linux DLAMI.
- After launching the instance, follow the instructions in [Connect to your instance](#) to connect to the instance

---

**Note:** You can also launch the instance from AWS CLI, please see [AWS CLI commands to launch inf2 instances](#).

---

### 2. Set up a development environment

#### Enable PyTorch-Neuron

---

##### Important:

For successful installation or update to next releases (Neuron 1.20.0 and newer):

- Uninstall `aws-neuron-dkms` by running: `sudo apt remove aws-neuron-dkms` or `sudo yum remove aws-neuron-dkms`
  - Install or upgrade to latest Neuron driver (`aws-neuron-dkms`) by following the “Setup Guide” instructions.
- 

#### PyTorch 1.9.1

##### Ubuntu DLAMI

---

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

---

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate
```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
```

(continues on next page)

(continued from previous page)

```

sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate

```

## PyTorch 1.8.1

### Ubuntu DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```

# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####
# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Install Neuron Driver
sudo apt-get install aws-neuronx-dkms --allow-change-held-packages -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#       Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate aws_neuron_pytorch_p36

```

## Amazon Linux DLAMI

**Note:** For a successful installation or update, execute each line of the instructions below separately or copy the contents of the code block into a script file and source its contents.

```
# Neuron is pre-installed on Deep Learning AMI (DLAMI), latest DLAMI version may not
↳ include latest Neuron versions
# To update to latest Neuron version, follow "Update to latest release" instruction on
↳ Neuron documentation

#####
↳ #####

# To install or update to Neuron versions 1.19.1 and newer from previous releases:
# - DO NOT skip 'aws-neuron-dkms' install or upgrade step, you MUST install or upgrade
↳ to latest Neuron driver
#####
↳ #####

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Install Neuron Driver
sudo yum versionlock delete aws-neuronx-dkms
sudo yum install aws-neuronx-dkms -y

#####
# Warning: If Linux kernel is updated as a result of OS package update
#           Neuron driver (aws-neuron-dkms) should be re-installed after reboot
#####

# Activate PyTorch
source activate aws_neuron_pytorch_p36
```

### 3. Set up Jupyter notebook

To develop from a Jupyter notebook see [setup-jupyter-notebook-steps-troubleshooting](#)

You can also run a Jupyter notebook as a script, first enable the ML framework Conda or Python environment of your choice and see [running-jupyter-notebook-as-script](#) for instructions.

*This document is relevant for:* Inf1

- [Compile with Framework API and Deploy on EC2 Inf1](#)
- [Compile with Framework API and Deploy on EC2 Inf2](#)

*This document is relevant for:* Inf1

*This document is relevant for:* Trn1, Trn2

## EC2 Flows- Training

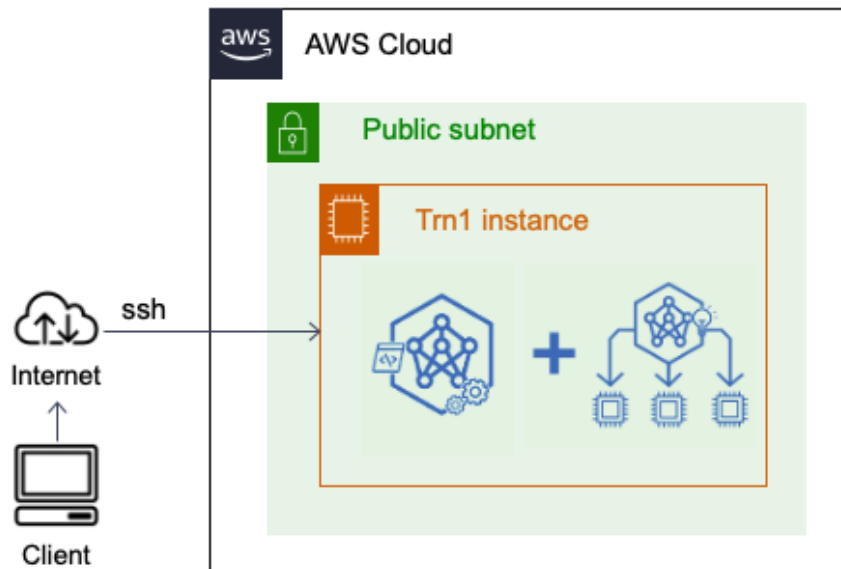
*This document is relevant for: Inf2, Trn1, Trn2*

### Train your model on EC2

#### Table of Contents

- *Description*
- *Setup Environment*
  - 1. *Launch an Trn1 Instance*
  - 2. *Set up a development environment*
    - \* *Enable PyTorch-Neuron*
  - 3. *Set up Jupyter notebook*

#### Description



You can use a single Trn1 instance as a development environment to compile and train Neuron models. In this developer flow, you provision an EC2 Trn1 instance using a Deep Learning AMI (DLAMI) and execute the two steps of the development flow in the same instance. The DLAMI comes pre-packaged with the Neuron frameworks, compiler, and required runtimes to complete the flow. Development happens through Jupyter Notebooks or using a secure shell (ssh) connection in terminal. Follow the steps below to setup your environment.

## Setup Environment

### 1. Launch an Trn1 Instance

- Please follow the instructions at [launch an Amazon EC2 Instance](#) to Launch an Trn1 instance, when choosing the instance type at the EC2 console. Please make sure to select the correct instance type. To get more information about Trn1 instances sizes and pricing see [Trn1 web page](#).
- Select your Amazon Machine Image (AMI) of choice, please note that Neuron support Ubuntu 18 AMI or Amazon Linux 2 AMI, you can also choose Ubuntu 18 or Amazon Linux 2 Deep Learning AMI (DLAMI)
- When launching a Trn1, please adjust your primary EBS volume size to a minimum of 512GB.
- After launching the instance, follow the instructions in [Connect to your instance](#) to connect to the instance

---

**Note:** If you are facing a connectivity issue during the model loading process on a Trn1 instance with Ubuntu, that could probably be because of Ubuntu limitations with multiple interfaces. To solve this problem, please follow the steps mentioned [here](#).

Users are highly encouraged to use DLAMI to launch the instances, since DLAMIs come with the required fix.

---

### 2. Set up a development environment

#### Enable PyTorch-Neuron

##### PyTorch 1.11.0

##### Ubuntu 20 AMI

---

**Note:**

- Instructions in this page only apply to setting up Neuron components on Linux host running Ubuntu or Amazon Linux AMI.
  - When launching a Trn1/Trn2, please adjust your primary EBS volume size to a minimum of 512GB.
- 

```
# Configure Linux for Neuron repository updates
. /etc/os-release

sudo tee /etc/apt/sources.list.d/neuron.list > /dev/null <<EOF
deb https://apt.repos.neuron.amazonaws.com ${VERSION_CODENAME} main
EOF
wget -q0 - https://apt.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-
  ↪ NEURON.PUB | sudo apt-key add -

# Update OS packages
sudo apt-get update -y
```

(continues on next page)

(continued from previous page)

```

# Install git
sudo apt-get install git -y

# Install OS headers
sudo apt-get install linux-headers-$(uname -r) -y

# Remove preinstalled packages and Install Neuron Driver and Runtime
sudo apt-get remove aws-neuron-dkms -y
sudo apt-get remove aws-neuronx-dkms -y
sudo apt-get remove aws-neuronx-oci-hook -y
sudo apt-get remove aws-neuronx-runtime-lib -y
sudo apt-get remove aws-neuronx-collectives -y
sudo apt-get install aws-neuronx-dkms=2.* -y
sudo apt-get install aws-neuronx-oci-hook=2.* -y
sudo apt-get install aws-neuronx-runtime-lib=2.* -y
sudo apt-get install aws-neuronx-collectives=2.* -y

# Install EFA Driver(only required for multi-instance training)

curl -O https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz
wget https://efa-installer.amazonaws.com/aws-efa-installer.key && gpg --import
↪aws-efa-installer.key
cat aws-efa-installer.key | gpg --fingerprint
wget https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz.sig &&
↪ gpg --verify ./aws-efa-installer-latest.tar.gz.sig

tar -xvf aws-efa-installer-latest.tar.gz
cd aws-efa-installer && sudo bash efa_installer.sh --yes
cd
sudo rm -rf aws-efa-installer-latest.tar.gz aws-efa-installer

# Remove pre-installed package and Install Neuron Tools
sudo apt-get remove aws-neuron-tools -y
sudo apt-get remove aws-neuronx-tools -y
sudo apt-get install aws-neuronx-tools=2.* -y

export PATH=/opt/aws/neuron/bin:$PATH

# Install Python venv and activate Python virtual environment to install
# Neuron pip packages.
sudo apt install python3.8-venv
python3.8 -m venv aws_neuron_venv_pytorch
source aws_neuron_venv_pytorch/bin/activate
pip install -U pip

# Install wget, awscli
pip install wget
pip install awscli

# Install packages from repos
python -m pip config set global.extra-index-url "https://pip.repos.neuron.

```

(continues on next page)



(continued from previous page)

```
↪amazonaws.com"

# Install Python packages - Transformers package is needed for BERT
python -m pip install torch-neuronx=="1.11.0.1.*" "neuronx-cc==2.*"
```

## Amazon Linux 2 AMI

### Note:

- Instructions in this page only apply to setting up Neuron components on Linux host running Ubuntu or Amazon Linux AMI.
- When launching a Trn1/Trn2, please adjust your primary EBS volume size to a minimum of 512GB.

```
# Configure Linux for Neuron repository updates

sudo tee /etc/yum.repos.d/neuron.repo > /dev/null <<EOF
[neuron]
name=Neuron YUM Repository
baseurl=https://yum.repos.neuron.amazonaws.com
enabled=1
metadata_expire=0
EOF
sudo rpm --import https://yum.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-
↪AWS-NEURON.PUB

# Install OS headers
sudo yum install kernel-devel-$(uname -r) kernel-headers-$(uname -r) -y

# Update OS packages
sudo yum update -y

# Install git
sudo yum install git -y

# Remove preinstalled packages and Install Neuron Driver and Runtime
sudo yum remove aws-neuron-dkms -y
sudo yum remove aws-neuronx-dkms -y
sudo yum remove aws-neuronx-oci-hook -y
sudo yum remove aws-neuronx-runtime-lib -y
sudo yum remove aws-neuronx-collectives -y
sudo yum install aws-neuronx-dkms-2.* -y
sudo yum install aws-neuronx-oci-hook-2.* -y
sudo yum install aws-neuronx-runtime-lib-2.* -y
sudo yum install aws-neuronx-collectives-2.* -y

# Install EFA Driver(only required for multi-instance training)
curl -O https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz
wget https://efa-installer.amazonaws.com/aws-efa-installer.key && gpg --import↪
↪aws-efa-installer.key
```

(continues on next page)

(continued from previous page)

```

cat aws-efa-installer.key | gpg --fingerprint
wget https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz.sig &&
↪ gpg --verify ./aws-efa-installer-latest.tar.gz.sig
tar -xvf aws-efa-installer-latest.tar.gz
cd aws-efa-installer && sudo bash efa_installer.sh --yes
cd
sudo rm -rf aws-efa-installer-latest.tar.gz aws-efa-installer

# Remove pre-installed package and Install Neuron Tools
sudo yum remove aws-neuron-tools -y
sudo yum remove aws-neuronx-tools -y
sudo yum install aws-neuronx-tools-2.* -y

export PATH=/opt/aws/neuron/bin:$PATH

# Install Python venv and activate Python virtual environment to install
# Neuron pip packages.
python3.7 -m venv aws_neuron_venv_pytorch
source aws_neuron_venv_pytorch/bin/activate
python -m pip install -U pip

# Install wget, awscli
pip install wget
pip install awscli

# Install packages from repos
python -m pip config set global.extra-index-url "https://pip.repos.neuron.
↪amazonaws.com"

# Install Python packages - Transformers package is needed for BERT
python -m pip install torch-neuronx=="1.11.0.1.*" "neuronx-cc==2.*"

```

### 3. Set up Jupyter notebook

To develop from a Jupyter notebook see [setup-jupyter-notebook-steps-troubleshooting](#)

You can also run a Jupyter notebook as a script, first enable the ML framework Conda or Python environment of your choice and see [running-jupyter-notebook-as-script](#) for instructions.

*This document is relevant for:* Inf2, Trn1, Trn2

- *Train your model on EC2*

*This document is relevant for:* Trn1, Trn2

## Inference

- *Compile with Framework API and Deploy on EC2 Inf1*
- *Compile with Framework API and Deploy on EC2 Inf2*

## Training

- *Train your model on EC2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### 5.6.2 PyTorch Neuron (torch-neuronx) Setup

**Note:** This Setup guide is relevant for Inf2 & Trn1 / Trn1n / Trn2 instances.

#### Table of contents

- *torch-neuronx setup on Ubuntu 22*
- *torch-neuronx setup on Amazon Linux 2023 (AL2023)*
- *torch-neuronx setup on Rocky Linux 9*

#### torch-neuronx setup on Ubuntu 22

Ubuntu 22 (Neuron Multi-Framework DLAMI)

Ubuntu 22 (Ubuntu22 AMI)

#### torch-neuronx setup on Amazon Linux 2023 (AL2023)

Amazon Linux 2023 (Amazon Linux 2023 AMI)

#### torch-neuronx setup on Rocky Linux 9

Rocky Linux 9 (Rocky Linux 9 AMI)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1*

### 5.6.3 PyTorch Neuron (torch-neuron) Setup

---

**Note:** This Setup guide is relevant for Inf1 instances.

---

#### Table of contents

- [torch-neuron setup on Ubuntu 20](#)
- [torch-neuron setup on Ubuntu 22](#)
- [torch-neuron setup on Amazon Linux 2023 \(AL2023\)](#)
- [torch-neuron setup on Rocky Linux 9](#)

#### torch-neuron setup on Ubuntu 20

Ubuntu 20 (Ubuntu20 AMI)

Ubuntu 20 (DLAMI Base AMI)

Ubuntu 20 (DLAMI Pytorch AMI)

#### torch-neuron setup on Ubuntu 22

Ubuntu 22 (Neuron Multi-Framework DLAMI)

Ubuntu 22 (Ubuntu22 AMI)

#### torch-neuron setup on Amazon Linux 2023 (AL2023)

Amazon Linux 2023 (Amazon Linux 2023 AMI)

#### torch-neuron setup on Rocky Linux 9

Rocky Linux 9 (Rocky Linux 9 AMI)

*This document is relevant for: Inf1*

*This document is relevant for: Inf2, Trn1, Trn2*

### 5.6.4 Tensorflow Neuron (tensorflow-neuronx) Setup

---

**Note:** This Setup guide is relevant for Inf2 & Trn1 / Trn1n instances.

---

#### Table of contents

- [tensorflow-neuronx setup on Ubuntu 22](#)
- [tensorflow-neuronx setup on Amazon Linux 2023 \(AL2023\)](#)

**tensorflow-neuronx setup on Ubuntu 22**

Ubuntu 22 (Neuron Multi-Framework DLAMI)

Ubuntu 22 (Ubuntu22 AMI)

**tensorflow-neuronx setup on Amazon Linux 2023 (AL2023)**

Amazon Linux 2023 (Amazon Linux 2023 AMI)

*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf1***5.6.5 Tensorflow Neuron (tensorflow-neuron) Setup****Note:** This Setup guide is relevant for Inf1 instances.**Table of contents**

- [tensorflow-neuron setup on Ubuntu 20](#)
- [tensorflow-neuron setup on Ubuntu 22](#)
- [tensorflow-neuron setup on Amazon Linux 2023 \(AL2023\)](#)

**tensorflow-neuron setup on Ubuntu 20**

Ubuntu 20 (Ubuntu20 AMI)

Ubuntu 20 (DLAMI Base AMI)

**tensorflow-neuron setup on Ubuntu 22**

Ubuntu 22 (Neuron Multi-Framework DLAMI)

Ubuntu 22 (Ubuntu22 AMI)

**tensorflow-neuron setup on Amazon Linux 2023 (AL2023)**

Amazon Linux 2023 (Amazon Linux 2023 AMI)

*This document is relevant for: Inf1**This document is relevant for: Inf1***5.6.6 MxNet Neuron (mxnet-neuron) Setup****Note:** This Setup guide is relevant for Inf1 instances.**Table of contents**

- [mxnet-neuron setup on Ubuntu 20](#)
- [mxnet-neuron setup on Ubuntu 22](#)

- [mxnet-neuron setup on Amazon Linux 2023 \(AL2023\)](#)

### **mxnet-neuron setup on Ubuntu 20**

Ubuntu 20 (Ubuntu20 AMI)

Ubuntu 20 (DLAMI Base AMI)

### **mxnet-neuron setup on Ubuntu 22**

Ubuntu 22 (Ubuntu22 AMI)

### **mxnet-neuron setup on Amazon Linux 2023 (AL2023)**

Amazon Linux 2023 (Amazon Linux 2023 AMI)

*This document is relevant for: Inf1*

This section walks you through the various options to install and upgrade Neuron. You have to install Neuron on Trainium and Inferentia powered instances to enable deep-learning acceleration.

Launching Inf/Trn instances on Amazon EC2

PyTorch NeuronX (`torch-neuronx`) Setup for Inf2 & Trn1

/ Trn1n/ Trn2 Instances

PyTorch Neuron (`torch-neuron`) Setup for Inf1 Instances

JAX Setup

for Inf2 & Trn1 / Trn1n Instances

Tensorflow Neuron (`tensorflow-neuronx`) Setup for Inf2 & Trn1 /

Trn1n Instances

Tensorflow Neuron (`tensorflow-neuron`) Setup for Inf1 Instances

MxNet

Neuron (`mxnet-neuron`) Setup for Inf1 Instances

PyTorch Neuron Setup Guides for Rocky Linux 9 (Inf2

& Trn1 / Trn1n) Instances

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## RUNTIME & TOOLS

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.1 NeuronX Runtime

NeuronX runtime consists of kernel driver and C/C++ libraries which provides APIs to access Inferentia and Trainium Neuron devices. The Neuron ML frameworks plugins for TensorFlow, PyTorch and Apache MXNet use the Neuron runtime to load and run models on the NeuronCores. Neuron runtime loads compiled deep learning models, also referred to as Neuron Executable File Format (NEFF) to the Neuron devices and is optimized for high-throughput and low-latency.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### 6.1.1 API Reference Guide

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### Developer's Guide - NeuronX Runtime

##### Table of contents

- *Introduction*
- *Required Software*
- *Brief Introduction to Neuron Hardware*
  - *Neuron Device*
  - *NeuronCore*
- *The Neuron Runtime Architecture*
  - *Application Interface Layer (The libnrt API)*
  - *Monitoring and Profiling*
  - *The NEFF format and NEFF Parser*
  - *Graph Walker and CPU Node Executor*
  - *User Mode Driver*

- \* *Memory Management*
- *Building the first Neuron application*
  - *Prerequisites*
  - *Getting a NEFF file*
  - *The Code*
  - *Code Breakdown*
    - \* *Initialization and cleanup*
    - \* *Loading the NEFF*
    - \* *Creating input/output tensors*
    - \* *Iterating through tensors in an `nrt_tensor_set_t`*
    - \* *Deallocating input/output tensors*
    - \* *Executing the NEFF*
- *The LIBNRT API*
  - *API Return Codes*
  - *Initialization, configuration and teardown*
    - \* *Environment variables used to configure the Runtime Library*
  - *The Model API*
    - \* *Environment variables used to configure a model being loaded*
  - *The Tensor API*
    - \* *The Tensorset API*
  - *The Execution API*
  - *The Profiling API*
  - *Other APIs*

## Introduction

This guide is intended to support a deeper understanding of the Neuron Runtime and how ML applications are built using the Runtime APIs directly. Most customers will not need this level of detail as the interactions with the Neuron Runtime are already taken care by popular ML Frameworks with built-in Neuron support such as torch-neuron and tensorflow-neuron. This guide is focused on the information you need to know when building custom frameworks that will call libnrt APIs directly from C/C++ apps.

---

**Note:** The next few paragraphs provide a brief introduction to the Neuron hardware and the Neuron Runtime architecture. Customers who'd rather skip this and jump straight to building their first ML application which runs without the aid of an ML framework, should go to [Building the first Neuron application](#).

---

The Neuron Runtime Library (libnrt) is the intermediate layer between Application + Framework and Neuron Driver + Neuron Device. It provides a C API for initializing the Neuron hardware, staging models and input data, executing inferences and training iterations on the staged models, and retrieving output data. The vast majority of ML applications running on Neuron will follow one of the following 3 architectural templates:



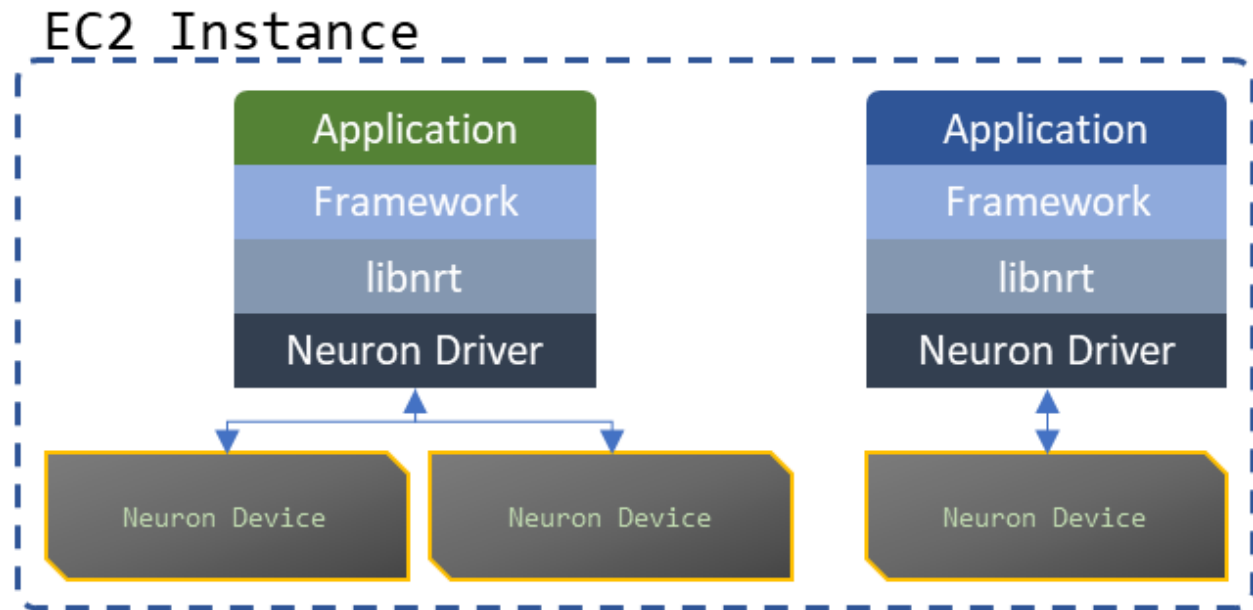


Fig. 6.1: Individual processes executing models on one or more Neuron Devices

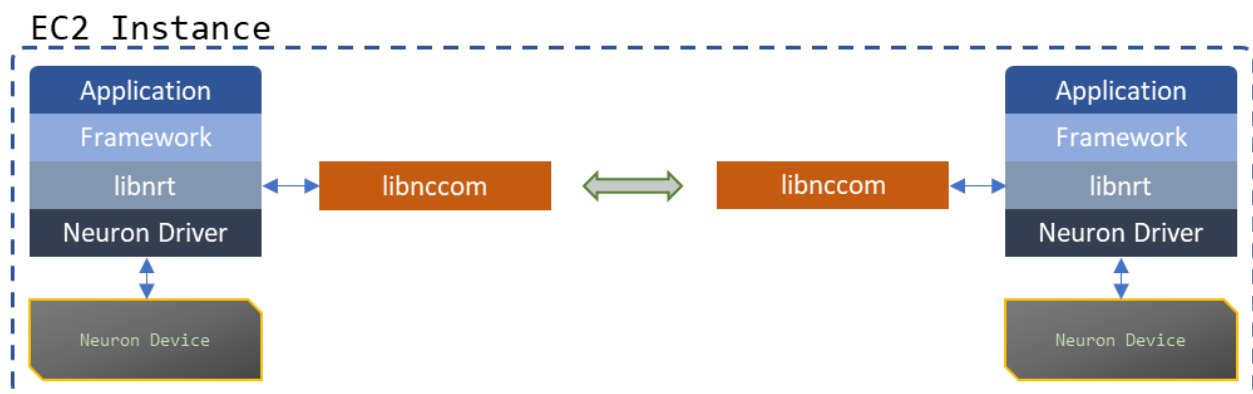


Fig. 6.2: Processes working together on executing models within the same instance - `libnccom` (The Neuron Collective Communication Library) handles inter-worker communication

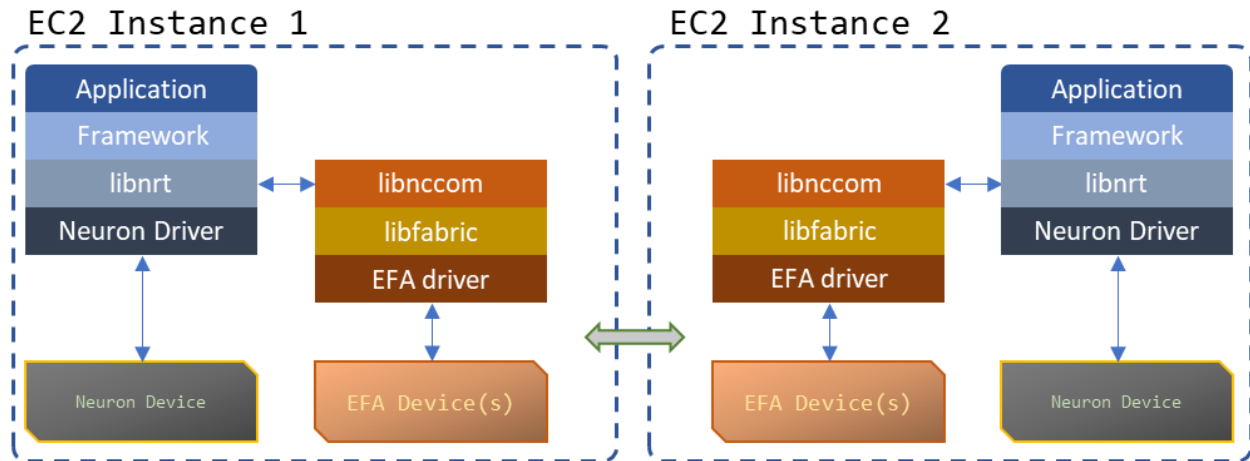


Fig. 6.3: Processes working together on executing models across multiple instances - libnccom, libfabric and the EFA driver handle communication

## Required Software

A more comprehensive guide to installing Neuron software can be found in the `torch_quick_start` guide.

The Neuron Runtime requires the Neuron Driver, which is provided by the `aws-neuronx-dkms` package:

AL2:

```
sudo yum install aws-neuronx-dkms
```

Ubuntu:

```
sudo apt-get install aws-neuronx-dkms
```

The Runtime Library consists of the `libnrt.so` and header files. These artifacts are version controlled and installed via the `aws-neuronx-runtime-lib` package. After installing the package, the binary (`libnrt.so`) is found in `/opt/aws/neuron/lib` and the needed header files are found in `/opt/aws/neuron/include`:

AL2:

```
sudo yum install aws-neuronx-runtime-lib
```

Ubuntu:

```
sudo apt-get install aws-neuronx-runtime-lib
```

For applications that use distributed training or distributed inferences, the Neuron Collective Communication Library is required:

AL2:

```
sudo yum install aws-neuronx-collectives
```

Ubuntu:

```
sudo apt-get install aws-neuronx-collectives
```

In case of multi-instance training, the EFA driver and the Libfabric library - provided by the EFA installer - need to be installed as well:

AL2 & Ubuntu:

```
curl -O https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz
wget https://efa-installer.amazonaws.com/aws-efa-installer.key && gpg --import aws-efa-
↪installer.key
cat aws-efa-installer.key | gpg --fingerprint
wget https://efa-installer.amazonaws.com/aws-efa-installer-latest.tar.gz.sig && gpg --
↪verify ./aws-efa-installer-latest.tar.gz.sig

tar -xvf aws-efa-installer-latest.tar.gz
cd aws-efa-installer && sudo bash efa_installer.sh --yes
cd
sudo rm -rf aws-efa-installer-latest.tar.gz aws-efa-installer
```

## Brief Introduction to Neuron Hardware

Neuron Machine Learning Accelerators (or Neuron Devices) are custom accelerators designed to efficiently execute Machine Learning workloads such as executing inference on a given model or running a distributed training job. Depending on the type of workload and its size, customers can opt for the following Neuron-equipped EC2 instances:

| Workload type | Neuron Device Name | Instance type(s)          | Devices Per Instance | Availability   |
|---------------|--------------------|---------------------------|----------------------|----------------|
| Inference     | Inferentia II (v3) | inf2.xlarge, inf2.8xlarge | 1                    | Available Now! |
| Inference     | Inferentia II (v3) | inf2.24xlarge             | 6                    | Available Now! |
| Inference     | Inferentia II (v3) | inf2.48xlarge             | 12                   | Available Now! |
| Inference     | Inferentia (v1)    | inf1.xlarge, inf1.2xlarge | 1                    | Available Now! |
| Inference     | Inferentia (v1)    | inf1.6xlarge              | 4                    | Available Now! |
| Inference     | Inferentia (v1)    | inf1.24xlarge             | 16                   | Available Now! |
| Training      | Trainium (v2)      | trn1.2xlarge              | 1                    | Available Now! |
| Training      | Trainium (v2)      | trn1.32xlarge             | 16                   | Available Now! |

## Neuron Device

Each Neuron Device consists of multiple execution units - called NeuronCores, a high throughput device memory, PCIe interfaces to the host CPU and to the other Neuron Devices and other components, depending on the Neuron Device version.

To get the number of NeuronCores per Neuron Device, the amount of Neuron Device memory and the way devices are directly connected, use the `neuron-ls` tool:

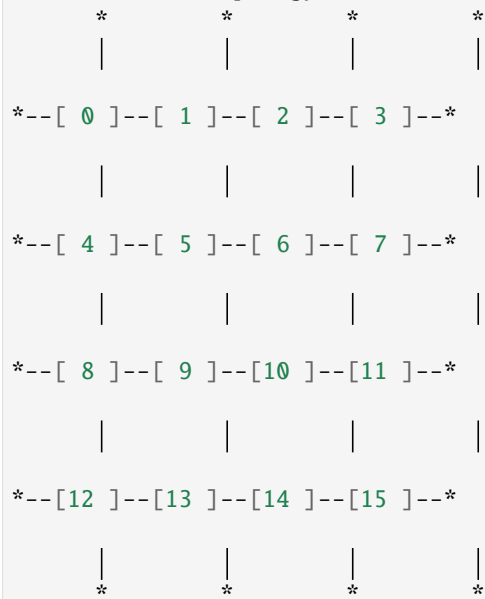
```
neuron-ls --topology
instance-type: trn1.32xlarge
instance-id: i-0633517e496256bf8
+-----+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | CONNECTED | PCI |
| DEVICE | CORES | MEMORY | DEVICES | BDF |
```

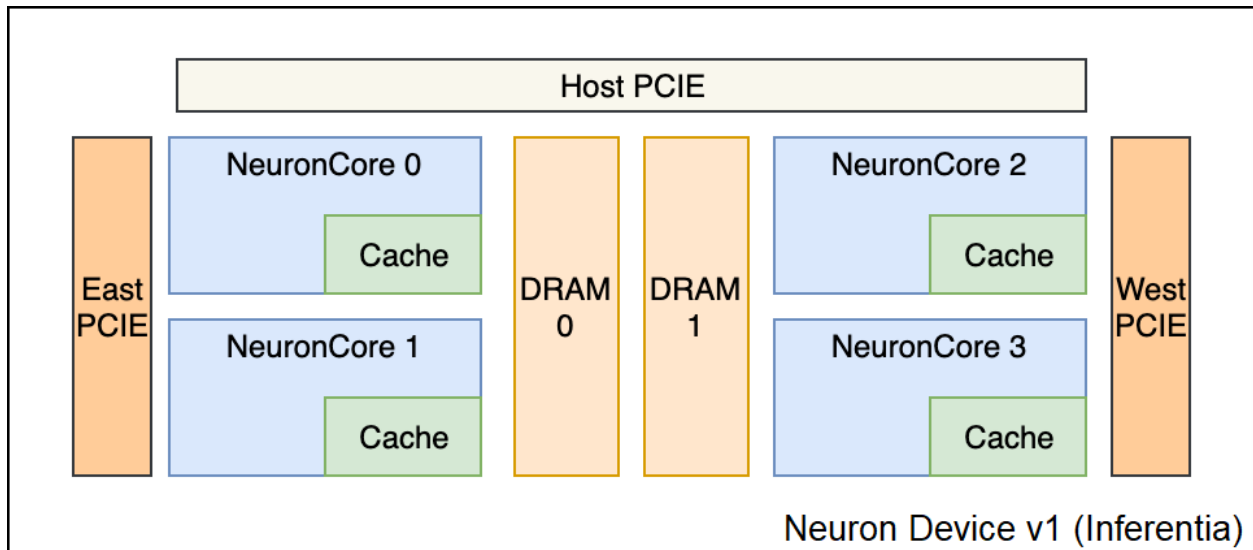
(continues on next page)

(continued from previous page)

|    |   |       |               |         |
|----|---|-------|---------------|---------|
| 0  | 2 | 32 GB | 12, 3, 4, 1   | 10:1c.0 |
| 1  | 2 | 32 GB | 13, 0, 5, 2   | 10:1d.0 |
| 2  | 2 | 32 GB | 14, 1, 6, 3   | a0:1c.0 |
| 3  | 2 | 32 GB | 15, 2, 7, 0   | a0:1d.0 |
| 4  | 2 | 32 GB | 0, 7, 8, 5    | 20:1b.0 |
| 5  | 2 | 32 GB | 1, 4, 9, 6    | 20:1c.0 |
| 6  | 2 | 32 GB | 2, 5, 10, 7   | 90:1b.0 |
| 7  | 2 | 32 GB | 3, 6, 11, 4   | 90:1c.0 |
| 8  | 2 | 32 GB | 4, 11, 12, 9  | 20:1d.0 |
| 9  | 2 | 32 GB | 5, 8, 13, 10  | 20:1e.0 |
| 10 | 2 | 32 GB | 6, 9, 14, 11  | 90:1d.0 |
| 11 | 2 | 32 GB | 7, 10, 15, 8  | 90:1e.0 |
| 12 | 2 | 32 GB | 8, 15, 0, 13  | 10:1e.0 |
| 13 | 2 | 32 GB | 9, 12, 1, 14  | 10:1b.0 |
| 14 | 2 | 32 GB | 10, 13, 2, 15 | a0:1e.0 |
| 15 | 2 | 32 GB | 11, 14, 3, 12 | a0:1b.0 |

## Neuron Device Topology

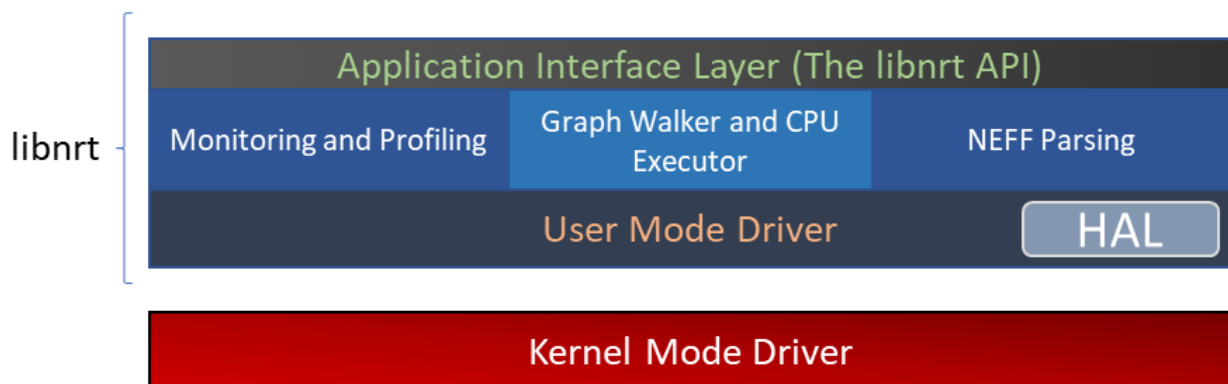




### NeuronCore

The NeuronCore is the primary execution unit within the accelerator. Each NeuronCore contains several execution engines (for different types of compute operations such as tensor-based, vector and scalar), DMA engines, and a local cache. A NeuronCore can operate independently or together with other NeuronCores, depending on the nature of the workload and the way a model is compiled and loaded to the NeuronCores in the accelerator. Each execution engine can access the cache and DRAM attached to the accelerator device. The primary form of data movement between the host CPU and the accelerator device, as well as between the device DRAM and NeuronCores, is Direct Memory Access (DMA). The use of DMA enables more efficient data movement.

### The Neuron Runtime Architecture



## Application Interface Layer (The libnrt API)

The Application Interface Layer allows applications and frameworks to use the available Neuron Devices to run inference or training workloads. A complete reference of the C interface can be found in [The LIBNRT API](#).

## Monitoring and Profiling

The Neuron Runtime is able to capture key execution metrics which can be read in real-time using `neuron-monitor` and `neuron-top`. `neuron-monitor` allows forwarding those metrics to CloudWatch or a Prometheus server, enabling fleet-wide monitoring - for more on that please refer to the `neuron-monitor` usage guide [Neuron Monitor User Guide](#). Profiling an execution is another feature of the Neuron Runtime - which provides an API for starting and stopping profiling, as well as saving the profile data to a file, which can be used by tools such as the Neuron Tensorboard. This API is documented in [The Profiling API](#) section.

## The NEFF format and NEFF Parser

A NEFF (*\*N*euron *\*E*xecutable *\*F*ile *\*F*ormat) is a single file container for all the artifacts needed to execute a model on one or more NeuronCores. A NEFF is the output of the Neuron Compiler (`neuron-cc`). It contains Neuron machine instructions, pseudo instructions (compiler-generated instructions which are parsed and replaced with Neuron instructions by the Neuron Runtime when the model loads), tensor information, model parameters and other components that support the model's execution on one or more NeuronCores. Operators that are not supported by Neuron can be compiled into CPU-executable binary and included into the NEFF as well.

The contents of a NEFF can be shown by using `neuron-packager` tool (which will be released soon).

Usually there is only one subgraph (which is executed on a single NeuronCore) in a NEFF:

NEFF Nodes:

| NODE    | Executor    | Name       | Variable     | Size    | Type | Format |          |
|---------|-------------|------------|--------------|---------|------|--------|----------|
| → Shape | Data Type   | TimeSeries |              |         |      |        |          |
| 1       | Neuron Core | sg00       | image:0      | 3259008 | IN   | NHWC   | [1 3 552 |
| → 984]  |             |            | net_output:0 | 1323972 | OUT  | NHWC   | [1 78 69 |
| → 123]  | false       |            |              |         |      |        |          |

In this example, there is a single subgraph, one input and one output:



Some NEFFs can have multiple subgraphs (which will be deployed by the runtime on separate NeuronCores) and multiple CPU operators, as exemplified below:

NEFF Nodes:

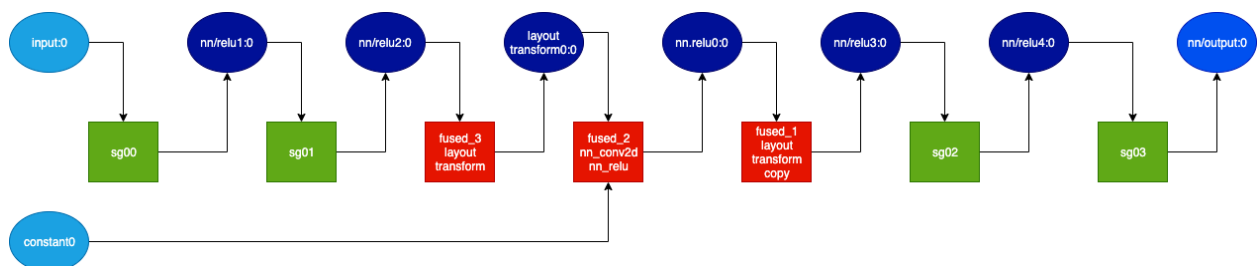
| NODE   | Executor | Name  | Variable  | Size       |
|--------|----------|-------|-----------|------------|
| → Type | Format   | Shape | Data Type | TimeSeries |

(continues on next page)

(continued from previous page)

|   |     |             |                               |       |  |                     |   |   |
|---|-----|-------------|-------------------------------|-------|--|---------------------|---|---|
|   | 1   | Neuron Core |                               | sg00  |  | input:0             | 2 | └ |
| ↪ | IN  | NHWC        | [1 1 1 1]                     |       |  | nn/relu1:0          | 2 | └ |
| ↪ | OUT | NHWC        | [1 1 1 1]                     | false |  | nn/relu1:0          | 2 | └ |
|   | 1   | Neuron Core |                               | sg01  |  | nn/relu1:0          | 2 | └ |
| ↪ | IN  | NHWC        | [1 1 1 1]                     |       |  | nn/relu2:0          | 2 | └ |
| ↪ | OUT | NHWC        | [1 1 1 1]                     | false |  | nn/relu2:0          | 2 | └ |
|   | 2   | CPU         | fused_3_layout_transform      |       |  | layout_transform0:0 | 0 | └ |
| ↪ | OUT |             | []                            |       |  | constant0           | 2 | └ |
|   | 4   | CPU         | fused_2_nn_conv2d_nn_relu     |       |  | nn.relu0:0          | 0 | └ |
| ↪ | IN  |             | [1 1 1 1] float16             |       |  | nn.relu0:0          | 0 | └ |
| ↪ | OUT |             | []                            |       |  | nn/relu3:0          | 0 | └ |
|   | 5   | CPU         | fused_1_layout_transform_copy |       |  | nn/relu3:0          | 0 | └ |
| ↪ | OUT |             | []                            |       |  | nn/relu3:0          | 2 | └ |
|   | 6   | Neuron Core |                               | sg02  |  | nn/relu4:0          | 2 | └ |
| ↪ | IN  | NHWC        | [1 1 1 1]                     |       |  | nn/relu4:0          | 2 | └ |
| ↪ | OUT | NHWC        | [1 1 1 1]                     | false |  | nn/relu4:0          | 2 | └ |
|   | 6   | Neuron Core |                               | sg03  |  | nn/relu4:0          | 2 | └ |
| ↪ | IN  | NHWC        | [1 1 1 1]                     |       |  | nn/output:0         | 2 | └ |
| ↪ | OUT | NHWC        | [1 1 1 1]                     | false |  |                     |   |   |

The output above can be summarized by the graph below:



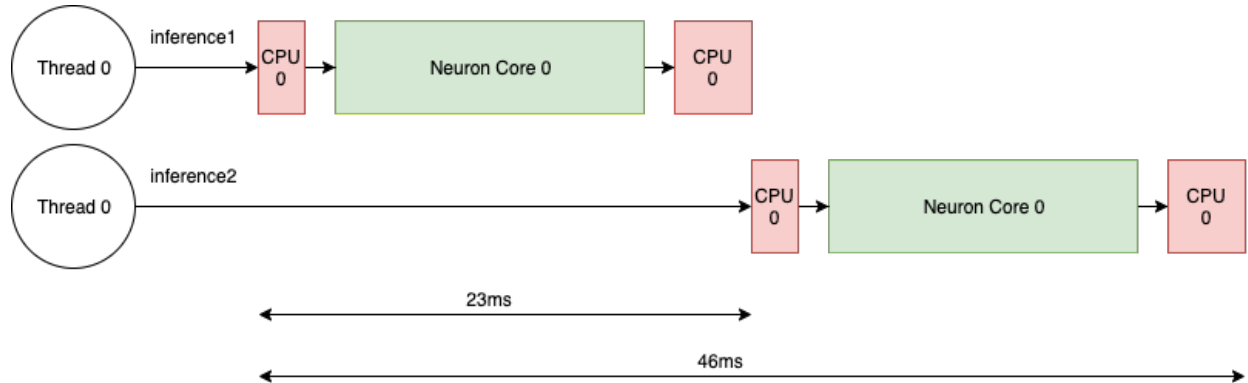
The nodes marked with dark blue are intermediate tensors that are handled internally by the Neuron Runtime. The other blue nodes are inputs/outputs. The green colored box indicates the operator is executed on the NeuronCore while the red color box indicates the execution is done on the CPU.

The NEFF layer in Neuron Runtime is responsible for parsing a NEFF, validating it, and translating pseudo instructions into hardware specific instructions and DMA descriptors.

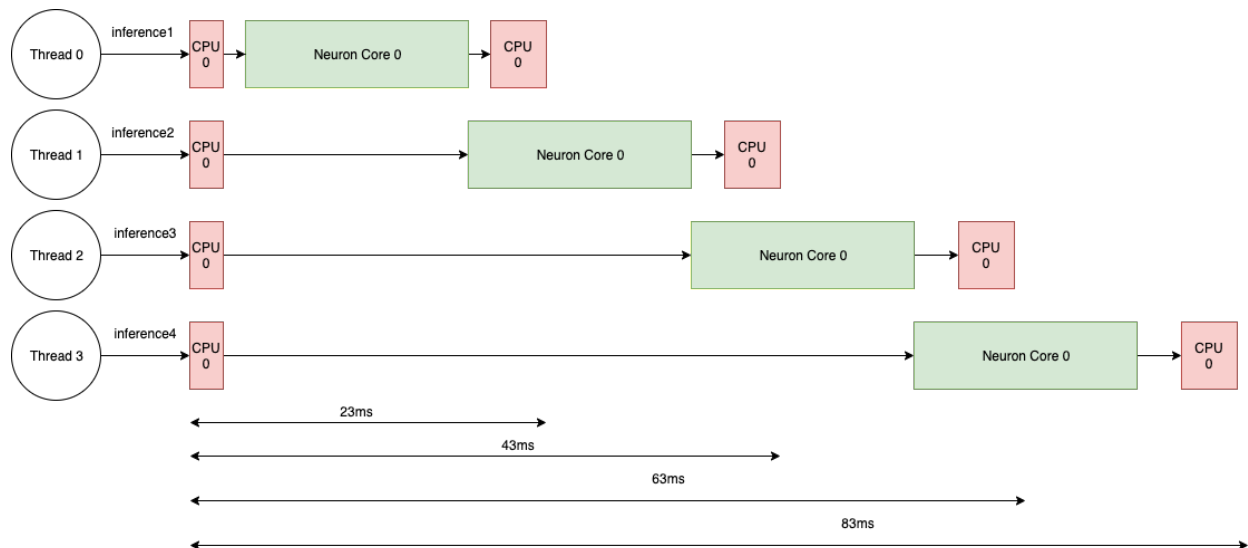
## Graph Walker and CPU Node Executor

As shown in the previous section, a NEFF can contain one or more nodes. During execution, the Neuron Runtime Graph Walker executes each node one by one and handles copying input and output between each of them. If a node needs to be executed by the CPU, then a corresponding library function, found in a .so file in the NEFF, is dynamically loaded using `dlopen()` during model load and executed during model execution. Since this library function is executed in the calling thread's context, the workload can be efficiently parallelized using a multi-threaded approach.

In the example below, each invocation of `nrt_execute()` would take 23ms: the first CPU node takes 1ms, the NeuronCore execution takes 20ms and the second CPU node takes 2 ms, so the total latency is 23ms and the throughput is 43 calls per second ( $1000/23$ ).



If multiple threads are used, subsequent executions would be pipelined inside the runtime, hence increasing the throughput in this case to ~50 ( $1000/20$ ).





## User Mode Driver

This is the lowest level component of the Neuron Runtime and handles programming the engines, managing memory, creating DMA descriptors to move data from host and device, handling notifications etc.

## Memory Management

The Neuron Runtime is responsible with managing Neuron Device and host memory for the running models. The application is responsible with deallocating every loaded model and allocated tensor so the proper deallocation method needs to be called. For more details, refer to *The LIBNRT API* documentation. Tools such as `neuron-top` and `neuron-monitor` can be used to determine the amount of memory being used at any given time.

## Building the first Neuron application

The simple application presented here will load a NEFF file, use the provided binary files' contents as input tensors (if a file wasn't provided for an input tensor, that input tensor will be zero-filled), and save the output tensors as binary files.

## Prerequisites

Building the application requires:

- a recent version of GCC
- installing the `aws-neuronx-runtime-lib` package as described in *Required Software*

Running the built application requires:

- a Neuron-equipped instance as shown in *Brief Introduction to Neuron Hardware*
- installing the `aws-neuronx-runtime-lib` and the `aws-neuronx-dkms` package as described in *Required Software*
- a NEFF file

## Getting a NEFF file

When running any workload through a Neuron framework, the compiled NEFFs will be placed in `/var/tmp/neuron-compile-cache`. Additionally, setting the `NEURON_FRAMEWORK_DEBUG` environment variable to 1 before running the workload will enable the compiled NEFFs to be written to the current directory.

## The Code

```
#include <stdbool.h>
#include <nrt/nrt.h>
#include <nrt/nrt_experimental.h>

#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include <time.h>
```

(continues on next page)

(continued from previous page)

```

#include <errno.h>
#include <sys/mman.h>
#include <sys/stat.h>
#include <pthread.h>
#include <fcntl.h>
#include <stdint.h>
#include <unistd.h>

// Function to mmap a file in the application's memory space,
// it will return a pointer to the mmapped memory and the size
// of the mmapped data will be written to *size
void *mmap_file(const char *filepath, size_t *size) {
    struct stat sb;
    int fd = open(filepath, O_RDONLY);
    if (fd < 0 || fstat(fd, &sb) != 0) {
        fprintf(stderr, "Unable to open %s: %s\n", filepath, strerror(errno));
        return MAP_FAILED;
    }
    *size = sb.st_size;
    return mmap(NULL, sb.st_size, PROT_READ, MAP_PRIVATE, fd, 0);
}

#define P_ERR(...) fprintf(stderr, __VA_ARGS__)

#define CHECK_RESULT(res, expected, ...) \
    if (res != expected) { \
        fprintf(stderr, __VA_ARGS__); \
        exit(-1); \
    }

// struct used to load input tensors from files
typedef struct {
    char *name;
    size_t size;
    void *data;
} input_tensor_info_t;

// simple container for input_tensor_info_t
typedef struct {
    input_tensor_info_t *entries;
    int entry_count;
} input_tensor_info_array_t;

// Allocate tensorsets and tensors based on the info_array and returns a valid tensorset
// in out_tset
// containing all the newly allocated tensors
NRT_STATUS allocate_tensors(nrt_tensor_info_array_t *info_array, nrt_tensor_usage_t
// usage_type, nrt_tensor_set_t **out_tset) {
    NRT_STATUS result;
    int tensor_idx;
    nrt_tensor_info_t *tensor_info = NULL;
    nrt_tensor_t *tensor = NULL;

```

(continues on next page)

(continued from previous page)

```

// We allocate a nrt_tensor_set which acts as a containers for nrt_tensors
result = nrt_allocate_tensor_set(out_tset);
if (result != NRT_SUCCESS) {
    P_ERR("Couldn't allocate %s tensorset\n", usage_type == NRT_TENSOR_USAGE_INPUT ?
↪ "input" : "output");
}

for (tensor_idx = 0; tensor_idx < info_array->tensor_count; tensor_idx++) {
    tensor_info = &info_array->tensor_array[tensor_idx];
    if (tensor_info->usage != usage_type) {
        continue;
    }
    // Allocate the tensor with the name and size found in tensor_info_array
    result = nrt_tensor_allocate(NRT_TENSOR_PLACEMENT_DEVICE, 0, tensor_info->size,
                                tensor_info->name, &tensor);
    if (result != NRT_SUCCESS) {
        P_ERR("Couldn't allocate tensor %s\n", tensor_info->name);
        return result;
    }
    // Finally add the tensors to the newly allocated tensor set
    result = nrt_add_tensor_to_tensor_set(*out_tset, tensor_info->name, tensor);
    if (result != NRT_SUCCESS) {
        P_ERR("Couldn't add tensor %s to tensorset\n", tensor_info->name);
        return result;
    }
}
return NRT_SUCCESS;
}

// Tensor iterator handler - returns false if the iteration needs to stop
typedef bool (*tensor_handler)(nrt_tensor_t *, nrt_tensor_info_t *, NRT_STATUS *, void_
↪ *);

// Iterates through all the tensors in the given tensorset, based on the data in info_
↪ array for the given usage_type
// and calls the handler function with the provided args pointer
// Will return the first error returned by a handler
NRT_STATUS iterate_tensors(nrt_tensor_set_t *tset, nrt_tensor_info_array_t *info_array,
↪ nrt_tensor_usage_t usage_type,
                           tensor_handler handler, void *args) {
    NRT_STATUS result = NRT_SUCCESS;
    NRT_STATUS final_result = NRT_SUCCESS;
    int tensor_idx;
    nrt_tensor_info_t *tensor_info = NULL;
    nrt_tensor_t *tensor = NULL;

    for (tensor_idx = 0; tensor_idx < info_array->tensor_count; tensor_idx++) {
        tensor_info = &info_array->tensor_array[tensor_idx];
        if (tensor_info->usage != usage_type) {
            continue;
        }
    }
}

```

(continues on next page)

(continued from previous page)

```

    result = nrt_get_tensor_from_tensor_set(tset, tensor_info->name, &tensor);
    if (result != NRT_SUCCESS) {
        P_ERR("Tensor %s not found in tensor set\n", tensor_info->name);
        continue;
    }
    result = NRT_SUCCESS;
    if ((*handler)(tensor, tensor_info, &result, args) == false) {
        return result;
    }
    if (final_result == NRT_SUCCESS && result != final_result) {
        final_result = result;
    }
}
return final_result;
}

// Tensor iteration handler that checks if a tensor has an input file associated with it
// based on the CLI args
bool handler_load_inputs(nrt_tensor_t *tensor, nrt_tensor_info_t *tensor_info, NRT_
↳STATUS *result, void* args) {
    NRT_STATUS res;
    int idx;
    input_tensor_info_array_t *info_array = (input_tensor_info_array_t *)args;
    bool input_found = false;

    for (idx = 0; idx < info_array->entry_count; idx++) {
        if (strcmp(info_array->entries[idx].name, tensor_info->name) != 0) {
            continue;
        }
        if (info_array->entries[idx].size != tensor_info->size) {
            P_ERR("Input file for tensor %s has incorrect size %lu, expected %lu\n",
                tensor_info->name, info_array->entries[idx].size, tensor_info->size);
            break;
        }
        res = nrt_tensor_write(tensor, info_array->entries[idx].data, 0, tensor_info->
↳size);
        if (res != NRT_SUCCESS) {
            P_ERR("Unable to write content to input tensor %s\n", tensor_info->name);
        } else {
            input_found = true;
        }
    }
    if (!input_found) {
        fprintf(stderr, "Input tensor %s will be zero-filled\n", tensor_info->name);
    }
    *result = NRT_SUCCESS;
    return true;
}

// Tensor iteration handler that saves outputs
bool handler_save_outputs(nrt_tensor_t *tensor, nrt_tensor_info_t *tensor_info, NRT_
↳STATUS *result, void* args) {

```

(continues on next page)

(continued from previous page)

```

static char filename[280];

int fd;
// Allocating a buffer large enough to read the entire tensor
void *tensor_data = malloc(tensor_info->size);

*result = NRT_SUCCESS;
if (tensor_data == NULL) {
    fprintf(stderr, "Unable to allocate memory for saving output tensor %s\n",
↪ tensor_info->name);
    *result = NRT_FAILURE;
    return true;
}
// Reading the tensor to the newly allocated buffer
*result = nrt_tensor_read(tensor, tensor_data, 0, tensor_info->size);
if (*result != NRT_SUCCESS) {
    fprintf(stderr, "Unable to read tensor %s\n", tensor_info->name);
    free(tensor_data);
    return true;
}

// Saving the tensor to a file
snprintf(filename, 280, "%s.out", tensor_info->name);
fd = open(filename, O_WRONLY | O_CREAT | O_TRUNC, 0644);
if (fd < 0) {
    fprintf(stderr, "Unable to open %s for writing\n", filename);
    free(tensor_data);
    *result = NRT_FAILURE;
    return true;
}
if (write(fd, tensor_data, tensor_info->size) != tensor_info->size) {
    *result = NRT_FAILURE;
    fprintf(stderr, "Unable to write tensor %s contents to file %s\n", tensor_info->
↪ name, filename);
}
close(fd);

free(tensor_data);
return true;
}

// Tensor iteration handler that deallocates tensors
bool handler_free_tensor(nrt_tensor_t *tensor, nrt_tensor_info_t *tensor_info, NRT_
↪ STATUS *result, void* args) {
    *result = NRT_SUCCESS;
    nrt_tensor_free(&tensor);
    return true;
}

int main(int argc, char *argv[]) {
    NRT_STATUS result;
    int idx = 0;

```

(continues on next page)

(continued from previous page)

```

int tensor_idx = 0;
void *neff_data = NULL;
size_t neff_size = 0;
void *input_data = NULL;

input_tensor_info_array_t input_tensor_info_array = {0};
input_tensor_info_t *current_input = NULL;

nrt_model_t *model = NULL;
nrt_tensor_set_t *inputs = NULL;
nrt_tensor_set_t *outputs = NULL;

nrt_tensor_t *tensor = NULL;
nrt_tensor_info_array_t *tensor_info_array = NULL;

if (argc < 2) {
    fprintf(stderr, "Incorrect number of args, usage: exec_test file.neff [input_1_
↪name] [input_1_file] ...\n");
    exit(-1);
}

// Try mmaping the NEFF file first, so we can fail fast if not found or
// mmap fails
neff_data = mmap_file(argv[1], &neff_size);
if (neff_data == MAP_FAILED) {
    fprintf(stderr, "Unable to map file %s\n", argv[1]);
    exit(-1);
}

// mmap input tensor files (if any provided) and fill the input_tensor_info array
if (argc > 3) {
    input_tensor_info_array.entries = malloc((argc - 2 / 2) * sizeof(input_tensor_
↪info_t));
    for (idx = 2; idx < argc; idx += 2) {
        if (idx + 1 >= argc) {
            break;
        }
        current_input = &input_tensor_info_array.entries[input_tensor_info_array.
↪entry_count];
        input_data = mmap_file(argv[idx + 1], &current_input->size);
        if (input_data == MAP_FAILED) {
            fprintf(stderr, "Unable to mmap inputs file %s\n", argv[idx + 1]);
            continue;
        }
        current_input->name = argv[idx];
        current_input->data = input_data;
        input_tensor_info_array.entry_count++;
    }
}

// Before calling any nrt API, nrt_init must be called
// Since this is not running as part of a framework, the correct parameter for

```

(continues on next page)

(continued from previous page)

```

↪ 'framework' is
    // NRT_FRAMEWORK_TYPE_NO_FW and the others can be empty strings
    result = nrt_init(NRT_FRAMEWORK_TYPE_NO_FW, "", "");
    CHECK_RESULT(result, NRT_SUCCESS, "NRTLIB could not be initialized, error: %d\n",
↪ (int)result);

    // Loading the NEFF
    printf("Loading NEFF\n");
    result = nrt_load(neff_data, neff_size, -1, -1, &model);
    CHECK_RESULT(result, NRT_SUCCESS, "Unable to load NEFF\n");

    // In order to allocate tensors, first we need to call nrt_get_model_tensor_info_
↪ which
    // will give us the model tensors' names and sizes in tensor_info_array
    printf("Getting IO tensor information\n");
    result = nrt_get_model_tensor_info(model, &tensor_info_array);
    CHECK_RESULT(result, NRT_SUCCESS, "Unable to get model tensor information\n");

    // Allocating tensors
    printf("Creating I/O data (%ld tensors)\n", tensor_info_array->tensor_count);
    result = allocate_tensors(tensor_info_array, NRT_TENSOR_USAGE_INPUT, &inputs);
    CHECK_RESULT(result, NRT_SUCCESS, "Error allocating input tensors\n");
    result = allocate_tensors(tensor_info_array, NRT_TENSOR_USAGE_OUTPUT, &outputs);
    CHECK_RESULT(result, NRT_SUCCESS, "Error allocating input tensors\n");

    // Loading input files (if provided)
    iterate_tensors(inputs, tensor_info_array, NRT_TENSOR_USAGE_INPUT, handler_load_
↪ inputs,
        (void*) &input_tensor_info_array);

    // Executing model using the tensors in the inputs tensorset and writing the outputs_
↪ to the tensors
    // in the outputs tensorset
    result = nrt_execute(model, inputs, outputs);
    CHECK_RESULT(result, NRT_SUCCESS, "Error during model execution: %d\n", result);

    // Saving outputs to files
    result = iterate_tensors(outputs, tensor_info_array, NRT_TENSOR_USAGE_OUTPUT,
↪ handler_save_outputs, NULL);
    if (result != NRT_SUCCESS) {
        P_ERR("Error saving outputs to files\n");
    }

    // Unloading the model
    result = nrt_unload(model);
    if (result != NRT_SUCCESS) {
        P_ERR("Unable to unload NEFF\n");
    }

    printf("Freeing tensors\n");
    iterate_tensors(inputs, tensor_info_array, NRT_TENSOR_USAGE_INPUT, handler_free_
↪ tensor, NULL);

```

(continues on next page)

(continued from previous page)

```

    iterate_tensors(outputs, tensor_info_array, NRT_TENSOR_USAGE_OUTPUT, handler_free_
↪ tensor, NULL);

    nrt_destroy_tensor_set(&inputs);
    nrt_destroy_tensor_set(&outputs);

    printf("Deallocating model tensor info\n");
    // We are done with the tensor_info_array, we can dispose of it
    nrt_free_model_tensor_info(tensor_info_array);

    printf("Deallocating inputs tensor info\n");
    // Unmapping the input files
    for (tensor_idx = 0; tensor_idx < input_tensor_info_array.entry_count; tensor_idx++)
↪ {
        munmap(input_tensor_info_array.entries[tensor_idx].data, input_tensor_info_array.
↪ entries[tensor_idx].size);
    }
    if (input_tensor_info_array.entries) {
        free(input_tensor_info_array.entries);
    }

    // Clean-up the runtime
    printf("Cleaning up the runtime\n");
    nrt_close();

    printf("DONE\n");
}

```

Building the example:

```
gcc run_neff.c -o run_neff -lnrt -pthread -I/opt/aws/neuron/include -L/opt/aws/neuron/lib
```

Running the example:

```
./run_neff my.neff [input_1] [input_1.bin] [input_2] [input_2.bin] ...
```

## Code Breakdown

### Initialization and cleanup

```

// ...
result = nrt_init(NRT_FRAMEWORK_TYPE_NO_FW, "", "");
// ...
nrt_close();

```

The Neuron Runtime is initialized by calling `nrt_init` and all applications should call `nrt_close` once they're done using it. For more details on these functions, go to the [Initialization, configuration and teardown](#) section.



## Loading the NEFF

Once the contents of a NEFF file have been mapped to virtual memory using `mmap` ...

```
// ...
void *mmap_file(const char *filepath, size_t *size) {
    struct stat sb;
    int fd = open(filepath, O_RDONLY);
    if (fd < 0 || fstat(fd, &sb) != 0) {
        fprintf(stderr, "Unable to open %s: %s\n", filepath, strerror(errno));
        return MAP_FAILED;
    }
    *size = sb.st_size;
    return mmap(NULL, sb.st_size, PROT_READ, MAP_PRIVATE, fd, 0);
}
// ...
neff_data = mmap_file(argv[1], &neff_size);
```

... the NEFF is loaded using `nrt_load`. The runtime will decide the optimal placement for the model - it will choose the best NeuronCore on which to deploy the model:

```
// ...
result = nrt_load(neff_data, neff_size, -1, -1, &model);
// ...
```

The call will return a valid model handle in `nrt_model_t*` which will subsequently be used for other calls to the Runtime API (such as `nrt_execute`).

For more details on the model API (including `nrt_load`), go to the [The Model API](#) section.

## Creating input/output tensors

The main container for tensors is the `nrt_tensor_set_t*`. Tensors (`nrt_tensor_t*`) are not passed directly to the NEFF execution function, `nrt_execute`, they have to be wrapped in a `nrt_tensor_set_t*`. The `allocate_tensors` function will allocate the tensorset and the tensors for the requested usage type (`NRT_TENSOR_USAGE_INPUT` or `NRT_TENSOR_USAGE_OUTPUT`) and return the tensorset containing the allocated tensors in `out_tset`.

```
NRT_STATUS allocate_tensors(nrt_tensor_info_array_t *info_array, nrt_tensor_usage_t
↪ usage_type, nrt_tensor_set_t **out_tset) {
    // ...
    // We allocate a nrt_tensor_set which acts as a containers for nrt_tensors
    result = nrt_allocate_tensor_set(out_tset);
    // ...

    for (tensor_idx = 0; tensor_idx < info_array->tensor_count; tensor_idx++) {
        tensor_info = &info_array->tensor_array[tensor_idx];
        if (tensor_info->usage != usage_type) {
            continue;
        }
        // ...
        // Allocate the tensor with the name and size found in tensor_info_array
        result = nrt_tensor_allocate(NRT_TENSOR_PLACEMENT_DEVICE, 0, tensor_info->size,
```

(continues on next page)

(continued from previous page)

```

                                tensor_info->name, &tensor);
    // ...
    // Finally add the tensors to the newly allocated tensor set
    result = nrt_add_tensor_to_tensor_set(*out_tset, tensor_info->name, tensor);
    // ...
}
// ...
}

```

### Iterating through tensors in an `nrt_tensor_set_t`

A helper function, `iterate_tensors` is used to iterate through the `nrt_tensor_t` in a tensorset and call the function handler for each of them. If the handler function returns `false` iteration ends. `iterate_tensors` returns the first error reported by the handler function.

```

// Tensor iterator handler - returns false if the iteration needs to stop
typedef bool (*tensor_handler)(nrt_tensor_t *, nrt_tensor_info_t *, NRT_STATUS *, void_
↪*);

NRT_STATUS iterate_tensors(nrt_tensor_set_t *tset, nrt_tensor_info_array_t *info_array, ↪
↪nrt_tensor_usage_t usage_type,
                           tensor_handler handler, void *args) {
    // ...
    for (tensor_idx = 0; tensor_idx < info_array->tensor_count; tensor_idx++) {
        // ...
        result = nrt_get_tensor_from_tensor_set(tset, tensor_info->name, &tensor);
        // ...
        if ((*handler)(tensor, tensor_info, &result, args) == false) {
            return result;
        }
        // ...
    }
}

```

### Deallocating input/output tensors

After the execution is complete, the tensors are deallocated using `iterate_tensors` and the tensorsets are deallocated using `nrt_destroy_tensor_set`:

```

iterate_tensors(inputs, tensor_info_array, NRT_TENSOR_USAGE_INPUT, handler_free_tensor, ↪
↪NULL);
iterate_tensors(outputs, tensor_info_array, NRT_TENSOR_USAGE_OUTPUT, handler_free_tensor,
↪ NULL);

nrt_destroy_tensor_set(&inputs);
nrt_destroy_tensor_set(&outputs);

```

The `handler_free_tensor` function simply deallocates the given tensor:

```

bool handler_free_tensor(nrt_tensor_t *tensor, nrt_tensor_info_t *tensor_info, NRT_
↪STATUS *result, void* args) {

```

(continues on next page)

(continued from previous page)

```
// ...
nrt_tensor_free(&tensor);
// ...
}
```

For more details on the tensor API, check out the *The Tensor API* and the *The Tensorset API* sections.

## Executing the NEFF

The NEFF is executed using a call to `nrt_execute`. If `nrt_execute` completes successfully, the output tensors are read and saved to files (one binary file per output tensor) using `iterate_tensors`:

```
// Executing model using the tensors in the inputs tensorset and writing the outputs to
↳ the tensors
// in the outputs tensorset
result = nrt_execute(model, inputs, outputs);
// ...
// Saving outputs to files
result = iterate_tensors(outputs, tensor_info_array, NRT_TENSOR_USAGE_OUTPUT, handler_
↳ save_outputs, NULL);
```

The iteration handler reads the tensor data and writes it to a file with the same name as the tensor:

```
bool handler_save_outputs(nrt_tensor_t *tensor, nrt_tensor_info_t *tensor_info, NRT_
↳ STATUS *result, void* args) {
    // ...
    void *tensor_data = malloc(tensor_info->size);
    // ...
    // Reading the tensor to the newly allocated buffer
    *result = nrt_tensor_read(tensor, tensor_data, 0, tensor_info->size);
    // ...

    // Saving the tensor to a file
    snprintf(filename, 280, "%s.out", tensor_info->name);
    fd = open(filename, O_WRONLY | O_CREAT | O_TRUNC, 0644);
    // ...
    if (write(fd, tensor_data, tensor_info->size) != tensor_info->size) {
        // ...
    }
    close(fd);
}
```

For more details on the execution API, go to the *The Execution API* section.

## The LIBNRT API

### API Return Codes

All API calls will return an `NRT_STATUS` value representing the return status of the call. In case of an error, an error message will also be logged (based on the logging settings, more on that in the next section). The table below contains all the possible error codes. Please note that some error codes only apply to certain API calls.

| Name  | Return Code | Error  |
|---|-------------|--|
| <code>NRT_SUCCESS</code>                    | 0           | Call was successful  |
| <code>NRT_FAILURE</code>                    | 1           | Generic failure  |
| <code>NRT_INVALID</code>                    | 2           | Invalid NEFF, bad instruction, bad DMA descriptor, input tensor name/size does not match the model, etc.   |
| <code>NRT_INVALID_HANDLE</code>             | 3           | Invalid handle (e.g. an invalid model handle)  |
| <code>NRT_RESOURCE</code>                   | 4           | Failed to allocate a resource for the requested operation  |
| <code>NRT_TIMEOUT</code>                    | 5           | Operation timed out  |
| <code>NRT_HW_ERROR</code>                   | 6           | Hardware failure   |
| <code>NRT_QUEUE_FULL</code>                 |             | Too many pending <code>nrt_execute()</code> requests. The runtime request queue is full. Cannot enqueue more <code>nrt_execute()</code> requests |
| <code>NRT_LOAD_NOT_ENOUGH</code>            | 7           | The number of available NeuronCores is insufficient for the requested operation  |
| <code>NRT_UNSUPPORTED</code>                | 8           | Unsupported NEFF version   |
| <code>NRT_UNINITIALIZED</code>              |             | Returned when attempting an API call when the library is not initialized   |
| <code>NRT_CLOSED</code>                     | 14          | Returned when attempting an API call after <code>nrt_close()</code> was called   |
| <code>NRT_EXEC_BADINPUT</code>              |             | Invalid input has been submitted to <code>nrt_execute()</code>   |
| <code>NRT_EXEC_COMPLETED_WITH_NUMERR</code> |             | Execution completed with numerical errors (produced NaN)   |
| <code>NRT_EXEC_COMPLETED_WITH_ERR</code>    |             | Execution completed with other errors, either logical (event double clear), or hardware (parity error)   |
| <code>NRT_EXEC_NC_BUSY</code>               |             | The neuron core is locked (in use) by another model/thread   |
| <code>NRT_OOB</code>                        | 1006        | One or more indirect memcopies and/or embedding updates are out of bound due to input corruptions  |
| <code>NRT_EXEC_HW1ERR</code>                |             | Unexpected hang in collectives operation due to hardware errors on this or other workers.  |
| <code>NRT_EXEC_HW1ERR_HBM_HUBM</code>       |             | HBM/HUBM suffered from an uncorrectable error and produced incorrect results   |

### Initialization, configuration and teardown

`NRT_STATUS nrt_init(nrt_framework_type_t framework, const char *fw_version, const char *fal_version)`

Initializes the Neuron Runtime's internal state and the Neuron hardware's state. This should be called before any other `nrt_*` call is attempted - although a small set of functions are exempt from this rule (for example `nrt_get_total_nc_count` and `get_nrt_version`). Any call to the NRT library API will return `NRT_FAILURE` if `nrt_init` has not been called beforehand and that API call requires it.

The runtime can be configured by setting the appropriate environment variable before this API call. The list of available environment variables is found in the [Environment variables used to configure the Runtime Library](#) section.

#### Parameters

- **framework** – Can be one of:

`NRT_FRAMEWORK_TYPE_INVALID`, // Invalid framework

`NRT_FRAMEWORK_TYPE_NO_FW`, // No framework `NRT_FRAMEWORK_TYPE_TENSORFLOW`,

```
// Tensorflow NRT_FRAMEWORK_TYPE_PYTORCH, // Pytorch
NRT_FRAMEWORK_TYPE_MXNET // Mxnet
```

This argument is used by our Neuron Tools to determine the type of application running, it has no other impact on the functioning of the runtime. Application using a custom framework or calling the Neuron Runtime directly should use `NRT_FRAMEWORK_TYPE_NO_FW`.

- **\*fw\_version** (const char) – version of the framework on top of which this runtime is running
- **\*fal\_version** (const char) – version of the framework adapter on top of which this runtime is running

Applications using `NRT_FRAMEWORK_TYPE_NO_FW` for the first argument should use two empty strings for the versions.

## Environment variables used to configure the Runtime Library

**NEURON\_RT\_LOG\_LOCATION=<CONSOLE/SYSLOG>, default=CONSOLE**

Chooses the output target for the Neuron Runtime logs (either console or syslog).

**NEURON\_RT\_LOG\_LEVEL=<ERROR/WARN/INFO/DEBUG/TRACE>, default=ERROR**

Specifies the logging verbosity for the Neuron Runtime library, from ERROR (least verbose), to TRACE (most verbose).

**NEURON\_RT\_NUM\_CORES=<n>**

Specifies how many NeuronCores are needed for the application. During `nrt_init` the requested number of NeuronCores are **exclusively** associated with the calling processes and become unavailable to any other process attempting to use them. If there aren't enough NeuronCores available, `nrt_init` will return an error. Once the owner process has called `nrt_close` or exited, the NeuronCores are released and become available to be associated with another process. By default, all NeuronCores present on the instance will be made available to the caller.

**NEURON\_RT\_VISIBLE\_CORES=<m,n,p-q>**

Similarly to the previous, it allows the calling process to get exclusive access to a set of NeuronCores, but it allows explicitly specifying which NeuronCores are available for the application based on their zero-based indices. This variable can be a list of NeuronCores, for example: `NEURON_RT_VISIBLE_CORES=3,4,5,6`, a range of NeuronCores, for example: `NEURON_RT_VISIBLE_CORES=3-6`, or a combination of both: `NEURON_RT_VISIBLE_CORES=3-5,6`. The resulting range must be contiguous, for example this is not valid: `NEURON_RT_VISIBLE_CORES=3,5,6` because 4 is missing from the list, and indices need to be provided in consecutive increasing order.

---

**Note:** If both `NEURON_RT_VISIBLE_CORES` and `NEURON_RT_NUM_CORES` are defined, `NEURON_RT_VISIBLE_CORES` will be used.

---

**NEURON\_RT\_ROOT\_COMM\_ID=<ip\_address:port>**

Mandatory for applications that run workloads containing Collective Communication operators, allows specifying the IP address and assign a port for the rank 0 worker in the Collective Compute worker pool. For example: `NEURON_RT_ROOT_COMM_ID=10.0.1.2:46820`.

**NEURON\_RT\_STOCHASTIC\_ROUNDING\_SEED=<value>**

Allows setting a value for the stochastic rounding seed. Has no effect on `inf1`.

**NEURON\_RT\_DEBUG\_MEMLOG\_MAX\_SIZE=<value>, default=1024\*1024**

Allows changing the number of entries in the memory allocations log. This log contains an entry for every allocation and deallocation and will be dumped to a file in case of a memory allocation failure in CSV format.

**NRT\_STATUS nrt\_close()**

Closes all the devices used by the application (as defined by `NEURON_RT_NUM_CORES/NEURON_RT_VISIBLE_CORES`) and cleans up the runtime state. Note that once `nrt_close` has been called, most `nrt_*` API calls will fail if attempted.

**The Model API**

**NRT\_STATUS nrt\_load**(const void \*neff\_bytes, size\_t size, int32\_t start\_nc, int32\_t nc\_count, nrt\_model\_t \*\*model)

Loads a NEFF file whose content is found in *neff\_bytes*, with the given size, placing it on `nc_count` NeuronCores starting with NeuronCore index *start\_nc*. If either `nc_count` or `start_nc` are -1, an optimal value for each will be determined automatically. The model can be configured using a list of environment variables read inside this API call which can be found in the [Environment variables used to configure a model being loaded](#) section. It returns a handle to the loaded model in the `nrt_model_t*` pointer if the call succeeds. The returned handle represents the loaded model and can be used with calls that operate on an `nrt_model_t*` (such as `nrt_execute`).

**Parameters**

- **neff\_bytes** – Pointer to existing NEFF file data
- **size** – Size of data in *neff\_bytes*
- **start\_nc** – Index of the NeuronCore on which to stage the model. The first NeuronCore owned by the application will always have the index 0 - for example, even if when setting `NEURON_RT_VISIBLE_CORES=3,4`, the two NeuronCores will be referred to as 0 and 1. If -1, an optimal index will be automatically determined (based on current NeuronCore usage).
- **nc\_count** – Number of NeuronCores on which to stage the model. If its value is a multiple of the amount of NeuronCores needed by the model, the model will be replicated on the number of NeuronCores specified in the argument. This feature is called **TBD** and it will be explained in detail in a separate section. If its value is -1, the model will be staged a single time, using the number of cores needed by a single instance of the model.
- **model** – Model handle returned by the call which can be passed to other functions that operate on models (such as `nrt_execute`).

**Environment variables used to configure a model being loaded**

**NEURON\_RT\_EXEC\_TIMEOUT=<n>, default=30 (inf1), default=600(trn1,inf2)**

Maximum of time, in seconds, allowed for one execution before timing out - which will cause the call to `nrt_execute` to fail and return `NRT_TIMEOUT`.

**NEURON\_RT\_VALIDATE\_HASH=<true/false>, default=false**

Verify the integrity of NEFF data being loaded by checking against a checksum found in the header.

**NEURON\_RT\_STOCHASTIC\_ROUNDING\_EN=<true/false>, default=false**

Enable stochastic rounding.

**NRT\_STATUS nrt\_load\_collectives**(const void \*neff\_bytes, size\_t size, int32\_t start\_nc, int32\_t nc\_count, uint32\_t g\_device\_id, uint32\_t g\_device\_count, nrt\_model\_t \*\*model)

Same as `nrt_load` (same environment variables can be used to configure the model), but must be used when loading NEFFs containing Collective Communication operators. Uses the same arguments as *nrt\_load*, but adds 2 extra ones.

**Parameters**

- **neff\_bytes** – Pointer to existing NEFF file data
- **size** – Size of data in **neff\_bytes**
- **start\_nc** – Index of NeuronCore on which to stage the model. If -1, an optimal index will be automatically determined (based on current NeuronCore usage).
- **nc\_count** – Number of NeuronCores on which to stage the model. If its value is a multiple of the amount of NeuronCores needed by the model, the model will be replicated on the number of NeuronCores specified in the argument. This feature is called **TBD** and it will be explained in detail in a separate section. If its value is -1, the model will be staged a single time, using the number of cores needed by a single instance of the model.
- **g\_device\_id** – Globally unique ID within the Collective Communication world associated with this model instance.
- **g\_device\_count** – Size of the Collective Communication world (total number of participating unique IDs).
- **model** – Model handle returned by the call which can be passed to other functions that operate on models (such as **nrt\_execute**).

**NRT\_STATUS nrt\_unload**(nrt\_model\_t \*model)

Unloads the given model and frees up device and host resources.

#### Parameters

- **model** – Pointer to model to unload. All data associated with the model is deleted, do not reuse the pointer or try to deallocate it afterwards. Do not call **nrt\_unload** again on the same **nrt\_model\_t\*** pointer (think of it as a call to **free()**).

**NRT\_STATUS nrt\_get\_model\_nc\_count**(const nrt\_model\_t \*model, uint32\_t \*nc\_count)

Gets the number of NeuronCores used by the model and writes that value at the address pointed by **nc\_count**.

#### Parameters

- **model** – Valid pointer to an **nrt\_model\_t**.
- **nc\_count** – If the call completes successfully, the pointed address will contain the number of NeuronCores used by the model.

**NRT\_STATUS nrt\_get\_model\_tensor\_info**(nrt\_model\_t \*model, nrt\_tensor\_info\_array\_t \*\*tensor\_info)

Gets input/output tensor information for a given loaded model.

#### Parameters

- **model** – Valid pointer to an **nrt\_model\_t**.
- **tensor\_info** – Pointer to a **nrt\_tensor\_info\_array\_t\*** which will contain the tensor information data. The function allocates memory for the structure internally which can only be correctly freed by calling **nrt\_free\_model\_tensor\_info**. The **nrt\_tensor\_info\_array\_t** struct and its dependencies are defined as follows:

```
typedef struct nrt_tensor_info_array {
    uint64_t tensor_count;           // Total number of input/
    //output tensors used by the model
    nrt_tensor_info_t tensor_array[]; // Array of tensor info_
    //representing those tensors
} nrt_tensor_info_array_t;

typedef struct nrt_tensor_info {
```

(continues on next page)

(continued from previous page)

```

    char name[NRT_TENSOR_NAME_MAX];    // Name of the tensor
    nrt_tensor_usage_t usage;           // Type of the tensor
    size_t size;                        // Tensor size in bytes
    nrt_dtype_t dtype;                  // Data type
    uint32_t *shape;                    // An array representing
↳ data shape
    uint32_t ndim;                      // The number of dimensions
↳ (number of elements in the shape array)
} nrt_tensor_info_t;

// Usage type definitions for tensors
typedef enum nrt_tensor_usage {
    NRT_TENSOR_USAGE_INPUT = 0,        // Tensor is used for input
    NRT_TENSOR_USAGE_OUTPUT,           // Tensor is used for output
} nrt_tensor_usage_t;

// Data type definitions for tensors
typedef enum nrt_dtype {
    NRT_DTYPE_UNKNOWN = 0,
    NRT_DTYPE_FLOAT32,
    NRT_DTYPE_FLOAT16,
    NRT_DTYPE_BFLOAT16,
    NRT_DTYPE_INT8,
    NRT_DTYPE_UINT8,
    NRT_DTYPE_INT16,
    NRT_DTYPE_UINT16,
    NRT_DTYPE_INT32,
    NRT_DTYPE_UINT32,
    NRT_DTYPE_INT64,
    NRT_DTYPE_UINT64
} nrt_dtype_t;

```

NRT\_STATUS **nrt\_free\_model\_tensor\_info**(nrt\_tensor\_info\_array\_t \*tensor\_info)

Frees a `nrt_tensor_info_array_t` allocated by a call to `nrt_get_model_tensor_info`. As with all deallocation functions, don't call it more than once on the same pointer.

#### Parameters

- **tensor\_info** – `nrt_tensor_info_array_t` to deallocate.

NRT\_STATUS **nrt\_get\_model\_instance\_count**(nrt\_model\_t \*model, uint32\_t \*instance\_count)

Returns the number of times this `nrt_model_t` is currently staged on the *NeuronDevice(s)* by writing it to the address pointed by `instance_count`. It will always be  $\geq 1$ . This value can be used to determine the number of threads that can optimally call `nrt_execute` on this `nrt_model_t`.

#### Parameters

- **model** – Valid pointer to an `nrt_model_t`.
- **instance\_count** – If the call completes successfully, the address will contain the instance count for this model



## The Tensor API

**NRT\_STATUS nrt\_tensor\_allocate**(nrt\_tensor\_placement\_t tensor\_placement, int logical\_nc\_id, size\_t size, const char \*name, nrt\_tensor\_t \*\*tensor)

Allocates a new tensor, placing it in either host virtual memory or device memory (based on the `tensor_placement` argument), on the specified NeuronCore index, of a given size, and attaches the given name to it - the name is only used for log messages. For applications running on Inferentia, `tensor_placement` should always be `NRT_TENSOR_PLACEMENT_VIRTUAL`. For all other cases, `NRT_TENSOR_PLACEMENT_DEVICE` should be used. If successful, the `tensor` address will contain a valid pointer to the newly allocated `nrt_tensor_t`. (deprecated) `tensor_placement` set to `NRT_TENSOR_PLACEMENT_HOST` will allocate tensors in physical host memory. Tensors allocated with `NRT_TENSOR_PLACEMENT_HOST` cannot be larger than 4MB, the Kernel physical page size limit. We restrict tensors to a single page of host memory to simplify the generation of DMA descriptors during pre-execution setup.

### Parameters

- **tensor\_placement** – Controls where the tensor will be placed, the definition of the `nrt_tensor_placement_t` enum is as follows:

```
typedef enum {
    NRT_TENSOR_PLACEMENT_DEVICE,    // the tensor is allocated
    ↪directly in device memory
    NRT_TENSOR_PLACEMENT_HOST,      // (deprecated) the tensor is
    ↪allocated in DMAable host memory (only for sizes < 4MB)
    NRT_TENSOR_PLACEMENT_VIRTUAL    // the tensor is allocated in
    ↪host memory
} nrt_tensor_placement_t;
```

- **logical\_nc\_id** (int) – Zero-based NeuronCore index on which to allocate the tensor (if `tensor_placement` is `NRT_TENSOR_PLACEMENT_DEVICE`) or to which associate the tensor for all other cases.
- **size** – Size for the new tensor.
- **name** – Name for the new tensor.
- **tensor** – If the call completes successfully, the address will contain a valid `nrt_tensor_t*` pointer.

**void nrt\_tensor\_free**(nrt\_tensor\_t \*\*tensor)

Frees a tensor allocated by a call to `nrt_tensor_allocate` and sets the `nrt_tensor_t*` pointer at address `tensor` to NULL.

### Parameters

- **tensor** – Pointer to a pointer to a previously allocated `nrt_model_t`. After the call returns, the `nrt_model_t*` pointer will be NULL.

**NRT\_STATUS nrt\_tensor\_read**(const nrt\_tensor\_t \*tensor, void \*buf, size\_t offset, size\_t size)

Reads `size` bytes of data from a given tensor, starting at `offset`, to `buf` starting at offset 0. `buf` needs to be allocated with a size of at least `size` bytes.

### Parameters

- **tensor** – Valid pointer to an `nrt_tensor_t`.
- **buf** – Buffer where to write read data, it needs to be at least `size` bytes in size.
- **offset** – Offset within the tensor from which to begin reading.

- **size** – Size to read.

`NRT_STATUS nrt_tensor_write(nrt_tensor_t *tensor, const void *buf, size_t offset, size_t size)`

Writes `size` bytes of data to a given tensor, starting at `offset`, from `buf` (starting at offset 0).

#### Parameters

- **tensor** – Valid pointer to an `nrt_tensor_t`.
- **buf** – Buffer containing `size` bytes of data to write to the tensor.
- **offset** – Offset within the tensor from which to begin writing.
- **size** – Size to write.

`size_t nrt_tensor_get_size(const nrt_tensor_t *tensor)`

Returns the size, in bytes, of the given tensor.

#### Parameters

- **tensor** – Valid pointer to an `nrt_tensor_t`.

#### Returns

Size in bytes of the given tensor.

`NRT_STATUS nrt_tensor_allocate_empty(const char *name, nrt_tensor_t **tensor)`

Allocates an empty tensor, i.e. the tensor structure w/o any attached storage.

#### Parameters

- **name** – Name for the new tensor.
- **tensor** – If the call completes successfully, the address will contain a valid `nrt_tensor_t*` pointer.

`NRT_STATUS nrt_tensor_attach_buffer(nrt_tensor_t *tensor, void *buffer, size_t size)`

Attaches a caller-supplied buffer to a tensor. Any storage previously attached to the tensor is detached and freed if was owned by the tensor. The attached buffer is managed by the caller and must persist through the entire lifetime of the tensor - calling `nrt_tensor_free` will not deallocate it. This changes the memory placement of the `nrt_tensor_t` to `NRT_TENSOR_PLACEMENT_VIRTUAL` regardless of the initial memory placement type.

#### Parameters

- **tensor** – Valid pointer to an `nrt_tensor_t`.
- **buffer** – Buffer of `size` bytes to attach to the tensor.
- **size** – Size of attached buffer.

`NRT_STATUS nrt_tensor_allocate_slice(const nrt_tensor_t *tensor_source, size_t offset, size_t size, const char *name, nrt_tensor_t **tensor_slice)`

Allocates a new `nrt_tensor_t` that doesn't have its own backing storage - instead, it will use a part (slice) of `tensor_source`'s storage, starting at `offset` with the given size. The shared backing storage is reference counted and it will not be deallocated until the last tensor using it is deallocated.

#### Parameters

- **tensor\_source** – Valid pointer to a `nrt_tensor_t` whose storage will be used by the new tensor.
- **offset** – Offset within the `tensor_source` used as origin for the 'slice'.
- **size** – Size of storage to be used by the new tensor.
- **name** – Name for the new tensor.

- **tensor\_slice** – If the call completes successfully, the address will contain a valid, newly allocated, `nrt_tensor_t*` pointer.

void **nrt\_tensor\_get\_va**(const `nrt_tensor_t` \*tensor)

Returns the virtual address for an allocated tensor.

#### Parameters

- **tensor** – Valid pointer to an `nrt_tensor_t`.

#### Returns

Pointer to host memory used by the tensor.

## The Tensorset API

Tensorsets are containers for tensors.

NRT\_STATUS **nrt\_allocate\_tensor\_set**(`nrt_tensor_set_t` \*\*result)

Allocates an empty `nrt_tensor_set_t` and places its address in `result`.

#### Parameters

- **result** – If the call completes successfully, this address will contain a pointer to a valid, newly allocated `nrt_tensor_set_t`.

void **nrt\_destroy\_tensor\_set**(`nrt_tensor_set_t` \*\*tensor\_set)

Frees a tensor set allocated by a call to `nrt_allocate_tensor_set` and sets the `nrt_tensor_set_t*` pointer at address `tensor_set` to NULL.

#### Parameters

- **tensor\_set** – Pointer to a pointer to a previously allocated `nrt_tensor_set_t`. After the call returns, the `nrt_tensor_set_t*` pointer will be NULL.

NRT\_STATUS **nrt\_add\_tensor\_to\_tensor\_set**(`nrt_tensor_set_t` \*tensor\_set, const char \*tensor\_name, `nrt_tensor_t` \*tensor)

Adds an `nrt_tensor` to a `tensor_set` under a given name. That name can be later used to retrieve the tensor.

#### Parameters

- **tensor\_set** – Pointer to a valid Tensorset where to add the tensor.
- **tensor\_name** – Name that will be used to access the added tensor in the container. Does not need to be the same as the `nrt_tensor_t`'s name.
- **tensor** – Pointer to a valid `nrt_tensor_t` to add to the Tensorset.

NRT\_STATUS **nrt\_get\_tensor\_from\_tensor\_set**(`nrt_tensor_set_t` \*tensor\_set, const char \*tensor\_name, `nrt_tensor_t` \*\*tensor)

Gets an `nrt_tensor` from the tensor set based on the name used when it was added by `nrt_add_tensor_to_tensor_set` and places its address at the address pointed by `tensor`. If the tensor is not found, `NRT_FAILURE` is returned and nothing gets written at the address pointed by `tensor`.

#### Parameters

- **tensor\_set** – Pointer to a valid Tensorset containing the tensor.
- **tensor\_name** – Name associated with the searched `nrt_tensor_t` when it was added to this Tensorset. Might be different from the `nrt_tensor_t`'s internal name.
- **tensor** – Address where the address of the found `nrt_tensor_t` will be placed.

## The Execution API

`NRT_STATUS nrt_execute(nrt_model_t *model, const nrt_tensor_set_t *input_set, nrt_tensor_set_t *output_set)`

Runs one execution of the given `nrt_model_t` using the provided input tensor set and writing the results to the provided output tensor set.

### Parameters

- **model** – Valid pointer to a `nrt_model_t` on which to run the execution.
- **input\_set** – Tensor set containing input data.
- **output\_set** – Tensor set where the output data will be written to.

`NRT_STATUS nrt_execute_repeat(nrt_model_t *model, const nrt_tensor_set_t *input_set, nrt_tensor_set_t *output_set, int repeat_count)`

Same as `nrt_execute` but it will repeat the execution `repeat_count` times using the outputs from the  $n - 1$ th iteration as inputs for the  $n$ th iteration. This requires a specially compiled NEFF and it's not a commonly used call.

### Parameters

- **model** – Valid pointer to a `nrt_model_t` on which to run the execution.
- **input\_set** – Tensor set containing input data.
- **output\_set** – Tensor set where the output data will be written to.
- **repeat\_count** – Number of times to repeat this execution.

## The Profiling API

`NRT_STATUS nrt_profile_start(nrt_model_t *model, const char *filename)`

Begins profiling of the execution of the given model. The profile data will be written to the file specified by the path in `filename`. The file will be truncated if it exists.

### Parameters

- **model** – Valid pointer to a `nrt_model_t` which will be profiled by the Neuron Runtime during execution.
- **filename** – Path to a file where the profile will be written. If the file already exists, it will be truncated.

`NRT_STATUS nrt_profile_stop(const char *filename)`

Ends profiling of the execution of a model and writes profile data to `filename`. `filename` needs to be the same path as the one used for `nrt_profile_start`.

### Parameters

- **filename** – Path to a file where the profile will be written. If the file already exists, it will be truncated.

## Other APIs

NRT\_STATUS **nrt\_get\_version**(nrt\_version\_t \*ver, size\_t size)

Fills a `nrt_version_t` struct with the provided size with version info. The `size` argument allows for backwards compatibility. if the struct changes in future releases.

### Parameters

- **\*ver** – Pointer to a `nrt_version_t` structure which is currently defined as:

```
typedef struct nrt_version {
    uint64_t rt_major;           // major version number
    uint64_t rt_minor;          // minor version number
    uint64_t rt_patch;           // patch version number
    uint64_t rt_maintenance;     // maintenance version number
    char rt_detail[RT_VERSION_DETAIL_LEN]; // runtime version_
    ↪description string
    char git_hash[GIT_HASH_LEN]; // runtime git hash
} nrt_version_t;
```

- **size** (size\_t) – Size of the `nrt_version_t` structure, should always be `sizeof(nrt_version_t)`

NRT\_STATUS **nrt\_get\_total\_nc\_count**(uint32\_t \*nc\_count)

Gets the total number of NeuronCores present on the current instance. The result is not affected by the values in `NEURON_RT_NUM_CORES` or `NEURON_RT_VISIBLE_CORES` and, in fact, this function can be called before calling `nrt_init`.

### Parameters

- **nc\_count** – If the call completes successfully, the address will contain the total number of NeuronCores present on the instance.

NRT\_STATUS **nrt\_get\_visible\_nc\_count**(uint32\_t \*nc\_count)

Gets the total number of NeuronCores available to the application after `nrt_init` has parsed the configuration environment variables `NEURON_RT_NUM_CORES` and `NEURON_RT_VISIBLE_CORES` (if provided).

### Parameters

- **nc\_count** – If the call completes successfully, the address will contain the total number of NeuronCores available to the application.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.1.2 Configuration Guide

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## NeuronX Runtime Configuration

NeuronX Runtime is responsible for executing ML models on Neuron Devices. NeuronX Runtime determines which NeuronCore will execute which model and how to execute it. Configuration of the NeuronX Runtime is controlled through the use of Environment variables at the process level. By default, Neuron framework extensions will take care of NeuronX Runtime configuration on the user's behalf. Explicit configurations are also possible when attempting to achieve a desired behavior.

This guide provides an overview of the different environment variables available to configure NeuronX Runtime behavior.

Table 6.1: Environment Variables

| Name                                       | Description   | Type                     | Expected Values   | Default Value                    | RT Version |
|--|---|--------------------------|---|----------------------------------|------------|
| NEURON_RT_NUM_CORES                        | Number of NeuronCores needed by the process   | Integer range (like 1-3) | Any value or range between 0 to Max NeuronCore in the system. | None                             | 2.0+       |
| NEURON_RT_NUM_CORES                        | Number of NeuronCores required by the process.  | Integer                  | A value from 1 to Max NeuronCore in the system.               | 0, which is interpreted as "all" | 2.0+       |
| NEURON_RT_LOG_LOCATION                     | Log file location   | string                   | console or syslog   | console                          | 2.0+       |
| NEURON_RT_LOG_LEVEL                        | Log verbose level   | string                   | ERROR, WARNING, INFO, DEBUG, TRACE                            | ERROR                            | 2.0+       |
| NEURON_RT_EXEC_TIMEOUT                     | Execution in seconds  | Integer                  | 0 to INT_MAX  | 30                               | 2.0+       |
| NEURON_RT_VALIDATE_HASH                    | Validate contents before loading into accelerator                                     | Boolean                  | TRUE or FALSE   | FALSE                            | 2.0+       |
| NEURON_RT_SHUFFLE_INSTANCE_SHARED_WEIGHTS  | Shuffle between multiple instance versions of the same model on different NeuronCores | Boolean                  | TRUE or FALSE   | FALSE                            | 2.11+      |
| NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS | Maximum number of asynchronous execution requests to be supported.                    | Integer                  | 0 to INT_MAX; 0 is disabled.                                  | 0                                | 2.15+      |

## NeuronCore Allocation

**Important:** NEURONCORE\_GROUP\_SIZES is being deprecated, if your application is using NEURONCORE\_GROUP\_SIZES please see [Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#) for more details.

By default, NeuronX Runtime initializes all the cores present in the system and reserves them for the current process.

---

**Note:** Once a NeuronCore is reserved for a process, it cannot be used by another process at all, until the process reserving that NeuronCore is terminated.

---

## Using NEURON\_RT\_VISIBLE\_CORES

For parallel processing, NEURON\_RT\_VISIBLE\_CORES can be used to control which NeuronCores each process would reserve. This variable is specified with a single NeuronCore index or an inclusive range value.

For example, if a process (myapp.py) requires one NeuronCore, then it can be started with NEURON\_RT\_VISIBLE\_CORES=0 to limit the process to NeuronCore 0. For parallel processing, multiple process can be started (without any change to myapp.py code) with different NEURON\_RT\_VISIBLE\_CORES values. Here is an example that runs myapp.py on inf1.xlarge in parallel across the four different NeuronCores available in the inf1.xlarge.

```
NEURON_RT_VISIBLE_CORES=0 myapp.py &
NEURON_RT_VISIBLE_CORES=1 myapp.py &
NEURON_RT_VISIBLE_CORES=2 myapp.py &
NEURON_RT_VISIBLE_CORES=3 myapp.py &
```

If myapp.py required 3 NeuronCores and was running on a inf1.6xlarge (16 NeuronCores maximum), the first instance of myapp.py could use NeuronCores 0-2, the next instance could use 3-5 and so on:

```
NEURON_RT_VISIBLE_CORES=0-2 myapp.py &
NEURON_RT_VISIBLE_CORES=3-5 myapp.py &
NEURON_RT_VISIBLE_CORES=6-8 myapp.py &
NEURON_RT_VISIBLE_CORES=9-11 myapp.py &
NEURON_RT_VISIBLE_CORES=12-14 myapp.py &
```

## Using NEURON\_RT\_NUM\_CORES

If NEURON\_RT\_NUM\_CORES is set to a value between 1 and the maximum number of NeuronCores in the instance, Neuron Runtime will attempt to automatically reserve the number of free NeuronCores specified for the process. The difference between NEURON\_RT\_VISIBLE\_CORES and NEURON\_RT\_NUM\_CORES is that, NEURON\_RT\_VISIBLE\_CORES specifies exact NeuronCores to allocate where as NEURON\_RT\_NUM\_CORES specifies the number of NeuronCores needed and Neuron Runtime selects free NeuronCores.

Using the same example earlier where myapp.py needed 3 cores, but \_which\_ 3 cores was of no concern, the same application could be executed in parallel up to 5 times on an inf1.6xlarge (16 NeuronCore max):

```
NEURON_RT_NUM_CORES=3 myapp.py &
NEURON_RT_NUM_CORES=3 myapp.py &
NEURON_RT_NUM_CORES=3 myapp.py &
NEURON_RT_NUM_CORES=3 myapp.py &
NEURON_RT_NUM_CORES=3 myapp.py &
```

Executing a 6th NEURON\_RT\_NUM\_CORES=3 myapp.py & in the above example would fail as there is only a single NeuronCore still free.

## Notes

1. Number of NeuronCores in a inferentia device is 4
2. Number of inferentia is depends on the instance size.
3. The NeuronCore index in NEURON\_RT\_VISIBLE\_CORES starts from 0 and ends at (number of NeuronDevices \* number of NeuronCores) - 1.
4. By default, NEURON\_RT\_NUM\_CORES is set to 0, which indicates to RT that all cores are to be used.
5. NEURON\_RT\_VISIBLE\_CORES takes precedence over NEURON\_RT\_NUM\_CORES. If specified, all cores within the range will be assigned to the owning process.

## Logging and debug-ability

By default, NeuronX Runtime logs to syslog with verbose level of *INFO* and only *ERROR* s are logged in console. The following code snippet shows ways to increase/decrease the log level.

```
NEURON_RT_LOG_LEVEL=INFO myapp.py           # Sets the log level for syslog and console to
↪ INFO
NEURON_RT_LOG_LOCATION=console NEURON_RT_LOG_LEVEL=QUIET myapp.py    # Completely
↪ disables console logging.
```

By default, NeuronX Runtime expects the NeuronCore to complete execution of any model with in 2 seconds. If NeuronCore didn't complete the execution within 2 seconds then runtime would fail the execution with timeout error. Most of the models takes few milliseconds to complete so 2 seconds(2000 milliseconds) is more than adequate. However if your model is expected to run more than 2 seconds then you can increase the timeout with NEURON\_RT\_EXEC\_TIMEOUT.

```
NEURON_RT_EXEC_TIMEOUT=5 myapp.py           # increases the timeout to 5 seconds
```

## Additional Logging Controls

NeuronX Runtime enables detailed control over logging behaviors, including the ability to set separate log levels and log locations for individual components. When NEURON\_RT\_LOG\_LEVEL is set globally, NeuronX Runtime combines the logs from all modules into a single stream. For instance, the logs from the modules TDRV and NMGR would appear in the same stream as shown in the example below

```
::
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_infer_status_notifications (FATAL-
RT-UNDEFINED-STATE) inference timeout (600000 ms) on Neuron Device 0 NC 0, waiting for execution
completion notification 2023-Jan-09 20:27:41.0600 15042:15042 ERROR NMGR:dlr_infer
```

However, it is possible to adjust the log level for individual components to capture more or less detail as required for specific debugging contexts. These individual components are - TDRV: the low level driver library - NMGR: the higher level manager library bridging the driver and runtime - NRT: the Neuron Runtime library responsible for loading and executing models that is exposed to end users and frameworks

To adjust the log level for individual components, use the environment variable NEURON\_RT\_LOG\_LEVEL\_<component>, where <component> is the identifier of the component (either TDRV, NMGR, or NRT). This allows for precise control over the verbosity of logs generated by each component, facilitating more targeted debugging. For example, the following sets different log levels for the TDRV and NMGR components.

```
::
export NEURON_RT_LOG_LEVEL_TDRV=DEBUG export NEURON_RT_LOG_LEVEL_NMGR=ERROR
```



Similarly, to specify separate log locations for individual components, use the environment variable `NEURON_RT_LOG_LOCATION_<component>`, following the same naming convention as for log levels. This feature enables logs from different components to be directed to separate files or destinations, making it easier to organize and analyze the log output. For example, the following sets different log locations for the TDRV and NMGR components.

```
::
    export NEURON_RT_LOG_LOCATION_TDRV=tdrv.log export NEURON_RT_LOG_LOCATION_NMGR=nmgr.log
```

## Checksum

To execute a model (NEFF), NeuronX Runtime needs to load the NEFF file into NeuronCore and run. Neuron Runtime provides a way to do checksum validation on each NEFF file while loading to validate the file is not corrupted. This option is off by default to avoid performance penalty during model load time (~50%).

```
NEURON_RT_VALIDATE_HASH=true myapp1.py      # enables model checksum validation while
↪ loading
NEURON_RT_VALIDATE_HASH=false myapp2.py     # disables(default) model checksum validation
↪ while loading
```

## Shared Weights (NEURON\_RT\_MULTI\_INSTANCE\_SHARED\_WEIGHTS)

By default, NeuronX Runtime will make copies of model weights when loading the same instance of a model to multiple NeuronCores. Changing this default to a weight sharing mechanism is possible with NeuronX Runtime 2.11 or higher by setting `NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS=TRUE`. Use of this flag will allow for more models to be loaded by reducing the memory requirements, but will potentially come at a cost of throughput by forcing the execution across cores to compete for memory bandwidth.

Note: the use of this flag requires the model to be loaded with the multi-instance feature (see *PyTorch Neuron (torch-neuron) Core Placement API [Beta]*).

See the [\[BERT tutorial with shared weights notebook\]](#) for an example of how this is used in Torch-Neuron.

```
NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS=TRUE myapp1.py      # enables model weight sharing
NEURON_RT_MULTI_INSTANCE_SHARED_WEIGHTS=FALSE myapp2.py     # disables(default) model
↪ weight sharing
```

## Asynchronous Execution (NEURON\_RT\_ASYNC\_EXEC\_MAX\_INFLIGHT\_REQUESTS)

A beta asynchronous execution feature which can reduce latency by roughly 12% for training workloads. Starting in Neuron Runtime version 2.15, the feature is available, but disabled. To enable the feature for possible improvement, recommendation is to set `NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS` to 3. Setting the number of inflight requests above 3 may lead to Out-Of-Memory (OOM) errors during execution. For developers using `libnrt.so` directly, please use `nrt_register_async_exec_callback` to register a callback for the `nrt` execution thread to post the execution status to. A default callback will be registered if one is not set by the developer.

```
NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS=3 myapp.py      # Up to 3 async exec requests
↪ at once.
NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS=0 myapp.py     # disables async execution
↪ (default behavior)
```

This document is relevant for: Inf1, Inf2, Trn1, Trn2

This document is relevant for: Inf1, Inf2, Trn1, Trn2

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.1.3 Misc (NeuronX Runtime)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### Neuron Runtime Troubleshooting on Inf1, Inf2 and Trn1

This document aims to provide more information on how to fix issues you might encounter while using the Neuron Runtime 2.x or above. For each issue we will provide an explanation of what happened and what can potentially correct the issue.

If your issue is not listed below or you have a more nuanced problem, contact us via [issues](#) posted to this repo, the [AWS Neuron developer forum](#), or through AWS support.

##### Table of contents

- *Generic Errors*
  - *Neuron Driver installation fails*
  - *Application fails to start*
  - *This Neuron Runtime (compatibility id: X) is not compatible with the installed aws-neuron-dkms package*
  - *Neuron Core is in use*
  - *Unsupported NEFF Version*
  - *Unsupported Hardware Operator Code*
  - *Insufficient Memory*
  - *Insufficient number of NeuronCores*
  - *Numerical Error*
  - *RuntimeError: module compiled against API version 0xf but this version of numpy is 0xe*
  - *Failure to initialize Neuron*
  - *An application is trying to use more cores that are available on the instance*
  - *Neuron Runtime execution fails at out-of-bound access*
- *Hardware Errors*
- *EFA and Collective Communication Errors*
  - *Missing aws-neuronx-collectives package*
  - *Missing efa installer package.*
  - *EFA is not enabled in trn1.32xlarge*
  - *Communication timeout*
  - *Communication errors.*
  - *EFA Kernel messages (dmesg) after process termination.*
  - *Failure to find bootstrap interface*

- *Name resolution failure*
- *Usage of Neuron Custom C++ Operators*
  - *Neuron Runtime timeout or GPUSIMD exception*
  - *FI\_EFA\_FORK\_SAFE*

## Generic Errors

### Neuron Driver installation fails

aws-neuron-dkms is a driver package which needs to be compiled during installation. The compilation requires kernel headers for the instance's kernel. `uname -r` can be used to find kernel version in the instance. In some cases, the installed kernel headers might be newer than the instance's kernel itself.

Please look at the aws-neuron-dkms installation log for message like the following:

```
Building for 4.14.193-149.317.amzn2.x86_64
Module build for kernel 4.14.193-149.317.amzn2.x86_64 was skipped since the
kernel headers for this kernel does not seem to be installed.
```

If installation log is not available, check whether the module is loaded.

```
$ lsmod | grep neuron
```

If the above has no output then that means aws-neuron-dkms installation is failed.

### Solution

1. Stop all applications using the NeuronCores.
2. Uninstall aws-neuron-dkms `sudo apt remove aws-neuron-dkms` or `sudo yum remove aws-neuron-dkms`
3. Install kernel headers for the current kernel `sudo apt install -y linux-headers-$(uname -r)` or `sudo yum install -y kernel-devel-$(uname -r) kernel-headers-$(uname -r)`
4. Install aws-neuron-dkms `sudo apt install aws-neuron-dkms` or `sudo yum install aws-neuron-dkms`

### Application fails to start

Neuron Runtime requires Neuron Driver(aws-neuron-dkms package) to access Neuron devices. If the driver is not installed then Neuron Runtime won't be able to access the Neuron devices and will fail with an error message in console and syslog.

If aws-neuron-dkms is not installed then the error message will be like the following:

```
2021-Aug-11 18:38:27.0917 13713:13713 ERROR NRT:nrt_init Unable to determine
Neuron Driver version. Please check aws-neuron-dkms package is installed.
```

If `aws-neuron-dkms` is installed but does not support the latest runtime then the error message will be like the following:

```
2021-Aug-11 19:18:21.0661 24616:24616 ERROR   NRT:nrt_init      This runtime requires_
↳ Neuron Driver version 2.0 or greater. Please upgrade aws-neuron-dkms package.
```

When using any supported framework from Neuron SDK version 2.5.0 and Neuron Driver (`aws-neuron-dkms`) versions 2.4 or older, Neuron Runtime will return the following error message:

```
2022-Dec-01 09:34:12.0559 138:138  ERROR   HAL:aws_hal_tpb_pooling_write_profile
↳ failed programming the engine
```

### Solution

Please follow the installation steps in *Setup Guide* to install `aws-neuronx-dkms`.

---

### This Neuron Runtime (compatibility id: X) is not compatible with the installed aws-neuron-dkms package

This error is caused by incompatibility between the Neuron Driver (`dkms` package) and the Runtime Library (`runtime-lib` package). The driver remains backwards compatible with older versions of Neuron Runtime, but newer versions of the Runtime might rely on the functionality that is only provided by a newer driver. In that case, an update to the newer driver is required.

In some cases the compatibility error persists even after the driver has been updated. That happens when the update process fails to reload the driver at the end of the update. Note that `$ modinfo neuron` will misleadingly show the new version because `modinfo` reads the version information for `neuron.ko` file that's been successfully replaced.

Reload failure happens because one of the processes is still using Neuron Devices and thus the driver cannot be reloaded.

### Solution

Check for any process that is still using the Neuron driver by running `lsmod`:

```
ubuntu@ip-10-1-200-50:~$ lsmod | grep neuron
neuron                237568  0
ubuntu@ip-10-1-200-50:~$
```

“Used by” counter, the second number, should be 0. If it is not, there is still a running process that is using Neuron. Terminate that process and either:

```
$ sudo rmmod neuron
$ sudo modprobe neuron
```

Or simply rerun the installation one more time. The driver logs its version in `dmesg`:

```
$ sudo dmesg
...
[21531.105295] Neuron Driver Started with Version:2.9.4.0-
↳ 8a6fdf292607dccc3b7059ebbe2fb24c60dfc7c4
```

A common culprit is a Jupyter process. If you are using Jupyter on the instance, make sure to terminate Jupyter process before updating the driver.

## Neuron Core is in use

A Neuron Core can't be shared between two applications. If an application started using a Neuron Core all other applications trying to use the NeuronCore would fail during runtime initialization with the following message in the console and in syslog:

```
2021-Aug-27 23:22:12.0323 28078:28078 ERROR   NRT:nrt_allocate_neuron_cores
↳ NeuronCore(s) not available - Requested:nc1-nc1 Available:0
```

## Solution

Terminate any other processes that are using NeuronCore and then try launching the application again. If you are using Jupyter, ensure that you only have a single Jupyter kernel attempting to access the NeuronCores by restarting or shutting-down any other kernels, which will release any NeuronCores that might be in use.

## Unsupported NEFF Version

While loading a model(NEFF), Neuron Runtime checks the version compatibility. If the version the NEFF is incompatible with Runtime then it would fail the model load with following error message:

```
NEFF version mismatch supported: 1.1 received: 2.0
```

## Solution

Use compatible versions of Neuron Compiler and Runtime. Updating to the latest version of both Neuron Compiler and Neuron Runtime is the simplest solution. If updating one of the two is not an option, please refer to the neuron-runtime-release-notes of the Neuron Runtime to determine NEFF version support.

## Unsupported Hardware Operator Code

While loading a model(NEFF), Neuron Runtime checks whether the hardware operators are supported or not. If unsupported, Neuron Runtime will display the following error messages:

```
2023-Jul-28 22:23:13.0357 101413:101422 ERROR   TDRV:translate_one_pseudo_instr_v2
↳ Unsupported hardware operator code 214 found in neff.
2023-Jul-28 22:23:13.0357 101413:101422 ERROR   TDRV:translate_one_pseudo_instr_v2
↳ Please make sure to upgrade to latest aws-neuronx-runtime-lib and aws-neuronx-
↳ collective; for detailed installation instructions visit Neuron documentation.
```

## Solution

Upgrade to latest Neuron Runtime and Neuron Collectives.

---

## Insufficient Memory

While loading a model(NEFF), Neuron Runtime reserves both device and host memory for storing weights, ifmap and ofmap of the Model. The memory consumption of each model is different. If Neuron Runtime is unable to allocate memory then the model load would fail with the following message in syslog

```
kernel: [XXXXX] neuron:mc_alloc: device mempool [0:0] total 1073741568 occupied_
↳ 960539030 needed 1272 available 768
```

## Solution

As the error is contextual to what's going on with your instance, the exact next step is unclear. Try unloading some of the loaded models which will free up device DRAM space. If this is still a problem, moving to a larger Inf1 instance size with additional NeuronCores may help.

---

## Insufficient number of NeuronCores

The NEFF requires more NeuronCores than available on the instance.

Check for error messages in syslog similar to:

```
NRT: 26638:26638 ERROR TDRV:db_vtpb_get_mla_and_tpb          Could not find VNC_
↳ id n
NRT: 26638:26638 ERROR NMGR:dlr_kelf_stage                  Failed to create_
↳ shared io
NRT: 26638:26638 ERROR NMGR:stage_kelf_models              Failed to stage_
↳ graph: kelf-a.json to NeuronCore
NRT: 26638:26638 ERROR NMGR:kmgr_load_nn_post_metrics      Failed to load NN:_
↳ xxxxxxxx, err: 2
```

## Solution

The NeuronCores may be in use by models you are not actively using. Ensure you've unloaded models you're not using and terminated unused applications. If this is still a problem, moving to a larger Inf1 instance size with additional NeuronCores may help.

---

## Numerical Error

Neuron Devices will detect any NaN generated during execution and report it. If Neuron Runtime sees NaNs are generated then it would fail the execution request with Numerical Error with the following message:

```
nrttd[nnnnn]: .... Error notifications found on NC .... INFER_ERROR_SUBTYPE_NUMERICAL
```

## Solution

This usually an indication of either error in the model or error in the input.

Report issue to Neuron by posting the relevant details on GitHub [issues](#).

---

## RuntimeError: module compiled against API version 0xf but this version of numpy is 0xe

This usually means that the numpy version used during compilation is different than the one used when executing the model. As of Neuron SDK release 2.15, numpy versions supported in Neuron SDK are following: numpy<=1.25.2, >=1.22.2. Check and confirm the right numpy version is installed and re-compile/execute the model.

---

## Failure to initialize Neuron

```
nd0 nc0 Timestamp program stop timeout (1000 ms)
nd0 nc0 Error while waiting for timestamp program to end on TPB eng 0
nd0 nc0 Failed to stop neuron core
nd0 nc0 Failed to end timestamp sync programs
TDRV not initialized
Failed to initialize devices, error:5
```

Previously executed application left Neuron devices in running state. Reset Neuron devices but reloading Neuron Driver. Note, this is a temporary workaround, future versions of Neuron will reset running devices automatically.

```
sudo rmmmod neuron; sudo modprobe neuron
```

---

## An application is trying to use more cores that are available on the instance

```
Could not open the nd1
```

Use properly sized instance. trn1.32xlarge has 32 Neuron Cores, trn1.2xlarge has 2 Neuron Cores.

---

## Neuron Runtime execution fails at out-of-bound access

When a Neuron Runtime execution encounters an out-of-bound access error, the runtime logs in the stdout console will display one of the following error messages:

```
2024-08-12 18:34:56,116::ERROR: 2024-Aug-12 18:34:56.067150 159612:159612 ERROR TDRV:
↳generate_custom_notification_msg nd0:nc0:h_model.id1107: Received notification.
↳generated at runtime: failed to run embedding table update, due to out-of-bound access.
2024-08-12 18:34:56,116::ERROR: 2024-Aug-12 18:34:56.067151 159602:159602 ERROR TDRV:
↳generate_custom_notification_msg nd0:nc1:h_model.id1109: Received notification.
↳generated at runtime: failed to run scatter/gather (indirect memory copy), due to out-
↳of-bound access.
```

### Cause of the Error

An out-of-bound access error typically indicates that incorrect inputs have been provided to the model.

### How to Debug

To troubleshoot this issue, you need to examine both the High-Level Operation (HLO) and all inputs. Neuron Runtime can automatically dump all inputs in binary format, which can be instrumental in debugging. To enable input dumping for each failed execution, set the following environment variable:

```
export NEURON_RT_DBG_DUMP_INPUTS_ON_ERR=<an NRT_STATUS value>
```

A complete set of NRT\_STATUS can be found under *The LIBNRT API Return Codes*.

Once this variable is set, Neuron Runtime generates a directory in the current working directory for each failed execution at this NRT\_STATUS value. The directory name follows this pattern:

```
input_dump_<runtime_generated_random_number>_h_nn_<runtime_generated_execution_id>
```

Inside each directory, you'll find all the inputs that led to this failure, stored in binary format. Additionally, the model name is saved in a separate file called model\_name.txt within the same directory.

To disable input dump, you can set the environment variable back to 0

```
export NEURON_RT_DBG_DUMP_INPUTS_ON_ERR=0
```

### Example: Debug an out-of-bound access execution

To debug an out-of-bound (OOB) execution, which returns an NRT\_STATUS code of 1006, both HLO and all inputs are required. By setting the NEURON\_RT\_DBG\_DUMP\_INPUTS\_ON\_ERR environment variable to 1006, you can capture the inputs leading to an OOB execution.

For example, when an OOB error occurs, Neuron Runtime creates a directory named input\_dump\_424238335\_h\_nn\_10001. Here, 424238335 is a randomly generated number by Neuron Runtime, and 10001 is the Neuron Runtime generated execution ID. All relevant inputs, labeled from input0 to input14, are saved in binary format within this directory.

```
ubuntu@ip-172-31-53-90:~$ NEURON_RT_DBG_DUMP_INPUTS_ON_ERR=1006 torchrun --nproc_per_
↳node=2 train_torchrun.py
.....
2024-Jun-26 00:32:47.943821 30294:32381 ERROR TDRV:generate_custom_notification_msg
↳ nd0:nc0:h_model.id1001: Received notification generated at runtime: failed to run
↳scatter/gather (indirect memory copy), due to out-of-bound access. isa instruction
↳line number = 11. model name = /home/ubuntu/token-seqlen1280-batch128-FullyUnrolled.
```

(continues on next page)



(continued from previous page)

```

→736.2.0.62758.0a0+44863561.93f365ce40ab99133659.pb.neff
.....
2024-Jun-26 00:32:47.948678 30294:32381 ERROR  NMGR:dlr_infer
→ Inference completed with err: 1006. mode->h_nn=1001, start_nc=0, nc_count=1
2024-Jun-26 00:32:50.801487 30294:32381 ERROR  TDRV:tensor_dump_inputs
→ 15 input tensors were dumped successfully to directory /home/ubuntu/input_dump_
→424238335_h_nn_10001. Model name is /home/ubuntu/token-seqlen1280-batch128-
→FullyUnrolled.736.2.0.62758.0a0+44863561.93f365ce40ab99133659.pb.neff
.....

ubuntu@ip-172-31-53-90:~$ ls -lt
total 3908900
drwxrwxr-x 2 ubuntu ubuntu 4096 Jun 26 00:32 input_dump_424238335_h_nn_10001
.....

ubuntu@ip-172-31-53-90:~$ ls -lt input_dump_424238335_h_nn_10001
total 1405192
-rw-r--r-- 1 ubuntu ubuntu 5242880 Jun 26 00:32 input14.bin
-rw-r--r-- 1 ubuntu ubuntu 5242880 Jun 26 00:32 input13.bin
-rw-r--r-- 1 ubuntu ubuntu 5242880 Jun 26 00:32 input12.bin
-rw-r--r-- 1 ubuntu ubuntu 5242880 Jun 26 00:32 input11.bin
-rw-r--r-- 1 ubuntu ubuntu 13967360 Jun 26 00:32 input10.bin
-rw-r--r-- 1 ubuntu ubuntu 81920 Jun 26 00:32 input8.bin
-rw-r--r-- 1 ubuntu ubuntu 4 Jun 26 00:32 input9.bin
-rw-r--r-- 1 ubuntu ubuntu 4 Jun 26 00:32 input6.bin
-rw-r--r-- 1 ubuntu ubuntu 81920 Jun 26 00:32 input7.bin
-rw-r--r-- 1 ubuntu ubuntu 16777216 Jun 26 00:32 input5.bin
-rw-r--r-- 1 ubuntu ubuntu 131072 Jun 26 00:32 input3.bin
-rw-r--r-- 1 ubuntu ubuntu 13967360 Jun 26 00:32 input4.bin
-rw-r--r-- 1 ubuntu ubuntu 16777216 Jun 26 00:32 input2.bin
-rw-r--r-- 1 ubuntu ubuntu 13967360 Jun 26 00:32 input1.bin
-rw-r--r-- 1 ubuntu ubuntu 1342177280 Jun 26 00:32 input0.bin
-rw-r--r-- 1 ubuntu ubuntu 9 Jun 26 00:32 model_name.txt

ubuntu@ip-172-31-53-96:~$ cat input_dump_424238335_h_nn_10001/model_name.txt
/home/ubuntu/token-seqlen1280-batch128-FullyUnrolled.736.2.0.62758.0a0+44863561.
→93f365ce40ab99133659.pb.neff

```

### Known Limitations

- **HLO Access:** Neuron Runtime does not have direct access to the HLO; it must be deduced from the model name.
- **Partial Input Dumps:** If a Neuron Runtime execution fails and an exception is raised to the Neuron Framework, other ongoing Neuron Runtime executions may be terminated by the Neuron Framework. This means only one set of inputs may be fully captured, while others may be incomplete if terminated prematurely.
  - An input dump folder is considered complete when the `model_name.txt` file is fully written, as Neuron Runtime saves all inputs first and then writes the `model_name.txt` file. So you might find out the folder with the complete set of inputs by searching for the `model_name.txt` file.

### Hardware Errors

For Trn and Inf instances, the following hardware errors are monitored by Neuron Runtime:

| Error Type                  | Description  | Behaviors  | Recommended Actions  |
|-----------------------------|--|--|--|
| SRAM Uncorrectable          | An on-chip SRAM encountered a parity error and produced incorrect results. | <ol style="list-style-type: none"> <li>Instance Retirement Notice: You will receive an <a href="#">EC2 instance retirement notice</a> within 15 minutes of experiencing this message. EKS, EC2 Auto Scaling Groups, and AWS ParallelCluster will react to these retirement notices according to their configured policies, but you can also automate responses to these notices yourself with <a href="#">EventBridge rules</a>.</li> <li>Neuron Runtime Behavior: Neuron Runtime will timeout and exit with NRT_EXEC_COMPLETED_WITH_ERR (1004) or NRT_EXEC_HW_ERR_NC_UE (1202) return code. You will see the following error message in runtime logs from stdout console: (FATAL-RT-UNDEFINED-STATE) [ND 0][NC 0] Uncorrectable memory error is detected, metadata: 0x16. Please terminate or stop/start this instance to prevent future impact from the hardware error.</li> </ol> | <ol style="list-style-type: none"> <li>Replace the EC2 instance by <a href="#">terminating it</a> or <a href="#">stopping then starting it</a>.</li> <li>Utilize <a href="#">Neuron Sysfs</a> and <a href="#">Neuron Monitor</a> to monitor the <code>sram_ecc_uncorrected</code> error counts.</li> </ol> |
| HBM Uncorrectable           | An HBM encountered an uncorrectable error and produced incorrect results.  | <ol style="list-style-type: none"> <li>Instance Retirement Notice: You will receive an <a href="#">EC2 instance retirement notice</a> within 15 minutes of experiencing this message. EKS, EC2 Auto Scaling Groups, and AWS ParallelCluster will react to these retirement notices according to their configured policies, but you can also automate responses to these notices yourself with <a href="#">EventBridge rules</a>.</li> <li>Neuron Runtime Behavior: Neuron Runtime will timeout and exit with NRT_TIMEOUT (5) or NRT_EXEC_HW_ERR_HBM_UE (1201) return code. You will see the following error message in runtime logs from stdout console: (FATAL-RT-UNDEFINED-STATE) Uncorrectable HBM memory error is detected. Execution results may be invalid. Please terminate or stop/start this instance to prevent future impact from the hardware error.</li> </ol>          | <ol style="list-style-type: none"> <li>Replace the EC2 instance by <a href="#">terminating it</a> or <a href="#">stopping then starting it</a>.</li> <li>Utilize <a href="#">Neuron Sysfs</a> and <a href="#">Neuron Monitor</a> to monitor the <code>mem_ecc_uncorrected</code> error counts.</li> </ol>  |
| DMA Aborts                  | A DMA engine encountered an unrecoverable error.                           | Neuron Runtime Behavior: Neuron Runtime will timeout and exit with NRT_TIMEOUT (5) or NRT_EXEC_HW_ERR_DMA_ABORT (1203) return code. You will see the following error messages in runtime logs from stdout console: [MLA 0][NC 0] DMA TX engine 0 is in an abort state or [MLA 0][NC 0] DMA RX engine 0 is in an abort state  | Replace the EC2 instance by <a href="#">terminating it</a> or <a href="#">stopping then starting it</a> .  |
| Hang on Collectives         | Possibly caused by a hardware error on another worker.                     | Neuron Runtime Behavior: Neuron Runtime will timeout and exit with NRT_TIMEOUT (5) or NRT_EXEC_HW_ERR_COLLECTIVES (1200) return code. You will see the following error messages in runtime logs from stdout console: (FATAL-RT-UNDEFINED-STATE) missing collectives status on Neuron Device 0 NC 0, model 0 - suspected hang in collectives operation 0 out of 100   | Search for SRAM Uncorrectable, HBM Uncorrectable, DMA Aborts, and Hang on Compute errors on the other workers, and implement the recommended actions on the affected   |
| <b>6.1. NeuronX Runtime</b> |  |  | worker. Afterwa <b>955</b><br>restart your workload and attempt again.   |

Upon any hardware errors, you should also expect to see the error message like the following in dmesg: `NEURON_HW_ERR=SRAM_UNCORRECTABLE_ERROR instance-id=i-0592464924bd45322 hostname=ip-172-31-61-252 nd-id=0 nc-id=0 serial-num=19fcda00f5ff6eb9 action=TERMINATE_INSTANCE`

## EFA and Collective Communication Errors

### Missing aws-neuronx-collectives package

**aws-neuronx-collectives** package is required to execute Collective Communication on a single instance and across multiple instances.

```
NCCL init error: Error opening libnccom.so, cannot use collective operations! Please set
LD_LIBRARY_PATH to library location. Error: libnccom.so: cannot open shared object
file: No such file or directory
Please make sure to install correct version of aws-neuronx-collectives; for detailed
installation instructions visit Neuron documentation
```

Install aws-neuornx-collectives package. If the installation used non-default destination set LD\_LIBRARY\_PATH.

### Missing efa installer package.

**efa-installer** package is required to execute Collective Communication across multiple instances.

```
Unable to run multi-instance workload. Ofi plugin is not installed or EFA is not enabled
```

Follow the directions to install efa-installer package. Make sure to add the path to to libfabric library to LD\_LIBRARY\_PATH

### EFA is not enabled in trn1.32xlarge

EFA is used as a transport for Collective Communication among multiple instances. EFA must be enabled on the instances used for multi-node training.

```
OFI plugin initNet() failed is EFA enabled?
```

Confirm that EFA is enabled by running `lspci` command and making sure there are eight EFA devices. For example:

```
[ec2-user@ip-10-0-13-247 ~]$ lspci -tv
--[0000:a0]--00.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|      +-01.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|      +-19.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|      +-1a.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|      +-1b.0 Amazon.com, Inc. NeuronDevice
|      +-1c.0 Amazon.com, Inc. NeuronDevice
|      +-1d.0 Amazon.com, Inc. NeuronDevice
|      +-1e.0 Amazon.com, Inc. NeuronDevice
|      \-1f.0 Amazon.com, Inc. NVMe SSD Controller
+--[0000:90]--00.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|      +-01.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|      +-19.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
```

(continues on next page)

(continued from previous page)

```

|          +-1a.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|          +-1b.0 Amazon.com, Inc. NeuronDevice
|          +-1c.0 Amazon.com, Inc. NeuronDevice
|          +-1d.0 Amazon.com, Inc. NeuronDevice
|          +-1e.0 Amazon.com, Inc. NeuronDevice
|          \-1f.0 Amazon.com, Inc. NVMe SSD Controller
+-[0000:20]--00.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|          +-01.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|          +-19.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|          +-1a.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|          +-1b.0 Amazon.com, Inc. NeuronDevice
|          +-1c.0 Amazon.com, Inc. NeuronDevice
|          +-1d.0 Amazon.com, Inc. NeuronDevice
|          +-1e.0 Amazon.com, Inc. NeuronDevice
|          \-1f.0 Amazon.com, Inc. NVMe SSD Controller
+-[0000:10]--00.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|          +-01.0 Amazon.com, Inc. Elastic Network Adapter (ENA)
|          +-19.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|          +-1a.0 Amazon.com, Inc. Elastic Fabric Adapter (EFA)
|          +-1b.0 Amazon.com, Inc. NeuronDevice
|          +-1c.0 Amazon.com, Inc. NeuronDevice
|          +-1d.0 Amazon.com, Inc. NeuronDevice
|          +-1e.0 Amazon.com, Inc. NeuronDevice
|          \-1f.0 Amazon.com, Inc. NVMe SSD Controller
\-[0000:00]--00.0 Intel Corporation 440FX - 82441FX PMC [Natoma]
|          +-01.0 Intel Corporation 82371SB PIIIX3 ISA [Natoma/Triton II]
|          +-01.3 Intel Corporation 82371AB/EB/MB PIIIX4 ACPI
|          +-03.0 Amazon.com, Inc. Device 1111
|          +-04.0 Amazon.com, Inc. NVMe EBS Controller
|          \-1f.0 Amazon.com, Inc. NVMe EBS Controller

```

Launch instances with EFA enabled and try again. If not planning to use the instances for multi-node training or running on trn1.2xlarge, this error message can be ignored.

## Communication timeout

Ranks exchange information during NEFF loading and before the start of the execution. The loading/execution cannot move forward until all ranks are ready.

```
Timeout waiting for RX (waited 120 sec) - retrying
```

```
Timeout waiting for incoming connection (waited 120 sec) - retrying
```

```
Connect to localhost:33666 failed - retrying
```

The communication timeouts are not fatal. The ranks will continue waiting forever. In most case the timeouts are caused by one of the ranks getting delayed, usually by recompilation of a graph. The execution is resumed after the graph is compiled (might take significant amount of time). It is possible to determine if compilation is in progress by checking the logs on all nodes.

Communication timeouts might also indicate that one of the nodes or ranks is hang. If that is the case, terminate the run and restart from the last known good check point.

## Communication errors.

```
RX, connection closed by remote peer
```

There could be other similar messages indicating that ranks failed to communicate.

One of the ranks or nodes encountered a problem and terminated. Terminate the run and restart from the last known check point.

## EFA Kernel messages (dmesg) after process termination.

```
[298850.502143] neuron:npid_detach: neuron:npid_detach: pid=90193, slot=0
[298850.919248] efa 0000:a0:1a:0 rdmap160s26: Failed to process command DEREG_MR (opcode_
↪8) comp_status 7 err -22
```

When a process that executed Collective Communication terminates it deregisters buffers that were registered with the networking stack. There is a race condition because the Neuron driver deregisters buffers owned by terminating process as part of the memory cleanup. The error is benign and will be removed in the future releases.

## Failure to find bootstrap interface

```
No interface found in the same subnet as remote address fe80::1461:22ff:fe33:b471<45015>
No usable listening interface found
```

Bootstrap code incorrectly trying to use link-local IPv6 address for communication. This error will be fixed in the next Neuron release. In the meantime, as a workaround, disable IPv6 on the instances.

```
sudo sysctl -w net.ipv6.conf.all.disable_ipv6=1
sudo sysctl -w net.ipv6.conf.default.disable_ipv6=1
```

## Name resolution failure

```
WARN Invalid NCCL_COMM_ID [compute1-dy-training-0-1.pcluster-trn1-24-pdx80-2n.pcluster:
↪41211], please use format: <ipv4>:<port> or [<ipv6>]:<port>
```

Verify that the name can be resolved by DNS by using nslookup or dig. Currently released version fails to resolve FQDN longer than 63 characters. This error will be fixed in the upcoming Neuron release. In the mean time use shorter names to ensure that FQDN length does not exceed the maximum of 63 characters.

## Usage of Neuron Custom C++ Operators

### Neuron Runtime timeout or GPSIMD exception

At this point, reset of Neuron Runtime is required after running a model which invoked a Neuron Custom C++ operator. Otherwise, a Neuron Runtime timeout or GPSIMD exception may occur.

Example Neuron Runtime timeout:

```

2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_tpb_status_notifications
↳ Missing infer_status notification: (end:1)
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_tpb_status_notifications
↳ Missing infer_status notification: (end:2)
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_tpb_status_notifications
↳ Missing infer_status notification: (end:3)
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_tpb_status_notifications
↳ Missing infer_status notification: (end:4)
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_tpb_status_notifications
↳ Missing infer_status notification: (end:0)
2023-Jan-09 20:27:41.0593 15042:15042 ERROR TDRV:exec_consume_infer_status_
↳ notifications (FATAL-RT-UNDEFINED-STATE) inference timeout (6000000 ms) on Neuron
↳ Device 0 NC 0, waiting for execution completion notification
2023-Jan-09 20:27:41.0600 15042:15042 ERROR NMGR:dlr_infer
↳ Inference completed with err: 5

```

Example GPSIMD exception:

```

2023-Jan-06 22:28:01.0845 137472:137472 ERROR TDRV:pool_stdio_queue_consume_all_entries
↳ Printing stderr from GPSIMD:
GPSIMD EXCEPTION OCCURRED: ILLEGAL INSTRUCTION
Subtype/Type/Cause: 0x201
Exception PC: 0x840001E8

```

## Solution

If either of the above errors are seen, and `NEURON_RT_RESET_CORES` is set to 0, either unset it or set it to 1. This will enable the default runtime behaviour of resetting NeuronCores when initializing applications. See [NeuronX Runtime Configuration](#) for more information.

Also note that the timeout period can be changed by setting `NEURON_RT_EXEC_TIMEOUT`. See [NeuronX Runtime Configuration](#) for more information.

## FI\_EFA\_FORK\_SAFE

Older Linux (<5.15) kernels require environment variable `FI_EFA_FORK_SAFE` to be set to 1 for the libfabric to operate correctly. Specifically Amazon Linux 2 uses 5.10 kernel and requires the variable to be set.

When the variable is not set multi-node collective communication will be disabled. Intra-node collective communication is still possible. The following error message will be logged the first time a model containing collective communication is loaded:

```

Linux kernel 5.10 requires setting FI_EFA_FORK_SAFE=1 environment variable. Multi-node
↳ support will be disabled.
Please restart with FI_EFA_FORK_SAFE=1 set."

```

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## NeuronX runtime FAQ

### Table of Contents

- *Where can I find information about Neuron Runtime 2.x (libnrt.so)*
- *What will happen if I will upgrade Neuron Framework without upgrading latest kernel mode driver?*
- *Do I need to recompile my model to use the Runtime Library?*
- *Do I need to change my application launch command?*
- *How do I restart/start/stop the NeuronX Runtime?*
- *How do I know which runtimes are associated with which Neuron Device(s)?*
- *What about RedHat or other versions of Linux and Windows?*
- *How can I take advantage of multiple NeuronCores to run multiple inferences in parallel?*

### Where can I find information about Neuron Runtime 2.x (libnrt.so)

See *Introducing Neuron Runtime 2.x (libnrt.so)* for detailed information about Neuron Runtime 2.x (libnrt.so).

### What will happen if I will upgrade Neuron Framework without upgrading latest kernel mode driver?

Application start would fail with the following error message: .. code:: bash

```
2021-Aug-11 19:18:21.0661 24616:24616 ERROR NRT:nrt_init This runtime requires Neuron Driver version 2.0 or greater. Please upgrade aws-neuron-dkms package.
```

### Do I need to recompile my model to use the Runtime Library?

No. Runtime 2.x supports all the models compiled with Neuron Compiler 1.x.

### Do I need to change my application launch command?

No.

### How do I restart/start/stop the NeuronX Runtime?

Since Neuron Runtime is a library, starting/stopping application would result in starting/stopping the Neuron Runtime.



## How do I know which runtimes are associated with which Neuron Device(s)?

*neuron-ls* and *neuron-top* can be used to find out applications using Neuron Devices.

## What about RedHat or other versions of Linux and Windows?

We don't officially support it yet.

## How can I take advantage of multiple NeuronCores to run multiple inferences in parallel?

Examples of this for TensorFlow and MXNet are found [here](#) and [here](#).

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## Neuron Runtime Release Notes

Neuron Runtime consists of a kernel mode driver and C/C++ libraries which provides APIs to access Neuron Devices. The runtime itself (libnrt.so) is integrated into the ML frameworks for simplicity of deployment. The Neuron Runtime supports training models and executing inference on the Neuron Cores.

### Table of contents

- *Known issues*
- *NEFF Support Table:*
- *Neuron Runtime Library [2.26.42.0]*
- *Neuron Runtime Library [2.25.57.0]*
- *Neuron Runtime Library [2.24.53.0]*
- *Neuron Runtime Library [2.23.112.0]*
- *Neuron Runtime Library [2.23.110.0]*
- *Neuron Runtime Library [2.22.19.0]*
- *Neuron Runtime Library [2.22.14.0]*
- *Neuron Runtime Library [2.21.41.0]*
- *Neuron Runtime Library [PATCH 2.20.22.0]*
- *Neuron Runtime Library [2.20.11.0]*
- *Neuron Runtime Library [2.19.5.0]*
- *Neuron Runtime Library [2.18.15.0]*
- *Neuron Runtime Library [2.18.14.0]*
- *Neuron Runtime Library [2.17.7.0]*
- *Neuron Runtime Library [2.16.14.0]*
- *Neuron Runtime Library [2.16.8.0]*

- *Neuron Runtime Library [2.15.14.0]*
- *Neuron Runtime Library [2.15.11.0]*
- *Neuron Runtime Library [2.14.8.0]*
- *Neuron Runtime Library [2.13.6.0]*
- *Neuron Runtime Library [2.12.23.0]*
- *Neuron Runtime Library [2.12.14.0]*
- *Neuron Runtime Library [2.11.43.0]*
- *Neuron Runtime Library [2.10.18.0]*
- *Neuron Runtime Library [2.10.15.0]*
- *Neuron Runtime Library [2.9.64.0]*
- *Neuron Runtime 2.x (libnrt.so) release [2.2.51.0]*
- *Neuron Runtime 2.x (libnrt.so) release [2.2.31.0]*
- *Neuron Runtime 2.x (libnrt.so) release [2.2.18.0]*
- *Neuron Runtime 2.x (libnrt.so) release [2.2.15.0]*

## Known issues

Updated : 06/19/2025

- The `nrt_tensor_allocate` APIs do not support more than 4 GB ( $\geq 4\text{GB}$ ) sizes. Passing in a size larger than or equal to 4GB will result in datatype overflow leading to undefined behavior.
- A hardware bug affecting **Trainium** and **Inferentia2** devices causes numerical errors to become “sticky” within the Neuron Core hardware. When a legitimate numerical error occurs during execution, the error state persists in the hardware, causing all subsequent executions to incorrectly report numerical errors even when the computations are valid. This sticky error state can only be resolved by restarting the application to clear the hardware.

## NEFF Support Table:

Use this table to determine the version of Runtime that will support the version of NEFF you are using. NEFF version is determined by the version of the Neuron Compiler.

| NEFF Version | Runtime Version Range | Notes                               |
|--------------|-----------------------|-------------------------------------|
| 0.6          | *                     | All versions of RT support NEFF 0.6 |
| 1.0          | $\geq 1.0.6905.0$     | Starting support for 1.0 NEFFs      |
| 2.0          | $\geq 1.6.5.0$        | Starting support for 2.0 NEFFs      |

## Neuron Runtime Library [2.26.42.0]

Date: 06/24/2025

### New in this release

- Added support for 8x8 collective groups (TP8 + CP8) on **TRN2** for **LNC=2**
- Added support for direct *State-Buffer* to *State-Buffer* collective ops for **LNC=1**
- Introduce RDH algorithm for inter-node collective communication
- Added support for loading NEFF with different world sizes in the same NRT process

### Improvements

- Reduced the average latency of 32x2 collective groups by 65%
- Reduced latency for intra-chip reduce scatter operations on **TRN2** instances by up to 20% for small transfers and 60% for medium to large transfers
- Improved latency for medium message sizes for intra-chip All Gather operations on **TRN2** by up to 60%
- Improved the debugging experience by adding logs which print out the value of timed-out, non-zero semaphores on **Trainium2** platforms
- Improved timeout error messages by displaying the NEFF program counters for the stuck Neuron Core
- Refined out-of-memory error messages to report a NEFF level memory breakdown table

### Bug fixes

- Fixed crash caused by race condition during the capture of system profiles
- Fixed various memory leaks that occur during *nrt\_close*

### Neuron Profiler 2.0 (Beta)

- Added option to filter the Neuron Cores to capture trace events on ([reference](#))
- Added option to filter the event types recorded when capturing system traces ([reference](#))
- Added new trace events to capture the latency of the collectives execution barrier

### Compatibility Changes

- This version of the Neuron runtime requires *aws-neuron-dkms* version 2.22 or later on **Trainium2** instances

### Neuron Runtime Library [2.25.57.0]

Date: 05/19/2025

#### New in this release

- Added `NEURON_RT_LOW_LATENCY_TASKS_CPU_AFFINITY` environment variable to allow users to set the thread affinity of low latency tasks that run on host cpu

#### Improvements

- Refined software notification queue overflow detection flow and improved error message
- Reduced latency for All-Reduce intra-chip collective (TP 4) by 50% for medium message sizes
- Improved error message when an execution request is passed a tensor allocated on an incorrect HBM
- Improved NEFF switch latency by up to 95% when using async mode
- Increased the number of different replica groups supported in the same NEFF on TRN2
- Explicitly limit the max number of in-flight async requests to the hard limit of 63

#### Bug fixes

- Fixed segfault that can occur when applications attempt to load a NEFF with an unsupported number of FMA source descriptors

### Neuron Profiler 2.0 (Beta)

- Added traces for Host <-> device data transfer events in system profiles
- Added pre/post execution hooks to system profiles
- Significant performance improvements in time taken by calls to `nrt_sys_trace_fetch_events()`

### Neuron Runtime Library [2.24.53.0]

Date: 04/03/2025

#### New in this release

- Removed support for Neuron Distributed Event Tracing (reference)

## Improvements

- Improved dynamic DMA descriptor generation performance by up to 3% for certain workloads
- Reduced collectives device memory footprint for large Neffs
- Improved device latency for memory bound workloads on TRN2
- Added support for profiling executions when NRT is launched in Async Execution Mode (`NEURON_RT_ASYNC_EXEC_MAX_INFLIGHT_REQUESTS > 0`)
- Added check to detect execution completion queue overflows

## Bug fixes

- Fixed bug introduced in NRT 2.23 where the runtime was incorrectly reporting executions that hit “Out of Bound” errors as successful executions
- Fixed segfault when encountering “out of memory” errors when starting profiles

## Neuron Profiler 2.0 (Beta)

- Reduced overhead of Neuron Profiler 2.0 to <1% of overall latency
- Added new `nrt_sys_trace_fetch_events` API to retrieve system trace events
- Added out of bound error events to system trace
- Removed the `NEURON_RT_INSPECT_DURATION_NSEC` and `NEURON_RT_INSPECT_START_OFFSET_NSEC` configuration options

## NKI

- Added dynamic DMA support for block scatter ops
- Added RangeSelect instruction Support for the Vector engine

## Neuron Runtime Library [2.23.112.0]

Date: 01/14/2025

## Bug fixes

- Fixed DMA abort errors on TRN2.

### Neuron Runtime Library [2.23.110.0]

Date: 12/20/2024

#### New in this release

- Added Trainium2 support
- Added runtime support to detect and fail on out-of-bound memory access in DMA operations
- Added support for 4-rank replica group on adjacent Neuron cores on TRN1/TRN1N
- Added new profiling API for capturing system and device profiles. This feature is currently in beta. See the Neuron Profiler 2.0 (Beta) documentation for usage. ([reference](#))

#### Improvements

- Reduced runtime host RAM utilization
- Improved Neff context switch overhead reducing latency by up to 500us
- **Split hardware errors into more granular categories**
  - NRT\_EXEC\_HW\_ERR\_HBM\_UE (1201)
  - NRT\_EXEC\_HW\_ERR\_NC\_UE (1202)
  - NRT\_EXEC\_HW\_ERR\_DMA\_ABORT (1203)
- **Updated runtime to breakdown DMA ring memory usage into more detailed categories**
  - dma rings io
  - dma rings spill
  - dma rings collectives
  - dma rings runtime
- Updated the `nrt_load` error path to print a clear error message when failing to load a collectives Neff instead of aborting

#### Bug fixes

- Fixed multiple memory corruptions and exhaustions on the collectives failure path
- Fixed bug where incorrect execution status was passed to the async execution callback

## End of Support

- Removed INF1 Support from Runtime library

## Neuron Runtime Library [2.22.19.0]

Date: 11/20/2024

### New in this release

- Minor improvements and bug fixes

## Neuron Runtime Library [2.22.14.0]

Date: 09/16/2024

### New in this release

- Improved the inter-node mesh algorithm to scale better for larger number of nodes and larger allreduce problem sizes

### Bug fixes

- Implemented a fix that differentiate between out-of-memory (OOM) conditions occurring on the host system versus the device when an OOM event occurs
- Resolved a performance issue with transpose operations, which was caused by an uneven distribution of work across DMA engines

## Neuron Runtime Library [2.21.41.0]

Date: 07/03/2024

### New in this release

- Improved collectives performance on small buffers
- Improved memory utilization by reducing the size of collective buffers
- Logging improvements including improvements for HW errors, out of bounds issues, and collectives
- Added fine grained NRT error return codes for execution errors ([reference](#))

### Bug fixes

- Fixed bug where runtime was incorrectly reporting instruction offsets to the profiler

### Neuron Runtime Library [PATCH 2.20.22.0]

Date: 04/01/2024

### Bug fixes

- Fixed a bug where setting *NEURON\_SCRATCHPAD\_PAGE\_SIZE* to a non-power of two value could lead to unnecessary Neuron memory allocations.
- Fixed messaging so that logs of benign numerical errors do not include a full dump of runtime state.
- Fixed a bug that was causing Neuron Collectives to consume excessive amount of Neuron memory, causing out of memory errors during model load.
- Fixed a bug where the Runtime would fail to report a hardware error while the status API reported instance retirement.
- Fixed a hang in Neuron Collectives that could occur when subgraphs running on different workers had a different number of replicas.

### Neuron Runtime Library [2.20.11.0]

Date: 02/13/2024

### New in this release

- Improved performance of collective communication operators (CC ops) by up to 30% for problem sizes smaller than 16MB. This is a typical size of CC ops when executing LLM inference.
- Added support for inter-node alltoall which is a MoE use case.
- Added NRT version check across all ranks to make sure all ranks are using the same runtime.
- **Improved logging on collectives timeout during model execution.**
  - “(FATAL-RT-UNDEFINED-STATE) missing collectives status on Neuron Device 0 NC 0, model model.neff - suspected hang in collectives operation 0 out of 32”
- **Log HBM uncorrectable errors on timeout if any occurred during model execution.**
  - “(FATAL-RT-UNDEFINED-STATE) encountered uncorrectable memory error on Neuron Device 0, execution results may be invalid. Please terminate or start/stop this instance to recover from bad hardware.”



## Bug fixes

- Fixed bug where metrics were undercounting the amount of memory used for a loaded model.
- Fixed bug which prevented the runtime from reporting more than 32 loaded models to metrics.
- Fixed replica group signature calculation check.

## Neuron Runtime Library [2.19.5.0]

Date: 12/21/2023

## New in this release

- **Added Out-of-bound error detection logic for Gather/Scatter operations**
  - Out-of-bound error message “failed to run scatter/gather (indirect memory copy), due to out-of-bound access” will be displayed on an OOB error
  - **The runtime execution will return an “Out of Bound” error return code in the case an OOB error occurs**
    - \* `NRT_EXEC_OOB = 1006`
- **Improved Neff not supported error message to list out runtime supported features vs features used by the Neff**
  - Example output: “NEFF version 2.0 uses unsupported features: [0x100000]. Neuron Runtime NEFF supported features map: [0x1ff]. Please update the aws-neuronx-runtime-lib package”
- **Increased limit of multicore custom ops functions**
  - Total number of CustomOps in a model has been increased to 10.
  - Note: these 10 ops have to reside in one .so, as a result, they either have to be all single core op or all multicore op.

## Neuron Runtime Library [2.18.15.0]

Date: 11/09/2023

## Bug fixes

- Removed unnecessary data collection during execution logging which could impact performance.

**Neuron Runtime Library [2.18.14.0]**

Date: 10/26/2023

**New in this release**

- Add beta Collectives barrier API (`nrt_barrier`) to `nrt_experimental.h`
- Improved error handling and logging for NaNs produced by intermediate calculations that do not affect output.
- Improved logging by surfacing model id on load and execution errors.
- Output a better error message when Neff fails to load due to JSON size issues, e.g. “File sg00/def.json size (8589934592) exceeds json parser maximum (4294967295)”

**Bug fixes**

- Fixed logging error message to specify Neuron Cores instead of Neuron Devices when loading unsupported collectives topology.
- Fixed segfault on error path when Neuron Device fails to initialize.

**Neuron Runtime Library [2.17.7.0]**

Date: 9/14/2023

**New in this release**

- Improved logging by printing out NEFF name in debug logs of `nrt_execute`

**Bug fixes**

- Fixed hang that would occur when running a NEFF which contains embedding update instructions in multiple functions.
- Fixed issue where the Neuron Runtime registered the same memory multiple times to an EFA device causing applications to exceed the number of physical pages that could be registered.
- Fixed `assert (void tvm::runtime::GraphRuntime::PatchDltDataPtr(DLTensor*, uint32_t*, size_t): Assertion `tensor_get_mem_type(grt->io_tensor) == NRT_TENSOR_MEM_TYPE_MALLOC' failed.)` that occurred on INF1 caused by an uninitialized pointer.
- Fixed potential hang that can occur when partial replica groups for collectives are present in a NEFF.

## Neuron Runtime Library [2.16.14.0]

Date: 9/01/2023

### Bug fixes

- Fixed a segfault on failure to complete Neuron Device initialization. New behavior will avoid the failure and escalate a fixed Neuron Runtime error code (NERR\_FAIL, code 0x1)
- Improved error messages around Neuron Device initialization failures.

## Neuron Runtime Library [2.16.8.0]

Date: 8/09/2023

### New in this release

- Add runtime version and capture time to NTFF
- Improved Neuron Device copy times for all instance types via async DMA copies
- Improved error messages for unsupported topologies (example below)  
global comm ([COMM ID]) has less channels than this replica group ([REPLICA GROUP ID]) :  
likely not enough EFA devices found if running on multiple nodes or CC not permitted on this group  
[[TOPOLOGY]]
- Improved logging message for collectives timeouts by adding rank id to trace logs (example below)  
[gid: [RANK ID]] exchange proxy tokens
- Improved error messages when loading NEFFs with unsupported instructions (example below)  
Unsupported hardware operator code [OPCODE] found in neff.  
Please make sure to upgrade to latest aws-neuronx-runtime-lib and aws-neuronx-collective; for detailed installation instructions visit Neuron documentation.

### Bug fixes

- Fixed “failed to get neighbor input/output addr” error when loading collectives NEFF compiled with callgraph flow and NEFF without callgraph flow.

## Neuron Runtime Library [2.15.14.0]

Date: 8/09/2023

**New in this release**

- Reduced the contiguous memory size requirement for initializing Neuron Runtime on trn1/inf2 instance families by shrinking some of the notification buffers. A particularly large decrease was the reduction of a 4MB error notification buffer down to 64K. Expectation is that under memory constrained or highly fragmented memory systems, the Neuron Runtime would come up more reliably than previous versions.

**Neuron Runtime Library [2.15.11.0]**

Date: 7/19/2023

**New in this release**

- Added beta asynchronous execution feature which can reduce latency by roughly 12% for training workloads. See Runtime Configuration guide for details on how to use the feature.
- AllReduce with All-to-all communication pattern enabled for 16 ranks on TRN1/TRN1N within the instance (intranode); choice of 16 ranks is limited to NeuronCores 0-15 or 16-31.
- Minor improvement in end-to-end execution latency after reducing the processing time required for benign error notifications.
- Reduced notification overhead by using descriptor packing improving DMA performance for memory bound workloads by up to 25%.
- Improved load speed by removing extraneous checks that were previously being performed during loads.
- Minor performance boost to CC Ops by removing the need to sort execution end notifications.
- Bumped profiling NTFF version to version 2 to remove duplicate information which may result in hitting protobuf limits, and avoid crashing when using an older version of Neuron tools to postprocess the profile. Please upgrade to Neuron tools 2.12 or above to view profiles captured using this version of the Neuron runtime.

**Neuron Runtime Library [2.14.8.0]**

Date: 6/14/2023

**New in this release**

- Added All-to-All All-Reduce support for Neuron Collective operations, which is expected to improve All-Reduce performance by 3-7x in most cases.
- Added more descriptive `NEURON_SCRATCHPAD_PAGE_SIZE` to eventually replace `NEURON_RT_ONE_TMPBUF_PAGE_SIZE_MB`
- Neuron Runtime is now getting the device BDF from Neuron Driver for internal use.

### Bug fixes

- Fixed rare race condition caused by DMA memory barrier not being set for certain data transfers leading to non-determinism in outputs
- Fixed NeuronCore latency not being counted properly in Neuron metrics
- Removed stack allocation of error notifications buffer when parsing error notifications, which may lead to stack overflows on smaller stack sizes.

### Neuron Runtime Library [2.13.6.0]

Date: 05/01/2023

### New in this release

- Added support for internal Neuron Compiler change, Queue Set Instances, which leads to reduced NEFF footprints on Neuron Devices. In some cases, the reduction is as much as 60% smaller DMA ring size.

### Bug fixes

- Fixed a rare fabric deadlock scenario (hang) in NeuronCore v2 related to notification events.
- Ensure tensor store writes are complete before synchronization event is set.

### Neuron Runtime Library [2.12.23.0]

Date: 04/19/2023

### Bug fixes

- Minor internal bug fixes.

### Neuron Runtime Library [2.12.14.0]

Date: 03/28/2023

### New in this release

- Added support for 16 channels and 16 EFA devices, which is required for enabling EC2 TRN1N instances with Neuron.
- Added support for hierarchical All-Reduce and Reduce-Scatter. These implementations are now used by default and provides up to 75% reduction in latency for 2MB buffers across 256 ranks.
- Added support for loading more than one Neuron Custom Operator library.
- Added support for loading multicore Neuron Custom Operators.
- Updated INF2 to support rank 1 topology.

- Minor improvement in model load time for small models (below 100MB).

### Neuron Runtime Library [2.11.43.0]

Date: 02/08/2023

#### New in this release

- Added support for Neuron Custom C++ operators as a beta feature. As of this release, usage of Custom C++ operators requires a reset of the Neuron Runtime after running a model which invoked a Neuron Custom C++ operator.
- Added support for a counter that enable measuring FLOPS on neuron-top and neuron-monitor.
- Added support for LRU cache for DMA rings.

#### Bug fixes

- Fixed load failures due to memory bounds checking for Neuron Collective Compute operations in Runtime during model load.
- Fixed an internal bug that was preventing Neuron Runtime metrics from posting.
- Fixed a bug that caused segfaults as a result of double frees and stack overflows.

### Neuron Runtime Library [2.10.18.0]

Date: 11/07/2022

#### New in this release

- Minor bug fixes and enhancements.

### Neuron Runtime Library [2.10.15.0]

Date: 10/26/2022

#### New in this release

- Changed default runtime behavior to reset NeuronCores when initializing applications. With this change, the resetting of the Neuron Driver after application crash is no longer necessary. The new reset functionality is controlled by setting environment variable: `NEURON_RT_RESET_CORES`, see [NeuronX Runtime Configuration](#) for more information.

## Bug fixes

- Fixed a bug where Stochastic Rounding was not being set for collective communication operators
- Fixed an issue with triggering DMA for large tensors
- Increased default execution timeout to 30 seconds
- Fixed IOQ resetting queue to incorrect ring id value
- Updated the Neuron driver for more reliable behavior of driver device reset. Driver no longer busy waits on reset or gets stuck waiting on reset, which caused kernel taints or caused driver unload attempts to fail.
- Fixed a bug that prevented collective communication over tensors larger than 2GB
- Fixed a bug that caused intermittent memory corruption when unloading a model
- Fixed a bug that caused the exhausting of EFA memory registration pool after multiple model reloads.

## Neuron Runtime Library [2.9.64.0]

Date: 10/10/2022

This release specifically adds support for training workloads on one or more EC2 TRN1 instances.

Required Neuron Driver Version: 2.5 or newer

## New in this release

- Broke out runtime into a separate package called aws-neuronx-runtime-lib.
- Added RUNPATH for discovery of libnrt.so, can be overridden with LD\_LIBRARY\_PATH.
- Added support for multiple collective compute operations, e.g. All-Reduce, Reduce-Scatter, All-Gather.
- Added Send/Recv operation support
- Added support for using multiple DMA engines with single pseudo embedding update instruction.
- Changed instruction buffer alignment to 32K.
- Reduced memory required during NEFF swapping.
- Enabled notifications for send/recv collectives operations.
- Added trace apis in support of execution profiling.
- Added support for TPB reset (default: off).
- Added version checking for libnccom (aws-neuronx-collectives).
- Added new runtime version API.
- Added 8-channel support for Trn1.
- Improved debug outputs.
- Added support for write combining on BAR4.
- Increased default execution timeout from 2 seconds to 30 seconds.
- Improved handling of zero-sized tensors

**Neuron Runtime 2.x (libnrt.so) release [2.2.51.0]**

Date: 03/25/2022

- Fixed an invalid memory access that could occur when unloading models.
- Reduced severity of logging for numerical errors from ERROR to WARN.
- Improved handling of models with numerous CPU operations to avoid inference failure due to memory exhaustion.

**Neuron Runtime 2.x (libnrt.so) release [2.2.31.0]**

Date: 01/20/2022

**New in the release**

- Changed error notifications from WARN to ERROR in cases when the causing problem is non-recoverable.
- Changed handling of inference timeouts (NERR\_TIMEOUT) to avoid failure when the timeout is related to a software thread scheduling conflict.

**Bug fixes**

- Increased the number of data queues in Neuron Runtime 2.x to match what was previously used in Neuron Runtime 1.x. The use of fewer number of data queues in Neuron Runtime 2.x was leading to crashes in a limited number of models.
- Fixed the way Neuron Runtime 2.x updates the inference end timestamp. Previously, Neuron Runtime 2.x update of the inference end timestamp would have lead to a negative latency statistics in neuron-monitor with certain models.

**Neuron Runtime 2.x (libnrt.so) release [2.2.18.0]**

Date: 11/05/2021

- Resolved an issue that affect the use of Neuron within container. In previous Neuron Runtime release (libnrt.so.2.2.15.0), when /dev/neuron0 was not used by the application, Neuron Runtime attempted and failed to initialize /dev/neuron0 because user didn't pass /dev/neuron0 to the container. this Neuron Runtime release (libnrt.so.2.2.18.0) allows customers to launch containers with specific NeuronDevices other than /dev/neuron0.

**Neuron Runtime 2.x (libnrt.so) release [2.2.15.0]**

Date: 10/27/2021



## New in this release

- *First release of Neuron Runtime 2.x* - In this release we are introducing Neuron Runtime 2.x which is a shared library named (`libnrt.so`) and replacing Neuron Runtime 1.x server (`neruon-rtd`) . Upgrading to `libnrt.so` improves throughput and latency, simplifies Neuron installation and upgrade process, introduces new capabilities for allocating NeuronCores to applications, streamlines container creation, and deprecates tools that are no longer needed. The new library-based runtime (`libnrt.so`) is integrated into Neuron's ML Frameworks (with the exception of MXNet 1.5) and Neuron Tools packages directly - users no longer need to install/deploy the `aws-neuron-runtimepackage`.

---

### Important:

- You must update to the latest Neuron Driver (`aws-neuron-dkms` version 2.1 or newer) for proper functionality of the new runtime library.
  - Read *Introducing Neuron Runtime 2.x (libnrt.so)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.
  - Read *Migrate your application to Neuron Runtime 2.x (libnrt.so)* for detailed information of how to migrate your application.
- 

This document is relevant for: Inf1, Inf2, Trn1, Trn2

This document is relevant for: Inf1, Inf2, Trn1, Trn2

## Neuron Driver Release Notes

### Table of contents

- *Known issues*
- *Neuron Driver release [2.22.2.0]*
- *Neuron Driver release [2.21.37.0]*
- *Neuron Driver release [2.20.74.0]*
- *Neuron Driver release [2.20.28.0]*
- *Neuron Driver release [2.19.64.0]*
- *Neuron Driver release [2.18.20.0]*
- *Neuron Driver release [2.18.12.0]*
- *Neuron Driver release [2.17.17.0]*
- *Neuron Driver release [2.16.7.0]*
- *Neuron Driver release [2.15.9.0]*
- *Neuron Driver release [2.14.5.0]*
- *Neuron Driver release [2.13.4.0]*
- *Neuron Driver release [2.12.18.0]*
- *Neuron Driver release [2.12.11.0]*
- *Neuron Driver release [2.11.9.0]*

- *Neuron Driver release [2.10.11.0]*
- *Neuron Driver release [2.9.4.0]*
- *Neuron Driver release [2.8.4.0]*
- *Neuron Driver release [2.7.33.0]*
- *Neuron Driver release [2.7.15.0]*
- *Neuron Driver release [2.6.26.0]*
- *Neuron Driver release [2.5.38.0]*
- *Neuron Driver release [2.3.26.0]*
- *Neuron Driver release [2.3.11.0]*
- *Neuron Driver release [2.3.3.0]*
- *Neuron Driver release [2.2.14.0]*
- *Neuron Driver release [2.2.13.0]*
- *Neuron Driver release [2.2.6.0]*
- *Neuron Driver release [2.1]*

## **Known issues**

Updated : 04/29/2022

- In rare cases of multi-process applications running under heavy stress a model load failure may occur. This may require reloading of the Neuron Driver as a workaround.

## **Neuron Driver release [2.22.2.0]**

Date: 06/24/2025

## **Bug Fixes**

- Added workaround for HW DGE descriptor fetching bug
- Fixed typos in certain error log messages

## **Upcoming Neuron driver 2.21 support changes for Inf1 instance users**

- Starting with Neuron Release 2.26, Neuron driver versions above 2.21 will only support non-Inf1 instances (such as Trn1, Inf2, or other instance types).
- Inf1 instance users, Neuron driver 2.21 and below will remain supported with regular security patches.
- Inf1 instance users are advised to pin the Neuron driver version to 2.21.\* in their installation script using following command:

```
sudo apt-get install aws-neuronx-dkms=2.21.* -y
```

**Neuron Driver release [2.21.37.0]**

Date: 05/19/2025

**New in this release**

- Added the ability for users to read power utilization for each neuron device via a sysfs interface. This interface shows the minimum, maximum and average power consumed by the device over the past minute, expressed as a percentage of the device's maximum power. ([reference](#))
- Added the ability for users to read the device utilization. This shows up as the microseconds between the start and end of the current execution on hardware. ([reference](#))

**Neuron Driver release [2.20.74.0]**

Date: 05/12/2025

**New in this release**

- Fixes DMA abort errors on Trainium2 that could occur in Neuron Runtime during specific workloads.

**Neuron Driver release [2.20.28.0]**

Date: 04/03/2025

**New in this release**

- This driver is required to run with Neuron Runtime 2.24 or later on Trainium2 machines. Included in the release is a bug fix to avoid device memory corruption issues leading to undefined Neuron Device behavior.

**Improvements**

- Improved interface between libnrt and the Driver resulting in stability improvements.

**Neuron Driver release [2.19.64.0]**

Date: 12/20/2024

### New in this release

- Added Trainium2 support

### Improvements

- Optimized HBM Memory allocation to reduce fragmentation. See [here](#) for more details.

### Neuron Driver release [2.18.20.0]

Date: 11/20/2024

### Bug Fixes

- This release addresses an issue with Neuron Driver that can lead to a user-space application either gaining access to kernel addresses or providing the driver with spoofed memory handles (kernel addresses) that can be potentially used to gain elevated privileges. We would like to thank [Cossack9989](#) for reporting and collaborating on this issue.

### Neuron Driver release [2.18.12.0]

Date: 09/16/2024

### New in this release

- Introduced a sysfs memory usage counter for DMA rings ([reference](#))

### Bug Fixes

- Resolved an issue where a memory allocation failure caused a hang due to the memory allocation lock not being released
- Resolved an issue where the driver was allocating more memory than needed for aligned device allocations

### Neuron Driver release [2.17.17.0]

Date: 07/03/2024

**New in this release**

- Improved detection and reporting of DMA errors
- Added more fine grained sysfs metrics to track memory allocation types
- Logging improvements

**Bug Fixes**

- Fixed compatibility issues for the Linux 6.3 kernel
- Resolved issue where device reset handling code was not properly checking the failure metric

**Neuron Driver release [2.16.7.0]**

Date: 04/01/2024

**Bug Fixes**

- Fixed installation issues caused by API changes in Linux 6.3 and 6.4 kernel distributions.
- Fixed an installation build failure when fault-injection is enabled in the kernel.
- Fixed an issue where sysfs total peak memory usage metrics can underflow
- Removed usage of sysfs\_emit which is not supported on Linux kernels <= v5.10-rc1

**Neuron Driver release [2.15.9.0]**

Date: 12/21/2023

**Bug Fixes**

- Release PCIe BAR4 on driver startup failure
- Fix container BDF indexing issues to support relative device ordering used by containers
- Remove incorrect error message in neuron\_p2p\_unregister\_va and harden P2P error checking

**Neuron Driver release [2.14.5.0]**

Date: 10/26/2023

### New in this release

- Show uncorrectable SRAM and HBM ECC errors on TRN1 and INF2
- Fixed double free on error path during driver startup

### Neuron Driver release [2.13.4.0]

Date: 9/14/2023

### New in this release

- Added sysfs support for showing connected devices on trn1.32xl, inf2.24xl, and inf2.48xl instances.

### Neuron Driver release [2.12.18.0]

Date: 9/01/2023

### Bug Fixes

- Added fixes required by Neuron K8 components for improving reliability of pod failures (see Neuron K8 release notes for more details).
- Added fixes required by Neuron K8 components to support zero-based indexing of Neuron Devices in Kubernetes deployments.

### Neuron Driver release [2.12.11.0]

Date: 8/28/2023

### New in this release

- Added FLOP count to sysfs (flop\_count)
- Added connected Neuron Device ids to sysfs (connected\_devices)
- Added async DMA copy support
- Suppressed benign timeout/retry messages

## Bug Fixes

- Allocated CC-Core to correct NeuronCore; splitting CC-Cores evenly between NeuronCores.

## Neuron Driver release [2.11.9.0]

Date: 7/19/2023

## New in this release

- Added support for creating batch DMA queues.

## Bug Fixes

- Error message, “ncdev is not NULL”, was being printed unnecessarily. Fixed.
- Fix DMA timeouts during NeuronCore reset of neighboring core caused by incorrect nc\_id (NeuronCore ID) assigned to reserved memory

## Neuron Driver release [2.10.11.0]

Date: 6/14/2023

## New in this release

- Added memory usage breakdown by category to the Neuron Sysfs nodes. New categories are code, misc, tensors, constants, and scratchpad. Please see the Sysfs page under Neuron Tools for more detailed description of each.
- Improved NeuronCore initialization (nrt\_init) performance by approximately 1 second.

## Bug Fixes

- Fixed small timing window during NeuronCore resets, which previously would timeout during memcpy
- Removed potential double free of memory when terminating the Neuron Driver.
- Fixed sysfs race condition, which was leading to Neuron Driver crash during termination.

## Neuron Driver release [2.9.4.0]

Date: 05/01/2023

### New in this release

- Added dma\_buf support, which is needed for future EFA implementations in the Linux kernel.
- Added new IOCTL to get Neuron Device BDF (used by Neuron Runtime)
- Added optional support for sysfs notify (off by default). See Neuron Sysfs documentation (under Neuron System Tools) for more details.

### Bug Fixes

- Fixed max DMA queue size constant to be the correct size - previous incorrect sizing had potential to lead to DMA aborts (execution timeout).

### Neuron Driver release [2.8.4.0]

Date: 03/28/2023

### New in this release

- Supports both Trn1n and Inf2 instance types.
- Renamed `NEURON_ARCH_INFERENTIA=>NEURON_ARCH_V1` and `NEURON_ARCH_TRN=>NEURON_ARCH_V2`
- Under sysfs nodes, the following changes were made:
  - Changed “infer” metrics to “execute” metrics
  - Added peak memory usage metric
  - Removed empty dynamic metrics directory
  - Removed refresh rate metric
  - Fixed arch type names in sysfs

### Bug Fixes

- Fixed minor memory leak when closing the Neuron Runtime.
- Fixed memory leaks on error paths in Neuron Driver.
- Added a workaround to resolve hangs when NeuronCore reset is ran while another core is performing DMA operations.



**Neuron Driver release [2.7.33.0]**

Date: 02/24/2023

**Bug Fixes**

- Added a retry mechanism to mitigate possible data copy failures during reset of a NeuronCore. An info log message will be emitted when this occurs indicating that the retry was attempted. An example:

```
kernel: [726415.485022] neuron:ndma_memcpy_wait_for_completion: DMA completion_
↳ timeout for UDMA_ENG_33 q0
kernel: [726415.491744] neuron:ndma_memcpy_offset_move: Failed to copy memory_
↳ during a NeuronCore reset: nd 0, src 0x100154480000, dst 0x100154500000, size_
↳ 523264. Retrying the copy.
```

**Neuron Driver release [2.7.15.0]**

Date: 02/08/2023

**New in this release**

- Added Neuron sysfs metrics under `/sys/devices/virtual/neuron_device/neuron{0,1, ...}/metrics/`

**Neuron Driver release [2.6.26.0]**

Date: 11/07/2022

**New in this release**

- Minor bug fixes and improvements.

**Neuron Driver release [2.5.38.0]**

Neuron Driver now supports INF1 and TRN1 EC2 instance types. Name of the driver package changed from `aws-neuron-dkms` to `aws-neuronx-dkms`. Please remove the older driver package before installing the newest one.

Date: 10/10/2022

### New in this release

- Support added for EC2 Trn1 instance types and ML training workloads.
- Added missing GPL2 LICENSE file.
- Changed package name to aws-neuronx-dkms (was previously minus the 'x').
- Security Update – blocked user space access to control registers and DMA control queues intended to be used by the Neuron Driver only.
- Added support for DMA Aborts to avoid hangs.
- Added support for TPB Reset.
- Added sysfs entries for triggering resets and reading core counts.
- Added write combining on BAR4.
- Added PCI Device ID update as part of install.
- Added handling for known duplicate device id error.

### Bug Fixes

- Fixed a null pointer free scenario.
- Fixed installation issue related to install without internet connectivity.

### Neuron Driver release [2.3.26.0]

Date: 08/02/2022

### Bug Fixes

- Security Update: Blocked user space access to control registers and DMA control queues intended to be used by the Neuron Driver only. Recommending upgrade to all customers.

### Neuron Driver release [2.3.11.0]

Date: 05/27/2022

### New in this release

- This driver is required to support future releases of the Neuron Runtime. Included in the release is both a bug fix to avoid a kernel crash scenario and an increased compatibility range to ensure compatibility with future versions of Neuron Runtime.

**Bug Fixes**

- Correction to huge aligned memory allocation/freeing logic that was previously susceptible to crashes in the kernel. The crash would bring down the OS. Recommending upgrade to all customers.

**Neuron Driver release [2.3.3.0]**

Date: 04/29/2022

**New in this release**

- Minor performance improvements on inference and loading of models.

**Bug Fixes**

- Reduced Host CPU usage when reading `hw_counters` metric from neuron-monitor
- Minor bug fixes.

**Neuron Driver release [2.2.14.0]**

Date: 03/25/2022

**New in this release**

- Minor updates

**Neuron Driver release [2.2.13.0]**

Date: 01/20/2022

**New in this release**

- Minor updates

**Neuron Driver release [2.2.6.0]**

Date: 10/27/2021

### New in this release

- Memory improvements made to ensure all allocations are made with 4K alignments.

### Resolved issues

- No longer delays 1s per NeuronDevice when closing Neuron Tools applications.
- Fixes a Ubuntu 20 build issue

### Neuron Driver release [2.1]

- Support is added for Neuron Runtime 2.x (`libnrt.so`).
- Support for previous releases of Neuron Runtime 1.x is continued with Driver 2.x releases.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### Neuron Collectives Release Notes

Neuron Collectives refers to a set of libraries used to support collective compute operations within the Neuron SDK. The collectives support is delivered via the `aws-neuronx-collectives` package and includes a pre-built version of the OFI plugin required for use of collectives with Elastic Fabric Adapter (EFA).

#### Table of contents

- *Neuron Collectives [2.26.43.0]*
- *Neuron Collectives [2.25.65.0]*
- *Neuron Collectives [2.24.59.0]*
- *Neuron Collectives [2.23.135.0]*
- *Neuron Collectives [2.23.133.0]*
- *Neuron Collectives [2.22.26.0]*
- *Neuron Collectives [2.21.46.0]*
- *Neuron Collectives [2.20.22.0]*
- *Neuron Collectives [2.20.11.0]*
- *Neuron Collectives [2.19.7.0]*
- *Neuron Collectives [2.18.18.0]*
- *Neuron Collectives [2.17.9.0]*
- *Neuron Collectives [2.16.16.0]*
- *Neuron Collectives [2.16.8.0]*
- *Neuron Collectives [2.15.16.0]*
- *Neuron Collectives [2.15.13.0]*

- *Neuron Collectives [2.14.9.0]*
- *Neuron Collectives [2.13.7.0]*
- *Neuron Collectives [2.12.35.0]*
- *Neuron Collectives [2.12.22.0]*
- *Neuron Collectives [2.11.47.0]*
- *Neuron Collectives [2.10.20.0]*
- *Neuron Collectives [2.9.86.0]*

### Neuron Collectives [2.26.43.0]

Date: 06/24/2025

#### Bug fixes

- Fixed various memory leaks which occur during process cleanup

### Neuron Collectives [2.25.65.0]

Date: 05/19/2025

#### New in this release

- Added multinode collectives support for Trainium2 instances without EFA devices

#### Improvements

- Minor performance improvement to network proxy handshake

#### Bug fixes

- Fixed memory leak clearing up communication devices during `nrt_close`

### Neuron Collectives [2.24.59.0]

Date: 04/03/2025

### Improvements

- Improved interface between `libnccom` and `libnrt` resulting stability improvements

### Neuron Collectives [2.23.135.0]

Date: 01/14/2025

### Improvements

- Aws-ofi-nccl: minor performance improvement

### Neuron Collectives [2.23.133.0]

Date: 12/20/2024

### New in this release

- Added Trainium2 support

### Improvements

- Improved startup times for large scale training jobs by up to 5 seconds
- Enhanced error logging for bootstrap failures

### Neuron Collectives [2.22.26.0]

Date: 09/16/2024

### New in this release:

- Added check to print out an error message on invalid `NEURON_RT_ROOT_COMM_ID` configurations

### Bug fixes

- Resolved an issue where the `libnccom.so` filename was versioned incorrectly as `libnccom.so.2.y.y`. Will be correctly versioned as `libnccom.so.2.22.26` in this release.

**Neuron Collectives [2.21.46.0]**

Date: 07/03/2024

**New in this release:**

- Bootstrap changes to improve application startup latency for large-scale workloads
- Logging improvements

**Neuron Collectives [2.20.22.0]**

Date: 04/01/2024

**New in this release:**

- minor bug fixes and enhancements

**Neuron Collectives [2.20.11.0]**

Date: 02/13/2024

**Bug Fixes**

- Require “libatomic” for rpm installs

**Neuron Collectives [2.19.7.0]**

Date: 12/21/2023

**New in this release**

- Improve collectives barrier latency from 500us to 40us

**Bug Fixes**

- Fix bug where proxy thread blocks the runtime from adding ops leading to an execution hang

### Neuron Collectives [2.18.18.0]

Date: 10/26/2023

New in this release: \* Bumped compatibility version to 17 to align with struct change in the nec.h header

### Neuron Collectives [2.17.9.0]

Date: 9/14/2023

New in this release: \* minor bug fixes and enhancements

### Neuron Collectives [2.16.16.0]

Date: 9/01/2023

New in this release: \* minor bug fixes and enhancements

### Neuron Collectives [2.16.8.0]

Date: 8/28/2023

New in this release:

- Improved error messages for unsupported topologies
- Improved timeout error messages for bootstrapInit

Bug Fixes: \* Fix bug where Linux kernel version check for SAFE\_FORK env variable was incorrectly requiring SAFE\_FORK to be set on kernel versions greater than 5

### Neuron Collectives [2.15.16.0]

Date: 8/09/2023

New in this release:

- minor bug fixes and enhancements

### Neuron Collectives [2.15.13.0]

Date: 7/19/2023

New in this release:

- AllReduce with All-to-all communication pattern enabled for 16 ranks on TRN1/TRN1N within the instance (intranode); choice of 16 ranks is limited to NeuronCores 0-15 or 16-31.

Bug Fixes:

- Fix incorrect mask calculation for 16 ranks when using NeuronCores 16-31
- Fix channels for 16 ranks to avoid failures in the runtime; restrict participating ranks to 0-15 or 16-31



**Neuron Collectives [2.14.9.0]**

Date: 6/14/2023

New in this release

- Added check for FI\_EFA\_FORK\_SAFE environment variable; now forcing the flag to be set to 1 for multinode runs executing on Linux kernels older than 5.15.

**Neuron Collectives [2.13.7.0]**

Date: 05/01/2023

New in this release

- Added support for dma\_buf - required for future EFA and Linux kernel updates.
- Reduced benign reporting of timeouts. Previous implementations reported “Timeout waiting for incoming connection” too frequently (log spam).

**Neuron Collectives [2.12.35.0]**

Date: 04/19/2023

Bug Fixes

- Fixed support for SOCKET\_IFNAME config that was affecting EKS users at scale on large training jobs.

**Neuron Collectives [2.12.22.0]**

Date: 03/28/2023

New in this release

- Added support for TRN1N.
- Added support for 16 channels and 16 EFA devices, which is required for enabling EC2 TRN1N instances with Neuron.
- Added support for hierarchical All-Reduce and Reduce-Scatter. These implementations are now used by default and provides up to 75% reduction in latency for 2MB buffers across 256 ranks.

**Neuron Collectives [2.11.47.0]**

Date: 02/08/2023

New in this release

- Added support for Inf2.

### Neuron Collectives [2.10.20.0]

Date: 10/10/2022

New in this release

- Improved logging to appear similar in style to Neuron Runtime

Bug Fixes

- Fixed memory registration to support 2GB+ sizes
- Fixed association of network devices to channels (removes previous hard-coding).

### Neuron Collectives [2.9.86.0]

Date: 10/10/2022

New in this release

- Added support for All-Reduce, Reduce-Scatter, All-Gather, and Send/Recv operations.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### API Reference Guide

- *Runtime API*

### Configuration Guide

- *Runtime Configuration*

### Misc

- *Troubleshooting on Inf1 and Trn1*
- *FAQ*
- *Neuron Runtime Release Notes*
- *Neuron Driver Release Notes*
- *Neuron Collectives Release Notes*

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 6.2 Monitoring Tools

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.2.1 Neuron Monitor User Guide

#### Table of contents

- *Overview*
- *Using neuron-monitor*
  - *Configuration file example*
  - *Neuron applications tagging*
  - *JSON objects and fields in the configuration file*
  - *Neuron Runtime-level metric groups*
  - *System-wide metric groups*
- *Execution model*
- *The JSON output format*
  - *instance\_info*
  - *neuron\_hardware\_info*
  - *neuron\_k8s\_info*
  - *Metric Groups*
- *Neuron application level metric groups*
  - *neuroncore\_counters*
  - *execution\_stats*
  - *memory\_used*
  - *neuron\_runtime\_vcpu\_usage*
- *System level metric groups*
  - *neuron\_hw\_counters*
  - *vcpu\_usage*
  - *memory\_info*
- *Companion scripts*
  - *neuron-monitor-cloudwatch.py*
  - *neuron-monitor-prometheus.py*
  - *neuron-monitor-k8s-info.py (Beta)*
- *Running neuron monitor in Kubernetes environment*

## Overview

**neuron-monitor** collects metrics and stats from the Neuron Applications running on the system and streams the collected data to `stdout` in JSON format. It is provided as part of the `aws-neuron-tools` package.

These metrics and stats are organized into **metric groups** which can be configured by providing a configuration file as described in [Using neuron-monitor](#)

When running, **neuron-monitor** will:

- Collect the data for the metric groups which, based on the elapsed time since their last update, need to be updated
- Take the newly collected data and consolidate it into a large report
- Serialize that report to JSON and stream it to `stdout` from where it can be consumed by other tools - such as the sample `neuron-monitor-cloudwatch.py` and `neuron-monitor-prometheus.py` scripts.
- Wait until at least one **metric group** needs to be collected and repeat this flow

---

**Note:** `neuron-monitor` fully supports the newly launched Trn2 instances.

---

## Using neuron-monitor

### neuron-monitor CLI

**neuron-monitor** [parameters]

`neuron-monitor` accepts the following optional parameters:

- `--verbose` (int) default=0: Can be 0 to 4, and controls the amount of debugging and verbose information sent to `stderr`; **0: no output, 4: maximum verbosity**
- `-c, --config-file` (string): Allows specifying a valid path to a neuron-monitor JSON configuration file

#### Example:

```
neuron-monitor -c monitor.conf
```

Not specifying any configuration file will enable collecting all the metric groups with a period of 5 seconds for all currently running Neuron applications.

## Configuration file example

Example of a configuration file which enables all available **metric groups** for every running Neuron application, with a global update period of 1 second and sets an update period of 2 seconds for the `"neuron_hw_counters"` metric group:

```
{
  "period": "1s",
  "neuron_runtimes": [
    {
      "tag_filter": ".*",
      "metrics": [
        {
          "type": "neuroncore_counters"
```

(continues on next page)

(continued from previous page)

```

    },
    {
      "type": "memory_used"
    },
    {
      "type": "neuron_runtime_vcpu_usage"
    },
    {
      "type": "execution_stats"
    }
  ]
},
"system_metrics": [
  {
    "type": "vcpu_usage"
  },
  {
    "type": "memory_info"
  },
  {
    "period": "2s",
    "type": "neuron_hw_counters"
  }
]
}

```

## Neuron applications tagging

In order to make application monitoring easier, Neuron applications can be tagged with a 255 character string which identifies that app. Tagging is done using the `NEURON_PROCESS_TAG` environment variable.

For example: `NEURON_PROCESS_TAG=my_app_1 python training.py` will associate the `my_app_1` tag with that Python application. If `NEURON_PROCESS_TAG` is not specified, the application's PID will be used as a TAG.

This tag will be used by neuron-monitor to filter Neuron applications.

## JSON objects and fields in the configuration file

- "neuron\_runtimes" - array of objects specifying which Neuron Applications to monitor and what metric groups are enabled for each of them
  - "tag\_filter" - a regex which will be used to filter Neuron applications tags in order to determine if they will be monitored (optional)
  - "metrics" - array of objects specifying which metric groups to capture for this Neuron application
    - \* "type" - type of metric group
- "period" - this field applies to **metric group** objects and sets the amount of time between two updates for that metric group
  - if can be specified as part of the **root** and/or **neuron\_runtime** objects where it applies to all their children, and/or as part of a **metric group** object

- if there's no period specified, a default value of **5 seconds** will be used
- "system\_metrics" - array of objects specifying which system level metric groups are enabled

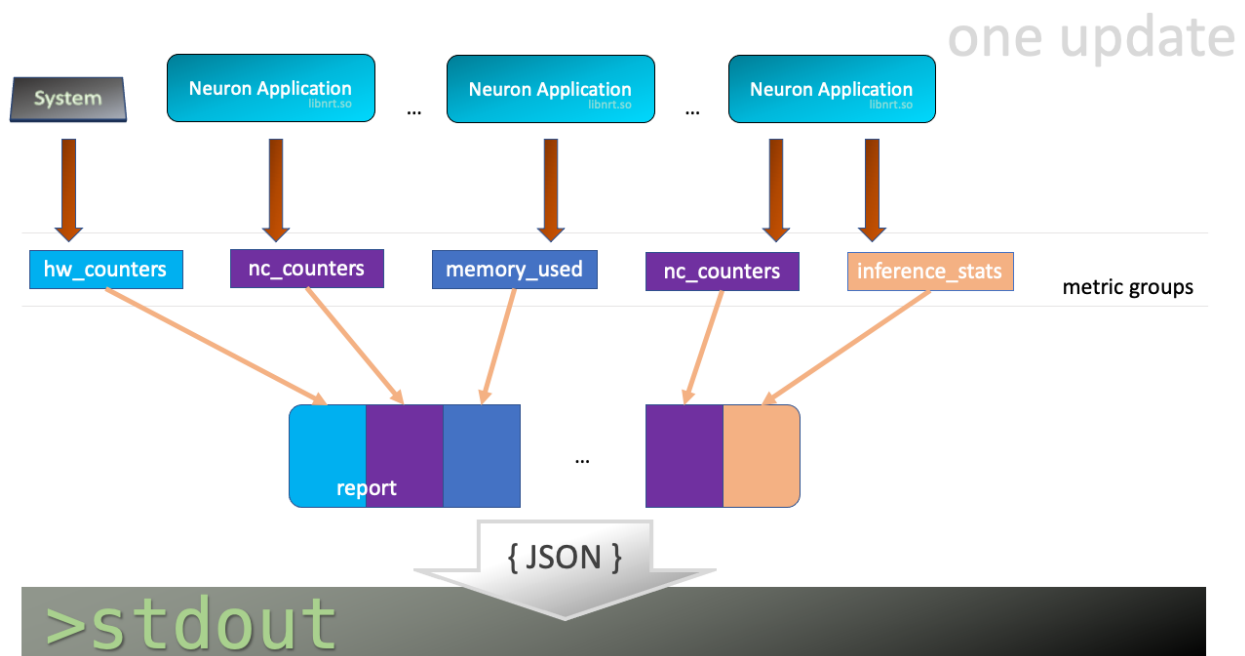
### Neuron Runtime-level metric groups

- *neuroncore\_counters* - NeuronCore related metrics
- *memory\_used* - data on the amount of memory used by the Neuron application
- *vcpu\_usage* - Neuron application vCPU utilization data
- *execution\_stats* - Neuron application execution stats, including error count and latency

### System-wide metric groups

- *vcpu\_usage* - system-wide vCPU usage
- *memory\_info* - system-wide memory usage
- *neuron\_hw\_counters* - counters for correctable and uncorrectable memory ecc events

### Execution model



neuron-monitor waits for one or more **metric groups** to be up for update, then collects the corresponding data, consolidates it into a report which is streamed to stdout as a JSON and goes back to waiting.

## The JSON output format

Whenever the report gets updated, a complete JSON is written to stdout. This is its structure:

```
{
  "neuron_runtime_data": [
    {
      "pid": 0,
      "address": "",
      "neuron_runtime_tag": "my_app_1",
      "error": "",
      "report": {
        "neuroncore_counters": {
          [...]
        },
        "execution_stats": {
          [...]
        },
        "memory_used": {
          [...]
        },
        "neuron_runtime_vcpu_usage": {
          [...]
        }
      }
    }
  ],
  "system_data": {
    "neuron_hw_counters": {
      [...]
    },
    "vcpu_usage": {
      [...]
    },
    "memory_info": {
      [...]
    }
  },
  "instance_info": {
    [...]
  },
  "neuron_hardware_info": {
    [...]
  },
  "neuron_k8s_info": {
    [...]
  }
}
```

- "neuron\_runtime\_data" is an array containing one entry per each Neuron application which passes the filter specified in the settings file
  - "pid" is the pid of this Neuron application
  - "neuron\_runtime\_tag" is the configured tag for the Neuron application

- "error" specifies any error that occurred when collecting data from this Neuron application
- "report" will contain the results for the Neuron application-level metric groups; their formats are described below
- "system\_data" has a similar structure to "neuron\_runtime\_data"'s "report" but only contains system-level metric groups (not associated to any Neuron application)

Regardless of the configuration, the following two JSON objects are always present in the output:

### instance\_info

Contains information about the instance on which neuron-monitor is running.

```
"instance_info": {
  "instance_name": "My_Instance",
  "instance_id": "i-0011223344556677a",
  "instance_type": "trn2n.48xlarge",
  "instance_availability_zone": "us-west-2b",
  "instance_availability_zone_id": "usw2-az2",
  "instance_region": "us-west-2",
  "ami_id": "ami-0011223344556677b",
  "subnet_id": "subnet-112233ee",
  "error": ""
}
```

Depending on when the instance was launched, the following fields might not be available:

- instance\_availability\_zone\_id : available only for instances launched in 2020-08-24 and later
- instance\_region : available only for instances launched on 2020-08-24 and later
- instance\_name : available only if instance\_region is set and aws-cli tools are installed

error will contain an error string if getting one of the fields, **except those mentioned above**, resulted in an error.

### neuron\_hardware\_info

Contains basic information about the Neuron hardware.

```
"neuron_hardware_info": {
  "neuron_device_type": "trainium2",
  "neuron_device_version": "v4",
  "neuroncore_version": "v3d",
  "neuron_device_count": 16,
  "neuron_device_memory_size": 103079215104,
  "neuroncore_per_device_count": 4,
  "logical_neuroncore_config": 2,
  "error": ""
}
```

- neuron\_device\_type: type of the Neuron Devices on the instance
- neuroncore\_version: version of the NeuronCores on the instance
- neuron\_device\_count : number of available Neuron Devices



- `neuron_device_memory_size`: total memory available on each Neuron Device
- `neuroncore_per_device_count` : number of NeuronCores present on each Neuron Device
- `logical_neuroncore_config` : the current Logical NeuronCore configuration
- `error` : will contain an error string if any occurred when getting this information (usually due to the Neuron Driver not being installed or not running).

The following JSON object is disabled by default, but can be made available if “k8s\_info” is enabled:

### neuron\_k8s\_info

Contains information about what Kubernetes pods/containers are using Neuron resources

```
"neuron_k8s_info": {
  "period": 15.030359284,
  "neuroncores_k8s_info": {
    "0": {
      "pod_name": "p0",
      "namespace": "n0",
      "container_name": ["c0"]
    },
    "1": {
      "pod_name": "p0",
      "namespace": "n0",
      "container_name": ["c0"]
    },
    ...
  },
  "neurondeVICES_k8s_info": {
    "0": {
      "pod_name": "p0",
      "namespace": "n0",
      "container_name": ["c0"]
    },
    ...
  },
  "error": ""
},
```

- `"neuroncores_k8s_info"` - object containing information on which Neuron cores are being used by Kubernetes pod/containers, indexed by Neuron core index: `"neuroncore_index": { neuroncore_k8s_data }`
  - `"pod_name"` - name of pod using Neuron core
  - `"namespace"` - namespace of pod using Neuron core
  - `"container_name"` - names of containers using Neuron core
- `"neurondeVICES_k8s_info"` - object containing information on which Neuron devices are being used by Kubernetes pod/containers, indexed by Neuron device index: `"neurondevice_index": { neurondevice_k8s_data }`
  - `"pod_name"` - name of pod using Neuron device
  - `"namespace"` - namespace of pod using Neuron device
  - `"container_name"` - names of containers using Neuron device

- "error" - will contain an error string if any occurred when getting this information

For more information on how to enable K8s information, see [neuron-monitor-k8s-info.py \(Beta\)](#).

## Metric Groups

Each **metric group** requested in the settings file will get an entry in the resulting output. The general format for such an entry is:

```
"metric_group": {
  "period": 1.015, // Actual captured period, in seconds
  "error": "",     // Error, if any occurred, otherwise an empty string
  [...]          // Metric group specific data
}
```

## Neuron application level metric groups

### neuroncore\_counters

```
"neuroncore_counters": {
  "period": 1.000113182,
  "neuroncores_in_use": {
    "0": {
      "neuroncore_utilization": 42.01,
      "flops": 1234567891011,
      "v3d": {
        "nc_v3.0": {
          "neuroncore_utilization": 21.01
        },
        "nc_v3.1": {
          "neuroncore_utilization": 63.01
        }
      }
    },
    "1": {
      "neuroncore_utilization": 42.02,
      "flops": 1234567891021,
      "v3d": {
        "nc_v3.2": {
          "neuroncore_utilization": 21.02
        },
        "nc_v3.3": {
          "neuroncore_utilization": 63.02
        }
      }
    },
    [...]
  },
  "error": ""
}
```

- "neuroncores\_in\_use" is an object containing data for all the NeuronCores that were active when the data was captured, indexed by NeuronCore index: "neuroncore\_index": { neuroncore\_data }
  - "neuroncore\_utilization" - NeuronCore utilization, in percent, during the captured period
  - "flops" - number of floating point operations per second during the captured period
  - "v3d" - only available on Trn2 - contains the utilization for every physical NeuronCore that makes up the current NeuronCore
- "error" - string containing any error that occurred when collecting the data

### execution\_stats

```

"execution_stats": {
  "period": 1.030613214,
  "error_summary": {
    "generic": 0,
    "numerical": 0,
    "transient": 0,
    "model": 0,
    "runtime": 0,
    "hardware": 0
  },
  "execution_summary": {
    "completed": 123,
    "completed_with_err": 0,
    "completed_with_num_err": 0,
    "timed_out": 0,
    "incorrect_input": 0,
    "failed_to_queue": 0
  },
  "latency_stats": {
    "total_latency": {
      "p0": 0.01100001,
      "p1": 0.01100002,
      "p25": 0.01100004,
      "p50": 0.01100008,
      "p75": 0.01100010,
      "p99": 0.01100012,
      "p100": 0.01100013
    },
    "device_latency": {
      "p0": 0.01000001,
      "p1": 0.01000002,
      "p25": 0.01000004,
      "p50": 0.01000008,
      "p75": 0.01000010,
      "p99": 0.01000012,
      "p100": 0.01000013
    }
  },
  "error": ""
},

```

- "error\_summary" is an object containing the error counts for the captured period indexed by their type
  - "generic" - generic execution errors
  - "numeric" - NAN errors encountered during execution
  - "transient" - recoverable errors, such as ECC corrections
  - "model" - model-related errors
  - "runtime" - Neuron Runtime errors
  - "hardware" - hardware errors such as uncorrectable ECC issues
- "execution\_summary" is an object containing all execution outcome counts for the captured period indexed by their type
  - "completed" - executions completed successfully
  - "completed\_with\_err" - executions that ended in an error other than a numeric error
  - "completed\_with\_num\_err" - executions that ended in a numeric error
  - "timed\_out" - executions that took longer than the Neuron Runtime configured timeout value
  - "incorrect\_input" - executions that failed to start due to incorrect input being provided
  - "failed\_to\_queue" - execution requests that were rejected due to Neuron Runtime not being able to queue them
- "latency\_stats" contains two objects containing latency percentiles, in seconds, for the data captured for the model executed during the captured period. If there are no models being executed during this time, the two objects will be null (i.e. "total\_latency": null)
  - "total\_latency" - percentiles, in seconds, representing latency for an execution as measured by the Neuron Runtime
  - "device\_latency" - percentiles, in seconds, representing execution time exclusively on the Neuron Device
- "error" - string containing any error that occurred when collecting the data

## memory\_used

```
"memory_used": {
  "period": 1.00001,
  "neuron_runtime_used_bytes": {
    "host": 6997643264,
    "neuron_device": 12519788544,
    "usage_breakdown": {
      "host": {
        "application_memory": 6996594688,
        "constants": 0,
        "dma_buffers": 1048576,
        "tensors": 0
      },
      "neuroncore_memory_usage": {
        "0": {
          "constants": 193986816,
          "model_code": 176285056,
          "model_shared_scratchpad": 0,

```

(continues on next page)

(continued from previous page)

```

        "runtime_memory": 0,
        "tensors": 20971520
    },
    "1": {
        "constants": 193986816,
        "model_code": 176285056,
        "model_shared_scratchpad": 0,
        "runtime_memory": 0,
        "tensors": 20971520
    },
    ...
}
}
"loaded_models": [
{
    "name": "neff",
    "uuid": "91f2f66e83ea419dace1da07617ad39f",
    "model_id": 10005,
    "is_running": false,
    "subgraphs": {
        "sg_00": {
            "memory_used_bytes": {
                "host": 20480,
                "neuron_device": 21001024,
                "usage_breakdown": {
                    "host": {
                        "application_memory": 20480,
                        "constants": 0,
                        "dma_buffers": 0,
                        "tensors": 0
                    },
                    "neuron_device": {
                        "constants": 20971520,
                        "model_code": 29504,
                        "runtime_memory": 0,
                        "tensors": 0
                    }
                }
            },
            "neuroncore_index": 0,
            "neuron_device_index": 12
        }
    },
    "error": ""
}
],
...
],
"error": ""
}

```

- "memory\_used" summarizes the amount of memory used by the Neuron application
  - "neuron\_runtime\_used\_bytes" - current amount of memory used by the Neuron application
    - \* "host" - total host DRAM usage in bytes

- \* "neuron\_device" - total Neuron device memory usage in bytes
- \* "usage\_breakdown" - a breakdown of the total memory usage in the other two fields
  - "host" - breakdown of the host memory usage
  - "application\_memory" - amount of host memory used by the application - this includes all allocations that are not included in the next categories
  - "constants" - amount of host memory used for constants during training (or weights during inference)
  - "dma\_buffers" - amount of host memory used for DMA transfers
  - "tensors" - amount of host memory used for tensors
  - "neuroncore\_memory\_usage" - a breakdown of memory allocated on the Neuron Devices and the NeuronCores for which it was allocated
    - "0" - "64" (for trn2-48xlarge) - NeuronCores for which the memory was allocated
    - "constants" - amount of device memory used for constants during training (or weights during inference)
    - "model\_code" - amount of device memory used for models' executable code
    - "model\_shared\_scratchpad" - amount of device memory used for the scratchpad shared by the models - a memory region reserved for the models'
    - internal variables and auxiliary buffers - "runtime\_memory" - amount of device memory used by the Neuron Runtime - "tensors" - amount of device memory used for tensors
- "loaded\_models" - array containing objects representing loaded models
  - "name" - name of the model
  - "uuid" - unique id for the model
  - "model\_id" - Neuron application-assigned ID for this model
  - "is\_running" - true if this model is currently started, false otherwise
  - "subgraphs" - object containing all the subgraphs for the model, indexed by their name:  
"subgraph\_name": { subgraph\_data }
  - \* "memory\_used\_bytes" - memory usage for this subgraph
    - "host" - total host DRAM usage in bytes
    - "neuron\_device" - total Neuron device DRAM usage in bytes
    - "usage\_breakdown" - a breakdown of memory allocated at load time for this model
    - "host" - breakdown of host memory allocated for this model
    - "application\_memory" - amount of host memory allocated for this model by the Neuron Runtime which doesn't fall in any of the next categories
    - "constants" - amount of host memory used for constants during training (or weights during inference)
    - "dma\_buffers" - host memory allocated for DMA transfers for this model
    - "tensors" - amount of device memory used for tensors at model load time
    - "neuron\_device" - a breakdown of device memory allocated for this model

- "constants" - amount of device memory used for constants during training (or weights during inference)
- "model\_code" - amount of device memory used for the model's executable code
- "runtime\_memory" - amount of device memory used by the Neuron Runtime for this model
- "tensors" - amount of device memory allocated for tensors at this model's load time
- \* "neuroncore\_index" - NeuronCore index on which the subgraph is loaded
- \* "neuron\_device\_index" - Neuron device index on which the subgraph is loaded
- "error" - string containing any error that occurred when collecting the data

### neuron\_runtime\_vcpu\_usage

```
"neuron_runtime_vcpu_usage": {
  "period": 1.030604818,
  "vcpu_usage": {
    "user": 42.01,
    "system": 12.34
  },
  "error": ""
}
```

- "vcpu\_usage" - object showing vCPU usage in percentages for the Neuron application during the captured period
  - "user" - percentage of time spent in user code by this Neuron Application
  - "system" - percentage of time spent in kernel code by this Neuron application
- "error" - string containing any error that occurred when collecting the data

## System level metric groups

### neuron\_hw\_counters

```
"neuron_hw_counters": {
  "period": 1.030359284,
  "neuron_devices": [
    {
      "neuron_device_index": 0,
      "mem_ecc_corrected": 0,
      "mem_ecc_uncorrected": 0,
      "sram_ecc_uncorrected": 0,
      "sram_ecc_corrected": 0
    }
  ],
  "error": ""
},
```

- "neuron\_devices" - array containing ECC data for all Neuron devices
  - "neuron\_device\_index" - Neuron device index

- "mem\_ecc\_corrected" - number of corrected ECC events in the Neuron device's DRAM
  - "mem\_ecc\_uncorrected" - number of uncorrected ECC events in the Neuron device's DRAM
  - "sram\_ecc\_uncorrected" - number of uncorrected ECC events in the Neuron device's SRAM
  - "sram\_ecc\_corrected" - number of corrected ECC events in the Neuron device's SRAM
- "error" - string containing any error that occurred when collecting the data

### vcpu\_usage

```
"vcpu_usage": {
  "period": 0.999974868,
  "average_usage": {
    "user": 32.77,
    "nice": 0,
    "system": 22.87,
    "idle": 39.36,
    "io_wait": 0,
    "irq": 0,
    "soft_irq": 0
  },
  "usage_data": {
    "0": {
      "user": 34.41,
      "nice": 0,
      "system": 27.96,
      "idle": 37.63,
      "io_wait": 0,
      "irq": 0,
      "soft_irq": 0
    },
    "1": {
      "user": 56.84,
      "nice": 0,
      "system": 28.42,
      "idle": 14.74,
      "io_wait": 0,
      "irq": 0,
      "soft_irq": 0
    },
    [...]
  },
  "context_switch_count": 123456,
  "error": ""
}
```

- each vCPU usage object contains the following fields:
  - "user" - percentage of time spent in user code
  - "nice" - percentage of time spent executing niced user code
  - "system" - percentage of time spent executing kernel code
  - "idle" - percentage of time spent idle



- "io\_wait" - percentage of time spent waiting for IO operations
- "irq" - percentage of time spent servicing hardware interrupts
- "soft\_irq" - percentage of time spent servicing software interrupts
- "average\_usage" - contains the average usage across all vCPUs during the captured period
- "usage\_data" - contains per vCPU usage during the captured period
- "context\_switch\_count" - contains the number of vCPU context switches during the captured period
- "error" - string containing any error that occurred when collecting the data

## memory\_info

```
"memory_info": {
  "period": 5.346411129,
  "memory_total_bytes": 49345835008,
  "memory_used_bytes": 16042344448,
  "swap_total_bytes": 0,
  "swap_used_bytes": 0,
  "error": ""
}
```

- "memory\_total\_bytes" - total size of the host memory, in bytes
- "memory\_used\_bytes" - amount of host memory in use, in bytes
- "swap\_total\_bytes" - total size of the host swap file, in bytes
- "swap\_used\_bytes" - amount of swap memory in use, in bytes

## Companion scripts

neuron-monitor is installed with three Python companion scripts: *neuron-monitor-cloudwatch.py*, *neuron-monitor-prometheus.py*, and *neuron-monitor-k8s-info.py* (Beta)

### neuron-monitor-cloudwatch.py

It requires Python3 and the `boto3` Python module. It is installed to: `/opt/aws/neuron/bin/neuron-monitor-cloudwatch.py`.

### Using neuron-monitor-cloudwatch.py

```
neuron-monitor | neuron-monitor-cloudwatch.py --namespace <namespace> --region <region>
```

For example:

```
neuron-monitor | neuron-monitor-cloudwatch.py --namespace neuron_monitor_test --region_
↪ us-west-2
```

### neuron-monitor-prometheus.py

It requires Python3 and the [Prometheus client Python module](#). It is installed to: `/opt/aws/neuron/bin/neuron-monitor-prometheus.py`.

### Using neuron-monitor-prometheus.py

```
neuron-monitor | neuron-monitor-prometheus.py --port <port>
```

For example:

```
neuron-monitor | neuron-monitor-prometheus.py --port 8008
```

The default value for `--port` is `8000`.

If your data visualization framework is Grafana, we provided a [Grafana dashboard](#) which integrates with Prometheus and this script.

### neuron-monitor-k8s-info.py (Beta)

It requires Python3 and the [gRPC Python package](#). It is installed to: `/opt/aws/neuron/bin/neuron-monitor-k8s-info.py`.

---

**Important:** This companion script is in Beta and is disabled by default.

It only works on EKS, and is currently not supported with EKS auto mode.

---

### Using neuron-monitor-k8s-info.py

```
neuron-monitor | neuron-monitor-prometheus.py --port <port> --enable-k8s-info | neuron-  
↪monitor-k8s-info.py --period <seconds>
```

For example:

```
neuron-monitor | neuron-monitor-prometheus.py --port 8008 --enable-k8s-info | neuron-  
↪monitor-k8s-info.py --period 30
```

The default value for `--period` is `15`.

### Running neuron monitor in Kubernetes environment

For running neuron monitor in Kubernetes environment, please refer to instructions [here](#).

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.2.2 Neuron Top User Guide

### Table of contents

- *Overview*
- *Using neuron-top*
  - *Command line arguments*
  - *User interface*

### Overview

`neuron-top` provides useful information about NeuronCore and vCPU utilization, memory usage, loaded models, and Neuron applications.

---

**Note:** `neuron-top` fully supports the newly launched `trn2` instances.

---



---

**Note:** If you are parsing `neuron-top` output in your automation environment, you can now replace it with `neuron-monitor` (*Neuron Monitor User Guide*) which outputs data in a standardized, easier to parse JSON format.

---

### Using neuron-top

#### Command line arguments

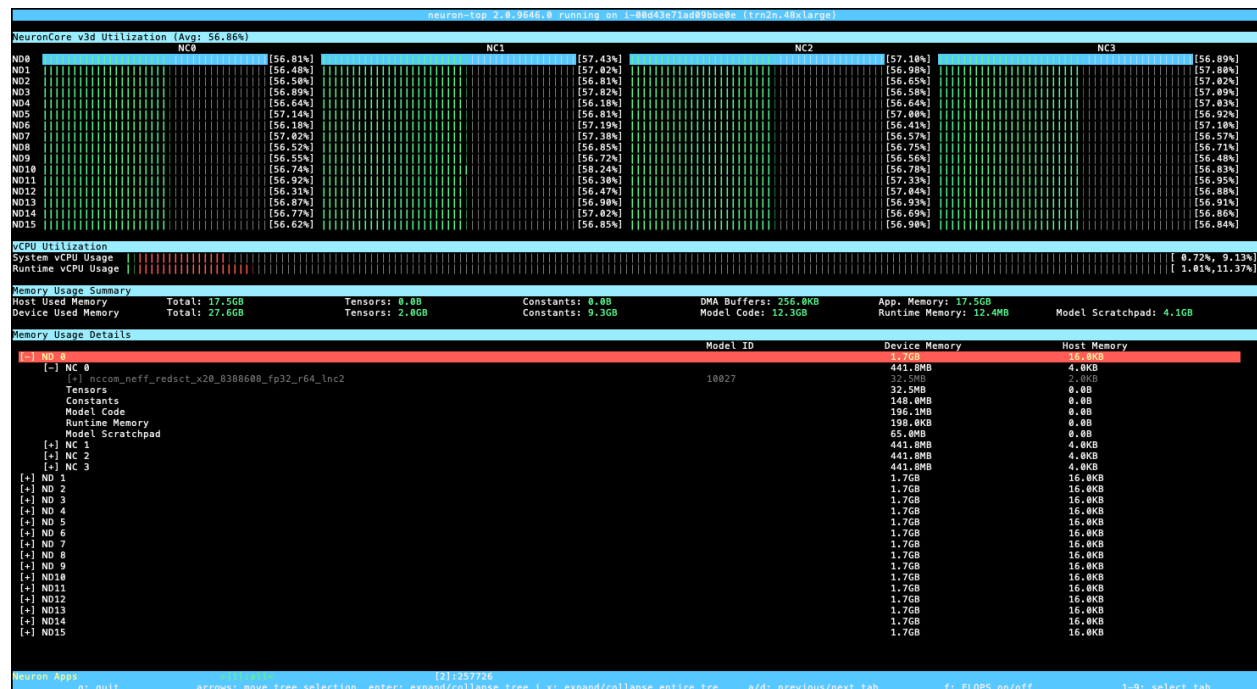
Launch `neuron-top` by simply typing its name in the shell: `neuron-top`.

#### User interface

The title section of the user interface shows the application's version number, EC2 instance ID, and the instance type on which it is running:

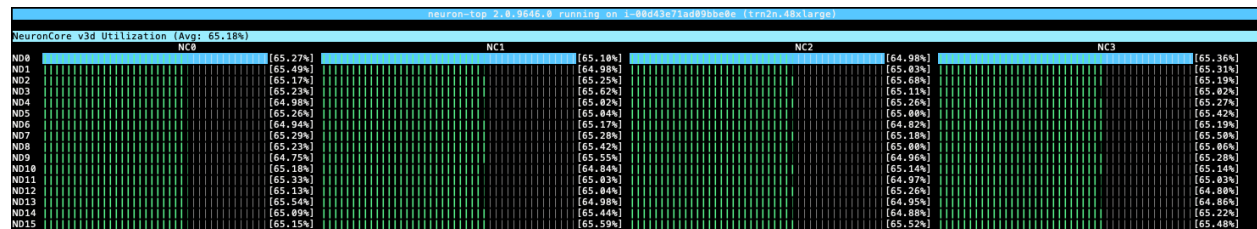
```
neuron-top 2.8.0.046.0 running on i-08d43c71ad09c0e0e (trn2n.4xlarge)
```

The rest of the user interface is divided in 4 sections. The data shown in these sections applies to the currently selected tab - which can be the 'all' tab, which aggregates data from all running Neuron processes, or a tab representing a single Neuron process:



- The NeuronCore <vers> Utilization section shows the NeuronCore utilization for the currently selected tab. <vers> is the version of the NeuronCores on the instance (for example, v2 for trn1 instances and inf2 instances, v3 for trn2 instances with LNC=1, v3d for trn2 instances with LNC=2)

Pressing the 'F' key will toggle between displaying utilization percentages - as seen in the previous image - and teraflops (trillion floating point operations per second), as seen in the image below:



- The VCPU Utilization section shows:
  - System vCPU usage - the two percentages are user% and system%
  - Runtime vCPU usage - same breakdown
- The Memory Usage Summary section provides a breakdown of the total memory usage on the Neuron Device as well as on the host:
  - Host Used Memory - amount of host memory used by the selected application (or an aggregate of all applications if 'All' is selected)
    - \* Total - total amount of host memory used
    - \* Tensors - amount of host memory used for tensors
    - \* Constants - amount of host memory used for constants (for applications running training) or weights (for applications running inferences)
    - \* DMA Buffers - amount of host memory used for DMA transfers

- \* **App. Memory** - amount of host memory used by the application that doesn't fall in any of the previous categories
- **Device Used Memory** - amount of device memory used by the selected application (or an aggregate of all applications if 'All' is selected)
  - \* **Total** - total amount of device memory used
  - \* **Tensors** - amount of device memory used for tensors
  - \* **Constants** - amount of device memory used for constants (for applications running training) or weights (for applications running inferences)
  - \* **Model Code** - amount of device memory used for storing model executable code
  - \* **Runtime Memory** - amount of device memory used by the Neuron Runtime (outside of the previous categories)
  - \* **Model Scratchpad** - amount of device memory used for the shared model scratchpad, a shared buffer used for internal model variables and other auxiliary buffers
- **Memory Usage Details** contains memory usage data organized as a tree which can be expanded/collapsed. The columns are:
  - **Model ID** - the Neuron Runtime identifier for this model instance
  - **Host Memory** - amount of host memory used
  - **Device Memory** - amount of device memory used

The tree view shows the amount of memory used for the same categories shown in the **Memory Usage Summary** but in this section they are attached to either a model (if the memory has been allocated at model load time for that model), or to a NeuronCore (if the memory can't be associated with a model, but has been allocated for that NeuronCore). The 'parent' shows the total amount of memory used - the sum of its children.

---

**Note:** The up/down/left/right keys can be used to navigate the tree view. The 'x' key expands/collapses the entire tree.

---

The bottom bar shows which Neuron process' data is currently displayed by highlighting its tag using a green font and marking it using a pair of '>', '<' characters. The 'all' tab shows an aggregated view of all the Neuron processes currently running on the instance.



---

**Note:** The '1'-'9' keys select the current tab. 'a'/'d' selects the previous/next tab on the bar.

---

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.2.3 Neuron LS User Guide

The `neuron-ls` command is a tool for managing Neuron devices in your instance. This command serves two key purposes: it identifies all Neuron devices present in the current instance and provides information about the processes running on each device along with the command that launched that process. To use this command, simply type `neuron-ls` in your terminal.

#### neuron-ls CLI

**neuron-ls** [options]

**Available options:**

- `--wide`, `-w`: Displays the table in a wider format.
- `--show-all-procs`, `-a`: Show all processes using the Neuron Devices, including processes that aren't using Neuron Runtime 2.x such as `neuron-monitor` or `neuron-ls` itself.
- `--topology`, `-t`: Display topology information about the system's Neuron Devices.
- `--json-output`, `-j`: Output in JSON format.

---

**Note:** `neuron-ls` fully supports the newly launched Trn2 instances.

---

#### Examples

`neuron-ls` is compatible with all Neuron instance types: `inf1`, `inf2`, `trn1` and `trn2`. These are a few examples on running the tool on a `trn2n.48xlarge`:

```
$ neuron-ls
instance-type: trn2n.48xlarge
instance-id: i-aabbccdd123456789
logical-neuroncore-config: 2
```

| NEURON<br>DEVICE | NEURON<br>CORES | NEURON<br>MEMORY | CONNECTED<br>DEVICES | PCI<br>BDF |
|------------------|-----------------|------------------|----------------------|------------|
| 0                | 4               | 96 GB            | 12, 3, 4, 1          | cc:00.0    |
| 1                | 4               | 96 GB            | 13, 0, 5, 2          | b5:00.0    |
| 2                | 4               | 96 GB            | 14, 1, 6, 3          | b6:00.0    |
| 3                | 4               | 96 GB            | 15, 2, 7, 0          | cb:00.0    |
| 4                | 4               | 96 GB            | 0, 7, 8, 5           | 6f:00.0    |
| 5                | 4               | 96 GB            | 1, 4, 9, 6           | 58:00.0    |
| 6                | 4               | 96 GB            | 2, 5, 10, 7          | 59:00.0    |
| 7                | 4               | 96 GB            | 3, 6, 11, 4          | 6e:00.0    |
| 8                | 4               | 96 GB            | 4, 11, 12, 9         | 9b:00.0    |
| 9                | 4               | 96 GB            | 5, 8, 13, 10         | 84:00.0    |
| 10               | 4               | 96 GB            | 6, 9, 14, 11         | 85:00.0    |
| 11               | 4               | 96 GB            | 7, 10, 15, 8         | 9a:00.0    |
| 12               | 4               | 96 GB            | 8, 15, 0, 13         | f8:00.0    |
| 13               | 4               | 96 GB            | 9, 12, 1, 14         | e1:00.0    |
| 14               | 4               | 96 GB            | 10, 13, 2, 15        | e2:00.0    |
| 15               | 4               | 96 GB            | 11, 14, 3, 12        | f7:00.0    |

```

$ neuron-ls --wide
instance-type: trn2n.48xlarge
instance-id: i-aabbccdd123456789
logical-neuroncore-config: 2
+-----+-----+-----+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | CONNECTED | PCI | PID | |
| COMMAND | RUNTIME | | | | |
| DEVICE | CORES | MEMORY | DEVICES | BDF | | |
| VERSION |
+-----+-----+-----+-----+-----+-----+
| 0 | 4 | 96 GB | 12, 3, 4, 1 | cc:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 1 | 4 | 96 GB | 13, 0, 5, 2 | b5:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 2 | 4 | 96 GB | 14, 1, 6, 3 | b6:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 3 | 4 | 96 GB | 15, 2, 7, 0 | cb:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 4 | 4 | 96 GB | 0, 7, 8, 5 | 6f:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 5 | 4 | 96 GB | 1, 4, 9, 6 | 58:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 6 | 4 | 96 GB | 2, 5, 10, 7 | 59:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 7 | 4 | 96 GB | 3, 6, 11, 4 | 6e:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 8 | 4 | 96 GB | 4, 11, 12, 9 | 9b:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 9 | 4 | 96 GB | 5, 8, 13, 10 | 84:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 10 | 4 | 96 GB | 6, 9, 14, 11 | 85:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 11 | 4 | 96 GB | 7, 10, 15, 8 | 9a:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 12 | 4 | 96 GB | 8, 15, 0, 13 | f8:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 13 | 4 | 96 GB | 9, 12, 1, 14 | e1:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 14 | 4 | 96 GB | 10, 13, 2, 15 | e2:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
| 15 | 4 | 96 GB | 11, 14, 3, 12 | f7:00.0 | 268911 | neuron-bench exec --run-
| as-cc-neff --warmup none --fixed-instance-count 64 --... | 2.0.0 |
+-----+-----+-----+-----+-----+-----+

```

```

$ neuron-ls --show-all-procs
instance-type: trn2n.48xlarge
instance-id: i-aabbccdd123456789
logical-neuroncore-config: 2
+-----+-----+-----+-----+-----+-----+-----+

```

(continues on next page)

(continued from previous page)

| NEURON     | NEURON | NEURON            | CONNECTED   | PCI     | PID    | COMMAND                  |
|------------|--------|-------------------|-------------|---------|--------|--------------------------|
| DEVICE     | CORES  | RUNTIME<br>MEMORY | DEVICES     | BDF     |        |                          |
|            |        | VERSION           |             |         |        |                          |
| 0          | 4      | 96 GB             | 12, 3, 4, 1 | cc:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 1          | 4      | 96 GB             | 13, 0, 5, 2 | b5:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 2          | 4      | 96 GB             | 14, 1, 6, 3 | b6:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 3          | 4      | 96 GB             | 15, 2, 7, 0 | cb:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 4          | 4      | 96 GB             | 0, 7, 8, 5  | 6f:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 5          | 4      | 96 GB             | 1, 4, 9, 6  | 58:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 6          | 4      | 96 GB             | 2, 5, 10, 7 | 59:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        | NA                |             |         |        |                          |
| 7          | 4      | 96 GB             | 3, 6, 11, 4 | 6e:00.0 | 268911 | neuron-bench exec --run- |
| as-cc-neff | --...  | 2.0.0             |             |         | 269192 | neuron-ls --show-all-    |
| procs      |        |                   |             |         |        |                          |

(continues on next page)



(continued from previous page)

```

↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 8          | 4          | 96 GB | 4, 11, 12, 9 | 9b:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 9          | 4          | 96 GB | 5, 8, 13, 10 | 84:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 10         | 4          | 96 GB | 6, 9, 14, 11 | 85:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 11         | 4          | 96 GB | 7, 10, 15, 8 | 9a:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 12         | 4          | 96 GB | 8, 15, 0, 13 | f8:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 13         | 4          | 96 GB | 9, 12, 1, 14 | e1:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 14         | 4          | 96 GB | 10, 13, 2, 15 | e2:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+
| 15         | 4          | 96 GB | 11, 14, 3, 12 | f7:00.0 | 268911 | neuron-bench exec --run-
↪as-cc-neff --... | 2.0.0      |
|           |           |       |              |         | 269192 | neuron-ls --show-all-
↪procs          | NA          |
+-----+-----+-----+-----+-----+-----+
↪-----+-----+

```

```
$ neuron-ls --topology
instance-type: trn2n.48xlarge
instance-id: i-aabbccdd123456789
logical-neuroncore-config: 2
```

| NEURON<br>DEVICE | NEURON<br>CORES | NEURON<br>MEMORY | CONNECTED<br>DEVICES | PCI<br>BDF |
|------------------|-----------------|------------------|----------------------|------------|
| 0                | 4               | 96 GB            | 12, 3, 4, 1          | cc:00.0    |
| 1                | 4               | 96 GB            | 13, 0, 5, 2          | b5:00.0    |
| 2                | 4               | 96 GB            | 14, 1, 6, 3          | b6:00.0    |
| 3                | 4               | 96 GB            | 15, 2, 7, 0          | cb:00.0    |
| 4                | 4               | 96 GB            | 0, 7, 8, 5           | 6f:00.0    |
| 5                | 4               | 96 GB            | 1, 4, 9, 6           | 58:00.0    |
| 6                | 4               | 96 GB            | 2, 5, 10, 7          | 59:00.0    |
| 7                | 4               | 96 GB            | 3, 6, 11, 4          | 6e:00.0    |
| 8                | 4               | 96 GB            | 4, 11, 12, 9         | 9b:00.0    |
| 9                | 4               | 96 GB            | 5, 8, 13, 10         | 84:00.0    |
| 10               | 4               | 96 GB            | 6, 9, 14, 11         | 85:00.0    |
| 11               | 4               | 96 GB            | 7, 10, 15, 8         | 9a:00.0    |
| 12               | 4               | 96 GB            | 8, 15, 0, 13         | f8:00.0    |
| 13               | 4               | 96 GB            | 9, 12, 1, 14         | e1:00.0    |
| 14               | 4               | 96 GB            | 10, 13, 2, 15        | e2:00.0    |
| 15               | 4               | 96 GB            | 11, 14, 3, 12        | f7:00.0    |

#### Neuron Device Topology

```

      *           *           *           *
      |           |           |           |

*--[ 0 ]--[ 1 ]--[ 2 ]--[ 3 ]--*

      |           |           |           |

*--[ 4 ]--[ 5 ]--[ 6 ]--[ 7 ]--*

      |           |           |           |

*--[ 8 ]--[ 9 ]--[10 ]--[11 ]--*

      |           |           |           |

*--[12 ]--[13 ]--[14 ]--[15 ]--*

      |           |           |           |
      *           *           *           *
```

#### Legend:

\*-- = Wrap-around link

```
$ neuron-ls -j
[
  {
    "neuron_device": 0,
    "bdf": "cc:00.0",
    "connected_to": [
      12,
      3,
      4,
      1
    ],
    "nc_count": 4,
    "logical_neuroncore_config": 2,
    "memory_size": 103079215104,
    "neuron_processes": [
      {
        "pid": 113985,
        "command": "neuron-bench exec --run-as-cc-neff --...",
        "neuron_runtime_version": "2.0.0"
      }
    ]
  },
  ...
  {
    "neuron_device": 15,
    "bdf": "f7:00.0",
    "connected_to": [
      11,
      14,
      3,
      12
    ],
    "nc_count": 4,
    "logical_neuroncore_config": 2,
    "memory_size": 103079215104,
    "neuron_processes": [
      {
        "pid": 113985,
        "command": "neuron-bench exec --run-as-cc-neff --...",
        "neuron_runtime_version": "2.0.0"
      }
    ]
  }
]
```

- instance-type: Type of instance on which neuron-ls is running.
- instance-id: EC2 ID of the instance on which neuron-ls is running.
- logical-neuroncore-config: (only available on trn2 instances) the current logical NeuronCore configuration; for more information refer to [Logical NeuronCore configuration](#)
- NEURON DEVICE / neuron\_device: Logical ID assigned to the Neuron Device.
- NEURON CORES / nc\_count: Number of NeuronCores present in the Neuron Device.

- NEURON MEMORY / `memory_size`: Amount DRAM memory in Neuron Device.
- CONNECTED DEVICES / `connected_to`: Logical ID of Neuron Devices connected to this Neuron Device.
- PCI BDF / `bdf`: PCI Bus Device Function (BDF) ID of the device.
- PID / `pid`: ID of the process using this NeuronDevice.
- COMMAND / `command`: Command used to launch the process using this Neuron Device.
- RUNTIME VERSION / `neuron_runtime_version`: Version of Neuron Runtime (if applicable) for the application using this Neuron Device.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.2.4 Neuron Sysfs User Guide

### Table of contents

- *Introduction*
- *Neuron Sysfs Filesystem Structure*
  - *High Level Overview*
  - *Description for Each Field*
  - *Read and Write to Sysfs*
  - *Note*
- *How to Troubleshoot via Sysfs*

### Introduction

The kernel provides a few ways in which userspace programs can get system information from the kernel space. Sysfs is one common way to do so. It is a virtual filesystem typically mounted on the `/sys` directory and contains information about hardware devices attached to the system and about drivers handling those devices. By navigating the hierarchical structure of the sysfs filesystem and viewing the information provided by its files and directories, you can gather valuable information that can help diagnose and resolve a wide range of hardware and system issues.

Thus a sysfs filesystem is set up per Neuron Device under `/sys/devices/virtual/neuron_device` to give you an insight into the Neuron Driver and Runtime at system level. By performing several simple CLIs such as reading or writing to a sysfs file, you can get information such as Runtime status, memory usage, Driver info etc. You can even create your own shell scripts to query Runtime and Driver statistics from sysfs and generate customized reports.

This user guide will first explain the Neuron sysfs structure and then introduce many ways where you can perform diagnostic works with Neuron sysfs.

## Neuron Sysfs Filesystem Structure

### High Level Overview

Here is the high level structure of the Neuron sysfs filesystem, where the total and present counters are not shown:

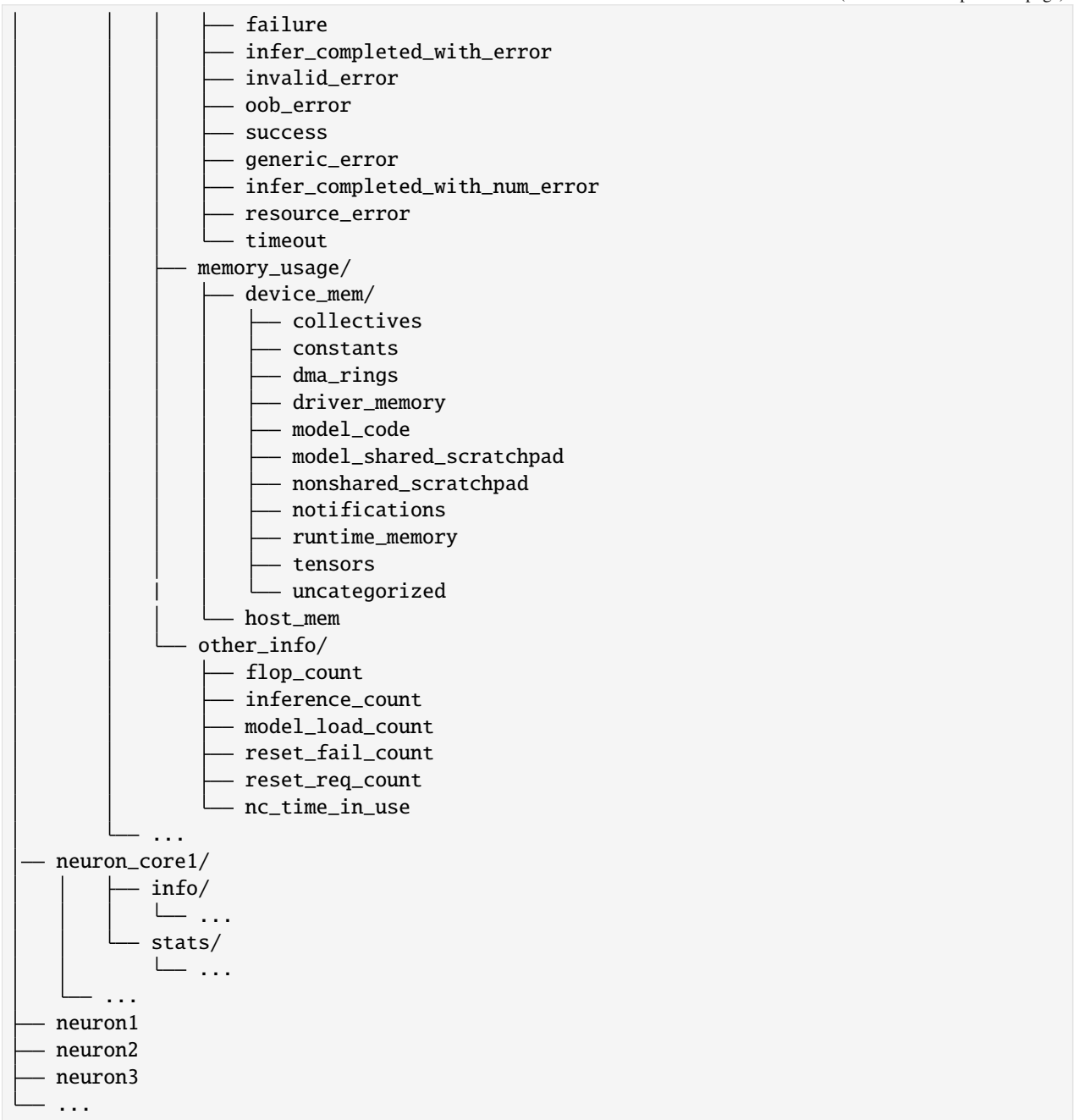
```

/sys/devices/virtual/neuron_device/
├── neuron0/
│   ├── subsystem
│   ├── uevent
│   ├── connected_devices
│   ├── core_count
│   ├── reset
│   ├── power/
│   │   ├── async
│   │   ├── control
│   │   ├── runtime_active_time
│   │   ├── runtime_active_kids
│   │   └── ...
│   ├── info/
│   │   ├── notify_delay
│   │   ├── serial_number
│   │   ├── architecture/
│   │   │   ├── arch_type
│   │   │   ├── device_name
│   │   │   └── instance_type
│   └── stats
│       ├── hardware
│       │   ├── mem_ecc_uncorrected
│       │   └── sram_ecc_uncorrected
│       ├── memory_usage
│       │   └── host_mem
│       │       ├── application_memory
│       │       ├── constants
│       │       ├── dma_buffers
│       │       ├── dma_rings
│       │       ├── driver_memory
│       │       ├── notifications
│       │       ├── tensors
│       │       └── uncategorized
│       └── power
│           └── utilization
├── neuron_core0/
│   ├── info/
│   │   └── architecture/
│   │       └── arch_type
│   └── stats/
│       └── status/
│           ├── exec_bad_input
│           ├── hw_error
│           ├── infer_failed_to_queue
│           ├── resource_nc_error
│           └── unsupported_neff_version

```

(continues on next page)

(continued from previous page)



Each Neuron Device is represented as a directory under `/sys/devices/virtual/neuron_device/`, where `neuron0/` represents the Neuron Device 0, `neuron1/` represents the Neuron Device 1, etc. Each NeuronCore is represented as a directory under a Neuron Device directory, represented as `neuron_core{0,1,2,...}`. Metrics such as Runtime and Driver info and statistics are collected as per NeuronCore in two directories under the NeuronCore directory, i.e. `info/` and `stats/`.

Most of the metrics belong to a category called “counter.” Each counter is represented as a directory, which holds two numerical values as two files: `total` and `present`. Each memory usage counter has an additional value called `peak`. The `total` value starts accumulating metrics when the Driver is loaded. The `present` value records the last changed metric value. The `peak` value records the max value so far. Each counter has the same filesystem structure like this:

```
/sys/devices/virtual/neuron_device/neuron0/neuron_core0/status/
├── exec_bad_input/
│   ├── total
│   └── present
├── hw_error/
│   ├── total
│   └── present
├── infer_failed_to_queue/
│   ├── total
│   └── present
└── ...
```

## Description for Each Field

**info/:** This directory stores general information about hardware and software. None of them are counter types.

- **notify\_delay:** The delay between notifications from the Neuron Device. Current settings are on (0) or off (-1). Off by default.
- **serial\_number:** The unique device identifier.
- **architecture/:** This directory stores hardware architecture information.
  - **arch\_type:** The architecture type of the Neuron Device. Sample architecture types are v1, v2, and v3. You can only read the value. You cannot change it.
  - **instance\_type:** The instance type of the Neuron Device. Sample instance types are Inf1, Inf2, and Trn1. You can only read the value. You cannot change it.
  - **device\_type:** The Neuron Device type. Sample Neuron Device types are Inferentia, Inferentia2, and Trainium1. You can only read the value. You cannot change it.

**stats/:** This directory stores Neuron Runtime and Driver statistics. It contains three subdirectories: **status/**, **memory\_usage/**, and **other\_info/**.

- **status/:** This directory stores the number of each return status of API calls. As explained in [The LIBNRT API Return Codes](#), every API call returns an NRT\_STATUS value, which represents the return status of that API call. Our sysfs filesystem stores all NRT\_STATUS as subdirectories under the **status/** directory. They all have the counter structure. Thus each NRT\_STATUS subdirectory holds two values (total and present) and records the number of times you receive a certain NRT\_STATUS. The following is description for each NRT\_STATUS subdirectory. You should see the description align with what is described in [The LIBNRT API Return Codes](#).
- **memory\_usage/:** This directory contains memory usage statistics for both device and host, represented as counters. In this directory, the total counters indicate the current memory usage, present counters represent the memory allocation or deallocation amount in the previous operation, and peak counters indicate the maximum memory usage observed. Additionally, this directory provides detailed breakdown statistics for device and host memory usage. These memory breakdown details correspond to the [Memory Usage Summary](#) section displayed on in Neuron Monitor.
  - **device\_mem/:** The amount of memory that Neuron Runtime uses for weights, instructions and DMA rings.
    - \* **This device memory per NeuronCore is further categorized into five types: collectives/, constants/, dma\_rings/, driver\_memory/, model\_code/, model\_shared\_scratchpad/, nonshared\_scratchpad/, notifications/, runtime\_memory/, tensors/, and uncategorized/. Each of these categories has total, present, and peak.**

- **collectives** - amount of device memory used for collective communication between workers
  - **constants** - amount of device memory used for constants (for applications running training) or weights (for applications running inferences)
  - **dma\_rings** - amount of device memory used for storing model executable code used for data movements
  - **driver\_memory** - amount of device memory used by the Neuron Driver
  - **model\_code** - amount of device memory used for storing model executable code
  - **model\_shared\_scratchpad** - amount of device memory used for the shared model scratchpad, a buffer shared between models on the same Neuron Core used for internal model variables and other auxiliary buffers
  - **nonshared\_scratchpad** - amount of device memory used for non-shared model scratchpad, a buffer used by a single model for internal model variables and other auxiliary buffers
  - **notifications** - amount of device memory used to store instruction level trace information used to profile workloads ran on the device
  - **runtime\_memory** - amount of device memory used by the Neuron Runtime (outside of the previous categories)
  - **tensors** - amount of device memory used for tensors
  - **uncategorized** - amount of device memory that does not belong in any other category in this list
- **host\_mem/**: The amount of memory that Neuron Runtime uses for input and output tensors.
    - \* The host memory per Neuron Device is further categorized into four types: **application\_memory/**, **constants/**, **dma\_buffers/**, **dma\_rings/**, **driver\_memory/**, **notifications/**, **tensors/**, **uncategorized/**. These categories provide more granular host memory classification compared to *Host Used Memory* section. Each of these categories has total, present, and peak
  - **hardware/**: Hardware statistics.
    - \* **mem\_ecc\_uncorrected**: The number of uncorrected ECC events in the Neuron device's DRAM.
    - \* **sram\_ecc\_uncorrected**: The number of uncorrected ECC events in the Neuron device's SRAM.
  - **power/**: Power statistics.
    - \* **utilization**: Reports per-minute power usage statistics as a percentage of max power in the following format:

<status>,<timestamp>,<min\_power>,<max\_power>,<avg\_power>

#### **Field descriptions:**

##### **status**

Indicates the sampling state in a string. Valid values are:

POWER\_STATUS\_VALID - Sampling successful

POWER\_STATUS\_NO\_DATA - No samples available

POWER\_STATUS\_INVALID - An internal sampling error occurred

##### **timestamp**

Time when the sample was collected in Unix epoch seconds (integer)



**min\_power**

Minimum power utilization during the sampling period (0.00-100.00%)

**max\_power**

Maximum power utilization during the sampling period (0.00-100.00%)

**avg\_power**

Average power utilization during the sampling period (0.00-100.00%)

The interface updates these statistics every minute based on continuous power sampling.

- **other\_info/**: This directory contains statistics that are not included by **status/** and **memory\_usage/**. None of them are counter types.
  - **flop\_count**: The number of flops. You can use it to calculate the TFLOP/s by **flop\_count** / time interval
  - **inference\_count**: The number of successful inferences
  - **model\_load\_count**: The number of successful model loads
  - **reset\_fail\_count**: The number of failed device resets
  - **reset\_req\_count**: The number of device resets requests
  - **nc\_time\_in\_use**: The time interval in microseconds between the start and the end of the current execution on hardware

Other fields:

- **connected\_devices**: The list of connected devices' ids. You should see the same output as **neuron-ls**'s **CONNECTED DEVICES**.
- **reset**: write to this file resets corresponding the Neuron Device.

## Read and Write to Sysfs

Reading a sysfs file gives the value for the corresponding metric. You can use the **cat** command to view the contents of the sysfs files.:

```
ubuntu@ip-xxx-xx-xx-xxx:~$ sudo cat /sys/devices/virtual/neuron_device/neuron0/neuron_
↪core0/stats/status/failure/total
0
ubuntu@ip-xxx-xx-xx-xxx:~$ sudo cat /sys/devices/virtual/neuron_device/neuron0/neuron_
↪core0/info/architecture/arch_type
NCv2
```

Sysfs metrics of counter type are write to clear. You can write any value to the file, and the metric will be set to 0:

```
ubuntu@ip-xxx-xx-xx-xxx:~$ echo 1 | sudo tee /sys/devices/virtual/neuron_device/neuron0/
↪neuron_core0/stats/status/failure/total
1
```

Writing to **reset** resets the corresponding Neuron Device. E.g. the below resets Neuron Device 0:

```
ubuntu@ip-xxx-xx-xx-xxx:~$ echo 1 | sudo tee /sys/devices/virtual/neuron_device/neuron0/
↪reset
1
```

## Note

All files under `/sys/devices/virtual/neuron_device/neuron0/power` such as `runtime_active_kids` or `runtime_status` are related to generic device power management. They are not created or controlled by our sysfs metrics. The word `runtime` in these files does not refer to Neuron Runtime.

## How to Troubleshoot via Sysfs

You can perform simple and easy tasks to troubleshoot your ML jobs with one or a few CLIs to read or write the sysfs filesystem. You can do aggregations across all the NeuronCores and all the Neuron Device to get a summarized view using your scripts.

You can also use the Sysfs notification feature to wait passively (without wasting CPU cycles) for changes to the values of Sysfs files. To use this feature, you need to implement a user-space program that calls the `poll()` function on the Sysfs file that you want to wait on. The `poll()` function has the following signature: `unsigned int (*poll) (struct file *, struct poll_table_struct *)`. By default, the Sysfs notification feature is turned off when the driver is loaded. To enable notifications, you can set the value of `/sys/devices/virtual/neuron_device/neuron0/info/notify_delay` to 0. To disable notifications, you can set it to -1. Please note that enabling this feature can impact performance.

Here is a sample user space program using `poll()`:

```
#include <fcntl.h>
#include <poll.h>
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>

int main(int argc, char * argv[])
{
    char readbuf[128];
    int attr_fd = -1;
    struct pollfd pfd;
    int retval = 0;
    ssize_t read_bytes;

    if (argc < 2) {
        fprintf(stderr, "Error: Please specify sysfs file path\n");
        exit(1);
    }
    attr_fd = open(argv[1], O_RDONLY, 0);
    if (attr_fd < 0) {
        perror(argv[1]);
        exit(2);
    }

    read_bytes = read(attr_fd, readbuf, sizeof(readbuf));
    if (read_bytes < 0) {
        perror(argv[1]);
        exit(3);
    }
    printf("%.5s", (int)read_bytes, readbuf);
}
```

(continues on next page)

(continued from previous page)

```

pfd.fd = attr_fd;
pfd.events = POLLERR | POLLPRI;
pfd.revents = 0;
while ((retval = poll(&pfd, 1, 100)) >= 0) {
    if (pfd.revents & (POLLERR | POLLPRI)) {
        pfd.revents = 0;

        lseek(attr_fd, 0, SEEK_SET);
        read_bytes = read(attr_fd, readbuf, sizeof(readbuf));
        if (read_bytes < 0) {
            perror(argv[1]);
            exit(4);
        }
        printf("%.5s", (int)read_bytes, readbuf);
    }
}
return 0;
}

```

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## 6.2.5 NCCOM-TEST User Guide

### Table of contents

- *Overview*
- *Using nccom-test*
  - *Output description*
  - *CLI arguments*
  - *Examples*

### Overview

**nccom-test** is a benchmarking tool for evaluating Collective Communication operations on AWS Trainium and Inferentia instances. It supports Trn1, Trn2, and Inf2 instance types. The tool can assess performance across multiple instances or perform quick environment sanity checks before running more complex workloads. While single-instance benchmarking is supported for all compatible instance types, multi-instance benchmarking is limited to Trainium instances (Trn1 and Trn2).

**Note:** On Inf2 instances, only single-instance benchmarking is supported. Running a multi-node nccom-test benchmark will result in an error.

## Using nccom-test

Here is a simple example which will run a 2 worker (ranks) all-reduce with a total size of 32MB:

```
nccom-test -r 2 allr
      size(B)    count(elems)    type    time(us)    algbw(GB/s)    busbw(GB/s)
      33554432      33554432    uint8      768          40.69          40.69
Avg bus bandwidth:    40.6901GB/s
```

## Output description

The command will output a table containing several columns containing performance metrics. There will be a line for every requested data size (by default the data size is 32MB as seen in the previous example).

| Column name       | Description   |
|-------------------|---|
| size(B)           | Size in bytes for the data involved in this operation   |
| count(elems)      | Number of elements in the data involved in this operation. For example, if <b>size(B)</b> is 4 and <b>type</b> is fp32, then <b>count</b> will be 1 since one single fp32 element has been processed.   |
| type              | Data type for the processed data. Can be: <b>uint8</b> , <b>int8</b> , <b>uint16</b> , <b>int16</b> , <b>fp16</b> , <b>bf16</b> , <b>int32</b> , <b>uint32</b> , <b>fp32</b>  |
| time(us)          | Time in microseconds representing the P50 of all durations for the Collective Communication operations executed during the benchmark.   |
| algbw(GB/s)       | Algorithm bandwidth in gibibytes (1GiB = 1,073,741,824 bytes) per second which is calculated as <b>size(B) / time(us)</b>   |
| busbw(GB/s)       | Bus bandwidth - bandwidth per data line in gibibytes per second - it provides a bandwidth number that is independent from the number of ranks (unlike <b>algbw</b> ). For a more in-depth explanation on bus Bandwidth, please refer to <a href="#">NVIDIA's nccl-tests documentation</a> . |
| Avg bus bandwidth | Average of the values in the busbw column   |

## CLI arguments

| Argument                              | Default value                  | Description   |
|---------------------------------------|--------------------------------|---|
| <cc operation>                        | N/A, required argument         | The type of Collective Communication operation to execute for this benchmark. Supported types: <ul style="list-style-type: none"> <li>• <code>all_reduce / allr</code>: All-Reduce</li> <li>• <code>all_gather / allg</code>: All-Gather</li> <li>• <code>reduce_scatter / redsct</code>: Reduce-Scatter</li> <li>• <code>sendrecv</code>: Send-Receive</li> <li>• <code>alltoall</code>: All-to-All</li> </ul>   |
| <code>-r, --nworkers</code>           | N/A, required argument         | Total number of workers (ranks) to use  |
| <code>-N, --nnodes</code>             | 1                              | Total number of nodes (instances) to use. The number of workers will be divided equally across all nodes. If this argument is greater than 1, the <b>NEURON_RT_ROOT_COMM_ID</b> environment variable needs to be set to the host address of the instance <b>nccom-test</b> is ran on, and a free port number (for example: <code>NEURON_RT_ROOT_COMM_ID=10.0.0.1:44444</code> ). Additionally, either <code>-s, --hosts</code> needs to be provided or a <code>~/hosts</code> file needs to exist - for more details refer to the <code>-s, --hosts</code> description below. |
| <code>-b, --minbytes</code>           | 32M                            | The starting size for the benchmark   |
| <code>-e, --maxbytes</code>           | 32M                            | The end size for the benchmark. <b>nccom-test</b> will run benchmarks for all sizes between <code>-b, --minbytes</code> and <code>-e, --maxbytes</code> , increasing the size by either <code>-i, --stepbytes</code> or <code>--f, --stepfactor</code> with every run.  |
| <code>-i, --stepbytes</code>          | $(--maxbytes - minbytes) / 10$ | Amount of bytes with which to increase the benchmark's size on every subsequent run. For example, for this combination of arguments: <code>-b 8 -e 16 -i 4</code> , the benchmark will be ran for the following sizes: 8 bytes, 12 bytes, 16 bytes.   |
| <code>-f, --stepfactor</code>         | N/A                            | Factor with which to increase the benchmark's size on every subsequent run. For example, for this combination of argument values: <code>-b 8 -e 32 -f 2</code> , the benchmark will be ran for the following sizes: 8 bytes, 16 bytes, 32 bytes.  |
| <code>-n, --iters</code>              | 20                             | Number of Collective Communication operations to execute during the benchmark.  |
| <code>-w, --warmup_iters</code>       | 5                              | Number of Collective Communication operations to execute as warmup during the benchmark (which won't be counted towards the result).  |
| <code>-d, --datatype</code>           | uint8                          | Data type for the data used by the benchmark. Supported types: <code>uint8, int8, uint16, int16, fp16, bf16, uint32, int32, fp32</code> . Input data will be zero filled, unless <code>--check</code> is provided (currently, only available for <code>--datatype fp32</code> ) in which case it will be filled by a repeated value of the requested type.  |
| <code>-c, --check</code>              | false                          | If provided, the correctness of the operations will be checked. This will not impact results (time, algbw and busbw) but will slightly increase the overall execution time.   |
| <code>-s, --hosts</code>              | N/A                            | Hosts on which to run execution. Checks <code>~/hosts</code> if not specified.  |
| <code>--non-interactive</code>        | true                           | Do not display any animation or progress indicator.   |
| <code>--report-to-json-file</code>    | N/A                            | Persist config and results to JSON file if a filepath is provided.  |
| <code>--show-input-output-size</code> | false                          | Print or save to JSON per rank input and output sizes in B.   |

**Note:** All arguments that take a size in bytes will also accept larger size units, for example: `-f 2048` can be written as `-f 2kb` or `-f 1048576` can be written as `-f 1MB`.

## Examples

**Note:** Performance data shown in these examples should not be considered up-to-date. For the latest performance data, please refer to the performance section.

### Single Instance Examples

- Quick environment validation

```
nccom-test -r 2 allr
      size(B)    count(elems)    type    time(us)    algbw(GB/s)
→busbw(GB/s)
      33554432      33554432    uint8      768      40.69
→40.69
Avg bus bandwidth:      40.6901GB/s
```

If a problem was found, it can be reported in two possible ways:

- Immediately:

```
nccom-test -r 2 allr
Neuron DKMS Driver is not running! Read the troubleshooting
→guide at: https://awsdocs-neuron.readthedocs-hosted.com/
→en/latest/neuron-runtime/nrt-troubleshoot.html#neuron-
→driver-installation-fails
```

- After a benchmark attempt:

```
nccom-test -r 2 allr
      size(B)    count(elems)    type    time(us)
→algbw(GB/s)    busbw(GB/s)
      33554432    Failure running neuron-bench - log file /
→tmp/nccom_test_log_7pqpdfjf.log
1 errors found - test failed
```

In this case, further information about the error can be found in the `neuron-bench` log file.

- 2 rank all-reduce on a single instance for sizes ranging from 1MiB to 1GiB with a step of 4x

```
nccom-test -r 2 --minbytes 1kb --maxbytes 1gb --stepfactor 4 --datatype
→fp32 allr
      size(B)    count(elems)    type    time(us)    algbw(GB/s)
→busbw(GB/s)
      1024      256    fp32      58      0.02
→ 0.02
```

(continues on next page)

(continued from previous page)

|                    |            |             |      |       |       |   |
|--------------------|------------|-------------|------|-------|-------|---|
|                    | 4096       | 1024        | fp32 | 58    | 0.07  | ✓ |
| → 0.07             |            |             |      |       |       |   |
|                    | 16384      | 4096        | fp32 | 58    | 0.26  | ✓ |
| → 0.26             |            |             |      |       |       |   |
|                    | 65536      | 16384       | fp32 | 58    | 1.05  | ✓ |
| → 1.05             |            |             |      |       |       |   |
|                    | 262144     | 65536       | fp32 | 60    | 4.07  | ✓ |
| → 4.07             |            |             |      |       |       |   |
|                    | 1048576    | 262144      | fp32 | 68    | 14.36 | ✓ |
| → 14.36            |            |             |      |       |       |   |
|                    | 4194304    | 1048576     | fp32 | 107   | 36.51 | ✓ |
| → 36.51            |            |             |      |       |       |   |
|                    | 16777216   | 4194304     | fp32 | 332   | 47.06 | ✓ |
| → 47.06            |            |             |      |       |       |   |
|                    | 67108864   | 16777216    | fp32 | 1214  | 51.48 | ✓ |
| → 51.48            |            |             |      |       |       |   |
|                    | 268435456  | 67108864    | fp32 | 4750  | 52.63 | ✓ |
| → 52.63            |            |             |      |       |       |   |
|                    | 1073741824 | 268435456   | fp32 | 18930 | 52.83 | ✓ |
| → 52.83            |            |             |      |       |       |   |
| Avg bus bandwidth: |            | 23.6671GB/s |      |       |       |   |

- 32 rank all-gather on a single instance for sizes ranging from 1KiB to 1MiB with a step of 8x, with correctness checking

```

nccom-test -r 32 --minbytes 1kb --maxbytes 1mb --stepfactor 8 --datatype fp32 --
→check allg
  size(B)    count(elems)    type    time(us)    algbw(GB/s)    busbw(GB/s)
  1024        256          fp32      151          0.01           0.01
  8192       2048          fp32      149          0.05           0.05
  65536      16384          fp32      150          0.41           0.39
  524288     131072          fp32      179          2.73           2.64
Avg bus bandwidth:    0.7731GB/s

```

- Reporting input and output size explicitly with --show-input-output-size

```

nccom-test -r 32 --minbytes 1kb --maxbytes 1mb --stepfactor 8 --datatype fp32 --check_
→allg --show-input-output-size
size(B)    count(elems)    total_input_size(B)    total_output_size(B)    type    time:
→avg(us)    algbw(GB/s)    busbw(GB/s)
  1024        256              32                    1024          fp32
→ 6.16        0.17            0.16
  8192       2048              256                    8192          fp32
→ 6.48        1.26            1.23
  65536      16384              2048                    65536          fp32
→ 8.17        8.02            7.77
  524288     131072              16384                    524288          fp32
→ 23.16       22.64           21.93
Avg bus bandwidth:    7.7715GB/s

```

- Example results as JSON with --report-to-json-file

```
nccom-test -r 32 --minbytes 1kb --maxbytes 1mb --stepfactor 8 --datatype fp32 --check_
```

```
→allg --report-to-json-file nccom-results.json
```

| size(B) | count(elems) | type | time:avg(us) | algbw(GB/s) | busbw(GB/s) |
|---------|--------------|------|--------------|-------------|-------------|
| 1024    | 256          | fp32 | 6.19         | 0.17        | 0.16        |
| 8192    | 2048         | fp32 | 6.55         | 1.25        | 1.21        |
| 65536   | 16384        | fp32 | 8.18         | 8.01        | 7.76        |
| 524288  | 131072       | fp32 | 23.11        | 22.69       | 21.98       |

```
Avg bus bandwidth: 7.7775GB/s
```

```
python3 -m json.tool nccom-results.json
```

```
{
  "results": [
    {
      "size(B)": 1024,
      "count(elems)": 256,
      "type": "fp32",
      "algbw(GB/s)": 0.16553675170497603,
      "busbw(GB/s)": 0.16036372821419553,
      "time:avg(us)": 6.19
    },
    {
      "size(B)": 8192,
      "count(elems)": 2048,
      "type": "fp32",
      "algbw(GB/s)": 1.2500906056270864,
      "busbw(GB/s)": 1.21102527420124,
      "time:avg(us)": 6.55
    },
    {
      "size(B)": 65536,
      "count(elems)": 16384,
      "type": "fp32",
      "algbw(GB/s)": 8.008982241741455,
      "busbw(GB/s)": 7.758701546687035,
      "time:avg(us)": 8.18
    },
    {
      "size(B)": 524288,
      "count(elems)": 131072,
      "type": "fp32",
      "algbw(GB/s)": 22.688776793562784,
      "busbw(GB/s)": 21.97975251876395,
      "time:avg(us)": 23.11
    }
  ]
}
```



## Multiple Instances Example

- 64 rank all-reduce on two instances for sizes ranging from 8 bytes to 1GiB with a step of 2x, running 50 ops

```
NEURON_RT_ROOT_COMM_ID=10.1.4.145:45654 nccom-test -r 64 -N 2 -b 8 -e 1GB -
```

| size(B)  | count(elems) | type | time(us) | algbw(GB/s) | busbw(GB/s) |
|----------|--------------|------|----------|-------------|-------------|
| 8        | 2            | fp32 | 520      | 0.00        | 0.00        |
| 16       | 4            | fp32 | 520      | 0.00        | 0.00        |
| 32       | 8            | fp32 | 523      | 0.00        | 0.00        |
| 64       | 16           | fp32 | 525      | 0.00        | 0.00        |
| 128      | 32           | fp32 | 553      | 0.00        | 0.00        |
| 256      | 64           | fp32 | 709      | 0.00        | 0.00        |
| 512      | 128          | fp32 | 782      | 0.00        | 0.00        |
| 1024     | 256          | fp32 | 840      | 0.00        | 0.00        |
| 2048     | 512          | fp32 | 881      | 0.00        | 0.00        |
| 4096     | 1024         | fp32 | 916      | 0.00        | 0.01        |
| 8192     | 2048         | fp32 | 1013     | 0.01        | 0.01        |
| 16384    | 4096         | fp32 | 1031     | 0.01        | 0.03        |
| 32768    | 8192         | fp32 | 1174     | 0.03        | 0.05        |
| 65536    | 16384        | fp32 | 1315     | 0.05        | 0.09        |
| 131072   | 32768        | fp32 | 1315     | 0.09        | 0.18        |
| 262144   | 65536        | fp32 | 1311     | 0.19        | 0.37        |
| 524288   | 131072       | fp32 | 1312     | 0.37        | 0.73        |
| 1048576  | 262144       | fp32 | 1328     | 0.74        | 1.45        |
| 2097152  | 524288       | fp32 | 1329     | 1.47        | 2.89        |
| 4194304  | 1048576      | fp32 | 1378     | 2.83        | 5.58        |
| 8388608  | 2097152      | fp32 | 1419     | 5.51        | 10.84       |
| 16777216 | 4194304      | fp32 | 2138     | 7.31        | 14.39       |
| 33554432 | 8388608      | fp32 | 2711     | 11.53       |             |

(continues on next page)

(continued from previous page)

|                    |            |            |      |       |       |   |
|--------------------|------------|------------|------|-------|-------|---|
| →22.69             | 67108864   | 16777216   | fp32 | 3963  | 15.77 | └ |
| →31.05             | 134217728  | 33554432   | fp32 | 6279  | 19.91 | └ |
| →39.19             | 268435456  | 67108864   | fp32 | 11954 | 20.91 | └ |
| →41.17             | 536870912  | 134217728  | fp32 | 21803 | 22.93 | └ |
| →45.15             | 1073741824 | 268435456  | fp32 | 41806 | 23.92 | └ |
| →47.09             |            |            |      |       |       |   |
| Avg bus bandwidth: |            | 9.3924GB/s |      |       |       |   |

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf1, Inf2, Trn1, Trn2

6.2.6 Neuron System Tools

Table of Contents

- *Neuron Tools [2.24.54.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.23.16.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.22.66.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.20.204.0]*
  - *New in the release*
  - *Neuron Profile 2.0 (Beta)*
- *Neuron Tools [2.19.0.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.18.3.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.17.1.0]*
  - *Bug fixes*

- *Neuron Tools [2.17.0.0]*
  - *New in the release*
  - *Bug fixes*
- *Neuron Tools [2.16.1.0]*
  - *New in the release*
  - *Bug fixes*
  - *Known issues*
- *Neuron Tools [2.15.4.0]*
- *Neuron Tools [2.14.6.0]*
- *Neuron Tools [2.13.4.0]*
- *Neuron Tools [2.12.2.0]*
- *Neuron Tools [2.11.10.0]*
- *Neuron Tools [2.10.1.0]*
- *Neuron Tools [2.9.5.0]*
- *Neuron Tools [2.8.2.0]*
- *Neuron Tools [2.7.2.0]*
- *Neuron Tools [2.6.0.0]*
- *Neuron Tools [2.5.19.0]*
- *Neuron Tools [2.5.16.0]*
- *Neuron Tools [2.4.6.0]*
- *Neuron Tools [2.1.4.0]*
- *Neuron Tools [2.0.790.0]*
- *Neuron Tools [2.0.623.0]*
- *Neuron Tools [2.0.494.0]*
- *Neuron Tools [2.0.327.0]*
- *Neuron Tools [2.0.277.0]*

## Neuron Tools [2.24.54.0]

Date: 6/24/2025

### New in the release

- Scratchpad memory usage visualization is now available in the Neuron Profiler UI. This feature shows scratchpad memory usage over time with tensor level visibility at each time slice. This view also provides the HLO name associate with the first use of each tensor. See [View Scratchpad Usage With Memory Tracker](#) for more details.
- Framework stack traces are now available in the Neuron Profiler UI. Users will now be able to see how instructions executed on device map to their source code. See [Framework Stack Trace](#) for more details.
- On-device collectives barriers are now shown in the Neuron Profiler UI. This annotation will make it more clear when there is overhead for collectives synchronization. See [Collectives](#) for an example.
- HBM throughput visualization over time is now shown in the Neuron Profiler UI, which reflects data movement where either the source or destination are HBM.
- Added option to filter the Neuron Cores to capture trace events on ([reference](#))
- Added option to filter the event types recorded when capturing system traces ([reference](#))
- Added a flag to `nccom-test` to get results in JSON (`--report-to-json-file <filename>`).
- Added a flag to `nccom-test` to explicitly show input and output sizes based on the operation (`--show-input-output-size`).

### Bug fixes

- Fixed instance id labeling in system profile view for framework events.
- Fixed issue in Neuron Profiler UI where the full data was not shown in the NEFF Nodes tab.

### Neuron Tools [2.23.16.0]

Date: 5/19/2025

### New in the release

- Improved Neuron Profiler performance, allowing users to view profile results 5x times faster on average.
- Improved error reporting with timeline support for error signatures via custom notifications in the Neuron Profiler UI. Added execution and out-of-bounds (OOB) error tracking in Neuron Profiler JSON outputs.
- Updated the default grouping for system profiles to include process ID.
- Added `neuron-monitor` companion script for collecting Kubernetes info in EKS. See [neuron-monitor-k8s-info.py \(Beta\)](#) for details.

## Bug fixes

- Fixed hang during data collection when running `nccom-test` across multiple instances.
- Fixed certain cases in Neuron Profiler where DMA sizes were always reported as 0 bytes.

## Neuron Tools [2.22.66.0]

Date: 2/14/2025

### New in the release

- `neuron-det` is no longer supported starting with this release. We recommend customers transition to Neuron Profiler 2.0 (Beta) for debugging runtime hangs and issues in large-scale settings. This tool offers the same runtime function level traces with improved ease of use and optimized performance. For more information about Neuron Profiler 2.0 (Beta), see [Neuron Profiler 2.0 \(Beta\) User Guide](#).
- Added several enhancements to the Neuron Profiler UI, including NeuronCore barrier annotations, a minimal default view to improve initial load performance, usability of updating markers, and better organization of view settings.
- Added new event types in the system profile for Neuron Profiler 2.0 (Beta) related to out-of-bounds execution errors, execution request submission, and model switch overhead.
- Updated system trace output format for Neuron Profiler 2.0 (Beta). Users will need to upgrade the `aws-neuronx-runtime-lib` and `aws-neuronx-tools` packages to the same Neuron SDK version to process and view the profiles.

## Bug fixes

- Fixed an issue in the Neuron Profiler UI where dependencies were misaligned in the timeline when highlighted.
- Fixed an issue where instruction dependency IDs were truncated in the Neuron Profiler JSON output.

## Neuron Tools [2.20.204.0]

Date: 12/20/2024

### New in the release

- Added support for Trn2 instance types.
- Added support for Logical Neuroncores. `neuron-top`, `neuron-monitor`, and `neuron-ls` now display and aggregate information per Logical Neuroncore based on LNC configuration.
- Added Neuron Profile 2.0 (Beta). See [Neuron Profiler 2.0 \(Beta\) User Guide](#) for more information.

## Neuron Profile 2.0 (Beta)

- System profiles featuring Neuron Runtime API trace and ML framework trace.
- Option to view system and device profiles using the Perfetto UI
- Support for native JAX and PyTorch profilers.
- Support for distributed workloads in environments such as EKS and ParallelCluster.
- Ability to drill down from high-level system profiles to low-level device profiles.
- Simplified experience for capturing profiles.

## Neuron Tools [2.19.0.0]

Date: 09/16/2024

### New in the release

- Added support for Neuron Kernel Interface (NKI). Please see [Profiling NKI kernels with Neuron Profile](#) for more info.
- Updated `neuron-profile` JSON output to include information regarding instruction dependencies, DMA throughput, and SRAM usage. See [Alternative output formats](#) on how to generate this output.
- Updated Neuron Profiler UI to display transpose information for DMAs (when applicable). Hover over the tooltip for further details (see [Features](#) on using tooltips).

### Bug fixes

- Fixed error handling in `neuron-top` to exit gracefully when passing an unknown argument

## Neuron Tools [2.18.3.0]

Date: 07/03/2024

### New in the release

- Profile captured with Neuron Runtime 2.20+ now includes annotations with additional information such as duration, size, and replica groups around collective operations.
- Running `neuron-profile capture` for workloads with collectives will now attempt to use the required number of workers if `-collectives-workers-per-node` or `-collectives-worker-count` is not set.
- Profiler UI now persists searched information in the URL and provides a summary of the search results.
- Updating sampling approach to show more representative data in the profiler UI when zoomed out.
- Updated groupings for displayed info on click in the profiler UI.
- Added `neuron_device_type` and `neuron_device_memory_size` to `neuron-monitor`'s hardware information output.

### Bug fixes

- Resolved issue where *NaN* would be seen in the JSON output of *neuron-profile* and result in parsing errors.
- Resolved inconsistent timeline display issues in profiler UI that depended on when the profile was processed.
- *neuron-profile view --output-format summary-text* will now display in a fixed order.
- Updated accuracy of pending DMA count in the profiler UI.
- Removed unnecessary calls to *exec* when capturing memory utilization metrics in *neuron-monitor*.

### Neuron Tools [2.17.1.0]

Date: 04/01/2024

### Bug fixes

- Fixed potential hang during synchronization step in *nccom-test*.

### Neuron Tools [2.17.0.0]

Date: 02/13/2024

### New in the release

- Added support to *neuron-profile* for collective communication operator improvements in Neuron SDK 2.17. See [Neuron Runtime Release Notes](#) for more info.
- Optimized count query for sampling in *neuron-profile* UI for up to 3x faster load performance.
- Introduced warning annotations in *neuron-profile* UI to automatically highlight potential performance issues. See the [Neuron Profile User Guide](#) for more info.

### Bug fixes

- Resolved issue of inaccurate execution time reported by *neuron-profile* as mentioned in Neuron Tools 2.16.1.0 release notes.
- Fixed NaN display errors in the *neuron-profile* UI.
- Fixed file naming issue when capturing collectives profiles with *neuron-profile*.

### Neuron Tools [2.16.1.0]

Date: 12/21/2023

#### New in the release

- First release of the Neuron Distributed Event Tracing tool `neuron-det` to visualize execution for multi-node workloads. Get started with the `neuron-det-ug`.
- `neuron-profile` now has the ability to capture multi-worker jobs. See the [Neuron Profile User Guide](#) for more info.
- Added terminology descriptions to `neuron-profile` summary statistics. To view through the CLI, use `neuron-profile view --terminology` To view in the UI, hover over the key in the summary.
- Added optional flags to `neuron-profile view` to change the InfluxDB bucket name (`--db-bucket <bucket name>`) and profile display name (`--display-name <name>`).

#### Bug fixes

- Fixed bug where GPSSimd summary values were missing in the profile summary.
- Fixed issue in `nccom-test` to no longer expect Neuron Device 0 in a container environemnt.
- Fixed issue in `nccom-test` to no longer require the instance launching `nccom-test` to be participating in the workload.

#### Known issues

- Execution time reported in `neuron-profile` is sometimes in-accurate due to a bug in how the time is captured. The bug will be address in upcoming Neuron releases.

### Neuron Tools [2.15.4.0]

Date: 10/26/2023

New in the release:

- Fixed bug in `neuron-profile` that may result in a crash when using the NeuronCore Pipeline feature on Inf1.
- Improved visibility of summary stats in the profiler UI with added groupings.
- Added support for alltoall CC operation in `nccom-test`.

### Neuron Tools [2.14.6.0]

Date: 09/15/2023

New in the release:

- Added legend in `neuron-ls` to clarify wrap around edges for topology view.
- Improved error messaging when passing invalid arguments to `neuron-profile view`.
- Fixed bug in `neuron-profile` that incorrectly calculated buffer utilization for more recently compiled NEFFs.



- Fixed bug in `neuron-profile` where the profile would sometimes include additional idle time while waiting for execution to start.
- Profiler output now includes HLO name in addition to framework layer names.
- `neuron-profile` view now has `--output-format json` option which will write to a file specified by `--output-file <name>` (default is `ntff.json`) instead of writing data to InfluxDB.

### Neuron Tools [2.13.4.0]

Date: 08/28/2023

New in the release:

- `--check` option of `nccom-test` now supports more data types (`fp16`, `bf16`, `(u)int8`, `(u)int16`, and `(u)int32` are now supported in addition to `fp32`)
- Fixed bug in `nccom-test` that would wait indefinitely for execution to end when running on multiple instances (`-N 2` and higher).
- Fixed bug in `neuron-profile` to prevent a crash during utilization calculation

### Neuron Tools [2.12.2.0]

Date: 7/19/2023

New in the release:

- Bumped the max supported profiling NTFF version to version 2 to resolve crashes when postprocessing NTFFs captured with newer versions of the Neuron Runtime Library. When viewing profiles captured using Neuron Runtime Library 2.15 or above, please upgrade tools to 2.12. This version of Neuron tools remains compatible with NTFF version 1.
- Bug fixes for `neuron-profile` related to the calculation of some summary stats.

### Neuron Tools [2.11.10.0]

Date: 6/14/2023

New in the release:

- `nccom-test` can now show multiple latency stats in the results table, such as average or percentiles, by specifying the `-s` option (for example: `-s p10 p99 avg p50`).
- First public support for `neuron-profile` as a standalone tool that can be used to profile executions on Neuron Devices. Visit the Neuron Tools documentation page for more details on how to use the Neuron Profiler.

### Neuron Tools [2.10.1.0]

Date: 05/01/2023

New in the release:

- Added new Neuron Collectives benchmarking tool, `nccom-test`, to enable benchmarking sweeps on various Neuron Collective Communication operations. See new `nccom-test` documentation under System Tools for more details.
- Expanded support for Neuron profiling to include runtime setup/teardown times and collapsed execution of NeuronCore engines and DMA. See Tensorboard release notes and tutorial for more details.

### Neuron Tools [2.9.5.0]

Date: 03/28/2023

New in the release:

- Updated neuron-top to show effective FLOPs across all NeuronCores.

### Neuron Tools [2.8.2.0]

Date: 02/24/2023

New in the release:

- Updated neuron-top to show aggregated utilization/FLOPs across all NeuronCores.

### Neuron Tools [2.7.2.0]

Date: 02/08/2023

New in the release:

- Added support for model FLOPS metrics in both neuron-monitor and neuron-top. More details can be found in the Neuron Tools documentation.

### Neuron Tools [2.6.0.0]

Date: 12/09/2022

This release adds support for profiling with the Neuron Plugin for TensorBoard on TRN1. Please check out the documentation *Neuron Plugin for TensorBoard (Trn1)*.

New in the release:

- Updated profile post-processing for workloads executed on TRN1

### Neuron Tools [2.5.19.0]

Date: 11/07/2022

New in the release:

- Minor bug fixes and improvements.

### Neuron Tools [2.5.16.0]

Date: 10/26/2022

New in the release:

- New neuron-monitor and neuron-top feature: **memory utilization breakdown**. This new feature provides more details on how memory is being currently used on the Neuron Devices as well as on the host instance.
- neuron-top's UI layout has been updated to accommodate the new **memory utilization breakdown** feature.
- neuron-monitor's inference\_stats metric group was renamed to execution\_stats. While the previous release still supported inference\_stats, starting this release the name inference\_stats is considered deprecated and can't be used anymore.

---

**Note:** For more details on the new **memory utilization breakdown** feature in `neuron-monitor` and `neuron-top` check out the full user guides: *Neuron Monitor User Guide* and *Neuron Top User Guide*.

---

Bug Fixes:

- Fix a rare crash in `neuron-top` when the instance is under heavy CPU load.
- Fix process names on the bottom tab bar of `neuron-top` sometimes disappearing for smaller terminal window sizes.

### Neuron Tools [2.4.6.0]

Date: 10/10/2022

This release adds support for both EC2 INF1 and TRN1 platforms. Name of the package changed from `aws-neuron-tools` to `aws-neuronx-tools`. Please remove the old package before installing the new one.

New in the release:

- Added support for ECC counters on Trn1
- Added version number output to `neuron-top`
- Expanded support for longer process tags in `neuron-monitor`.
- Removed hardware counters from the default `neuron-monitor` config to avoid sending repeated errors - will add back in future release.
- `neuron-ls` - Added option `neuron-ls --topology` with ASCII graphics output showing the connectivity between Neuron Devices on an instance. This feature aims to help in understanding pathways between Neuron Devices and in exploiting code or data locality.

Bug Fixes:

- Fix `neuron-monitor` and `neuron-top` to show the correct Neuron Device when running in a container where not all devices are present.

### Neuron Tools [2.1.4.0]

Date: 04/29/2022

- Minor updates

### Neuron Tools [2.0.790.0]

Date: 03/25/2022

- `neuron-monitor`: fixed a floating point error when calculating CPU utilization.

### Neuron Tools [2.0.623.0]

Date: 01/20/2022

New in the release:

- `neuron-top` - Added “all” tab that aggregates all aggregate all running Neuron processes into a single view.
- `neuron-top` - Improved startup time to approximately 1.5 seconds in most cases.
- `neuron-ls` - Removed header message about updating tools from `neuron-ls` output

Bug fixes:

- `neuron-top` - Reduced single CPU core usage down to 0.7% from 80% on `inf1.xlarge` when running `neuron-top` by switching to an event-driven approach for screen updates.

### Neuron Tools [2.0.494.0]

Date: 12/27/2021

- Security related updates related to log4j vulnerabilities.

### Neuron Tools [2.0.327.0]

Date: 11/05/2021

- Updated Neuron Runtime (which is integrated within this package) to `libnrt 2.2.18.0` to fix a container issue that was preventing the use of containers when `/dev/neuron0` was not present. See details here [neuron-runtime-release-notes](#).

### Neuron Tools [2.0.277.0]

Date: 10/27/2021

New in this release:

- Tools now support applications built with Neuron Runtime 2.x (`libnrt.so`).

---

#### Important:

- You must update to the latest Neuron Driver (`aws-neuron-dkms` version 2.1 or newer) for proper functionality of the new runtime library.
  - Read *[Introducing Neuron Runtime 2.x \(libnrt.so\)](#)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.
  - Read *[Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#)* for detailed information of how to migrate your application.
- 

- Updates have been made to `neuron-ls` and `neuron-top` to significantly improve the interface and utility of information provided.
- Expands `neuron-monitor` to include additional information when used to monitor latest Frameworks released with Neuron 1.16.0.

**`neuron_hardware_info`** Contains basic information about the Neuron hardware.

```
"neuron_hardware_info": {
  "neuron_device_count": 16,
  "neuroncore_per_device_count": 4,
  "error": ""
}
```

- neuron\_device\_count : number of available Neuron Devices
- neuroncore\_per\_device\_count : number of NeuronCores present on each Neuron Device
- error : will contain an error string if any occurred when getting this information (usually due to the Neuron Driver not being installed or not running).
- neuron-cli entering maintenance mode as it's use is no longer relevant when using ML Frameworks with an integrated Neuron Runtime (libnrt.so). see maintenance\_mxnet\_1\_5 for more information.
- For more information visit neuron-tools

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.3 Profiling Tools

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.3.1 Neuron Profile User Guide

#### Table of contents

- *Overview*
- *Installation*
  - *Ubuntu*
- *Capturing a profile*
- *Capturing profiles for multi-worker jobs*
  - *Capturing profiles for multi-node jobs*
- *Processing and viewing the profile results*
  - *Viewing a single profile*
  - *Viewing profiles for multi-worker jobs*
  - *Viewing multiple profiles*
  - *Accessing the profiles*
  - *Alternative output formats*
- *Understanding a Neuron profile*
  - *Timeline*

- *Features*
  - *Performance Warnings*
  - *Collectives*
  - *Event Details*
- *Framework Stack Trace*
  - *Searching Profiles*
  - *Hardware Errors*
- *View Scratchpad Usage With Memory Tracker*
- *Viewing Profiles with Perfetto*
  - *Viewing Large Profiles In Perfetto*
  - *Showing Dependencies In Perfetto*
- *CLI reference*
- *FAQ*
  - *Difference between TensorE and TensorMatrixE Rows in Timeline*
- *Troubleshooting*
  - *InfluxDB not installed*
  - *Too many open files*
  - *When viewing UI “FATAL - Failed metadata query”*
  - *Visual Artifacts when viewing profiles*

## Overview

`neuron-profile` is a tool to profile and analyze performance of a ML model compiled with the Neuron compiler and run on NeuronDevices.

---

**Note:** Please use the `aws-neuronx-tools` package from Neuron SDK 2.11 or higher.

---

`neuron-profile` helps developers identify performance bottlenecks and optimize their workloads for NeuronDevices. `neuron-profile` provides insights into NeuronDevice activity including the instructions executed on each compute engine (ex. Tensor engine, Vector engine, etc.), DMA data movement activity, and performance metrics such as engine utilization, DMA throughput, memory usage, and more. NeuronDevice activity is collected by the `neuron-profile capture` command which runs the model with tracing enabled. Profiling typically has near zero overhead because NeuronDevices have dedicated on-chip hardware profiling.

Additionally, `neuron-profile` supports Neuron Kernel Interface (NKI) developers in profiling their kernels. For more information, please refer to [Profiling NKI kernels with Neuron Profile](#)

---

**Note:** This page refers to the existing Neuron Profiler feature set focused on capturing and viewing device profiles (hardware activity during graph execution on NeuronCores). Neuron Profiler 2.0 is a set of new features currently in beta that enhance and simplify capturing and viewing profiles. It is not a replacement for the features described on this page. To learn more about Neuron Profiler 2.0, please refer to the [Neuron Profiler 2.0 \(Beta\) User Guide](#).

---

## Installation

`neuron-profile` comes as part of the `aws-neuronx-tools` package, and will be installed to `/opt/aws/neuron/bin`.

The Neuron web profile viewer utilizes InfluxDB OSS 2.x to store time series data for the profiled workloads after post processing. Please follow the instructions provided at <https://portal.influxdata.com/downloads/> for the correct OS. A sample installation of Neuron Profile and InfluxDB is provided below.

## Ubuntu

```
# Install Neuron Profile
. /etc/os-release
sudo tee /etc/apt/sources.list.d/neuron.list > /dev/null <<EOF
deb https://apt.repos.neuron.amazonaws.com ${VERSION_CODENAME} main
EOF

wget -qO - https://apt.repos.neuron.amazonaws.com/GPG-PUB-KEY-AMAZON-AWS-NEURON.PUB | \
  sudo apt-key add -
sudo apt-get update -y
sudo apt-get install aws-neuronx-runtime-lib aws-neuronx-dkms -y
sudo apt-get install aws-neuronx-tools -y

# Install InfluxDB
wget -q https://repos.influxdata.com/influxdata-archive_compat.key
echo '393e8779c89ac8d958f81f942f9ad7fb82a25e133faddaf92e15b16e6ac9ce4c influxdata-
  archive_compat.key' | sha256sum -c && cat influxdata-archive_compat.key | gpg --
  dearmor | sudo tee /etc/apt/trusted.gpg.d/influxdata-archive_compat.gpg > /dev/null
echo 'deb [signed-by=/etc/apt/trusted.gpg.d/influxdata-archive_compat.gpg] https://repos.
  influxdata.com/debian stable main' | sudo tee /etc/apt/sources.list.d/influxdata.list

sudo apt-get update && sudo apt-get install influxdb2 influxdb2-cli -y
sudo systemctl start influxdb
influx setup
# Fill in the information to finish the setup
```

## Capturing a profile

The `neuron-profile` tool can both capture and post-process profiling information. `neuron-profile` takes a compiled model (a NEFF), executes it, and saves the profile results to a NTFF (`profile.ntff` by default). For this example, we assume a NEFF is already available as `file.neff`

```
$ neuron-profile capture -n file.neff -s profile.ntff
```

## Capturing profiles for multi-worker jobs

`neuron-profile` can capture profiles for collectives-enabled NEFFs running across multiple NeuronCores, NeuronDevices, or even nodes. This is useful for understanding performance and communication overheads when deploying larger distributed models.

The following example, performs a distributed run across all NeuronDevices and NeuronCores on an `inf2.24xlarge` instances, capturing profiles for all 12 workers (one for each NeuronCore).

```
$ neuron-profile capture -n file.neff --collectives-workers-per-node 12 -s output/  
↪profile.ntff
```

A profile is saved for each worker in the output directory.

```
$ ls output  
profile_rank_0.ntff  profile_rank_2.ntff  profile_rank_6.ntff  profile_rank_1.ntff  ↪  
↪profile_rank_3.ntff  profile_rank_7.ntff  
profile_rank_10.ntff  profile_rank_4.ntff  profile_rank_8.ntff  profile_rank_11.ntff ↪  
↪profile_rank_5.ntff  profile_rank_9.ntff
```

It is also possible to run a distributed job while only capturing a profile for a specific worker instead of all workers. To do that, use the `--collectives-profile-id` option.

```
$ neuron-profile capture -n file.neff --collectives-profile-id 5 --collectives-workers-  
↪per-node 12 -s output/profile.ntff  
$ ls output  
profile_rank_5.ntff
```

## Capturing profiles for multi-node jobs

For multi-node jobs, `neuron-profile` must be invoked on each node using the `collectives-worker-start-id` to specify the global index of the first worker on the given node. For example, for a two node job with a total of four workers and two workers per node, the following commands are run on each node.

```
# on node 0  
$ neuron-profile capture -n file.neff --collectives-worker-start-id 0 --collectives-  
↪workers-per-node 2 --collectives-worker-count 4  
# on node 1  
$ neuron-profile capture -n file.neff --collectives-worker-start-id 2 --collectives-  
↪workers-per-node 2 --collectives-worker-count 4
```

`neuron-profile` saves the profile for a worker on the node where that worker was launched. So in the case above, `profile_rank_0.ntff` and `profile_rank_1.ntff` are saved to node 0, and `profile_rank_2.ntff` and `profile_rank_3.ntff` are saved to node 1.



## Processing and viewing the profile results

To analyze and view the collected profiling data, use the `view` subcommand of `neuron-profile`. This command performs two main functions: it post-processes the profiling data and starts up an HTTP server. Once the server is running, you can access the profiling results through your web browser. Please note: Chrome is the officially supported browser for viewing profiling results

**Note:** Profiles can be processed and viewed on another machine without Neuron devices. The `aws-neuronx-tools` package needs to be installed so that you can run `neuron-profile view`. To process the profile on another instance, you need to copy the NEFF and NTFF files from your Inf or Trn instance to that instance.

### Viewing a single profile

The first way to invoke `neuron-profile view` is to pass both the NEFF and the NTFF to this command. It will post-process these artifacts and print out a direct link to the profile view.

```
$ neuron-profile view -n file.neff -s profile.ntff
View profile at http://localhost:3001/profile/n_fdc71a0b582ee3009711a96e59958af921243921
ctrl-c to exit
```

### Viewing profiles for multi-worker jobs

Profiles from multi-worker jobs (i.e. more than one NeuronCore) can either be viewed individually or in a combined collectives view. Since profile data is often similar between workers and processing profile data for all workers can be time-consuming, it is recommended to first explore the profile for a single worker or small subset of workers. Viewing the profile for a specific worker is the same as for single-worker profiles.

```
$ neuron-profile view -n file.neff -s output/profile_rank_5.ntff
View profile at http://localhost:3001/profile/n_fdc71a0b582ee3009711a96e59958af921243921
```

To view the profile for multiple workers, pass the directory containing all worker profiles to `neuron-profile`.

```
$ neuron-profile view -n file.neff -d output
View profile at http://localhost:3001/profile_cc/p_
→9a69d907e1350100c9b03745eaa67aa7422842ed
```

file\_r12



When viewing profiles with the combined collectives view you can easily switch between the timelines of different workers by clicking the “Rank <x>” tabs.

Note: the “CC Aggregated View” currently shows no data. This will be populated in an upcoming release.

## Viewing multiple profiles

Alternatively, when post-processing multiple profiles, it may be desirable to have a persistent server running while processing results in the background. In this case, we can skip passing arguments to the command, which will direct users to the main page listing all available profiles.

```
$ neuron-profile view
View a list of profiles at http://localhost:3001/
```

In a separate window, we can kick off the post-processing without launching another server by passing the `--ingest-only` flag.

```
$ neuron-profile view -n file.neff -s profile.ntff --ingest-only
Profile "n_47cf9972d42798d236caa68952d0d29a76d8bd66" is ready to view
```

`n_47cf9972d42798d236caa68952d0d29a76d8bd66` is the bucket where the data is stored. We can find this profile at `localhost:3001/profile/<bucket>`.

## Accessing the profiles

If `neuron-profile view` is run on a remote instance, you may need to use port forwarding to access the profiles.

From the local machine, SSH to the remote instance and forward ports 3001 (the default `neuron-profile` HTTP server port) and 8086 (the default InfluxDB port). Then in the browser, go to `localhost:3001` to view the profiles.

```
$ ssh <user>@<ip> -L 3001:localhost:3001 -L 8086:localhost:8086
```

## Alternative output formats

Besides the web view mentioned above, `neuron-profile` also supports other output formats such as `summary-text` and `summary-json` for viewing overall metrics of the profile, as well as `json` for a parsable alternative.

## Profile summary

You can see a summary of each profile using the command `neuron-profile view --output-format summary-text -n file.neff -s output/profile_rank_<i>.</i>ntff`. This output includes summary metrics and fields for the NeuronCore (`nc_idx`) and NeuronDevice (`nd_idx`) on which the worker was run. For example, the following shows worker 5 used core 1 on device 3 and took 0.017 seconds (17 ms) to run the model.

```
$ neuron-profile view --output-format summary-text -n file.neff -s output/profile_rank_5.
↪ntff | grep -e "nd_idx" -e "nc_idx" -e "total_time"
nc_idx      1
nd_idx      2
total_time  0.017
```

This summary is also available as JSON using `--output-format summary-json`.

## JSON

You can also view the profile summary and all post-processed profiler events together as a single JSON. To do that, use the `--output-format json` option.

```
$ neuron-profile view --output-format json --output-file profile.json -n file.neff -s_
↪output/profile_rank_5.ntff
$ cat profile.json
{
  "summary": [
    {
      "total_time": 0.017,
      "event_count": 11215
      [...]
    }
  ],
  "instruction": [
    {
      "timestamp": 10261883214,
      "duration": 148,
      "label": "TensorMatrix",
```

(continues on next page)

(continued from previous page)

```

        "hlo_name": "%add.1 = add(%dot, %custom-call.44)",
        "opcode": "MATMUL",
        "operands": "S[5] (Tensor)++@complete acc_flags=3 row_grp=q0_
→src=fp16@0x5600[1,0,0][3,1,1] dst=0x20000000[1,0,0][3,1,1] 3*128 "
    },
    [...]
]
}

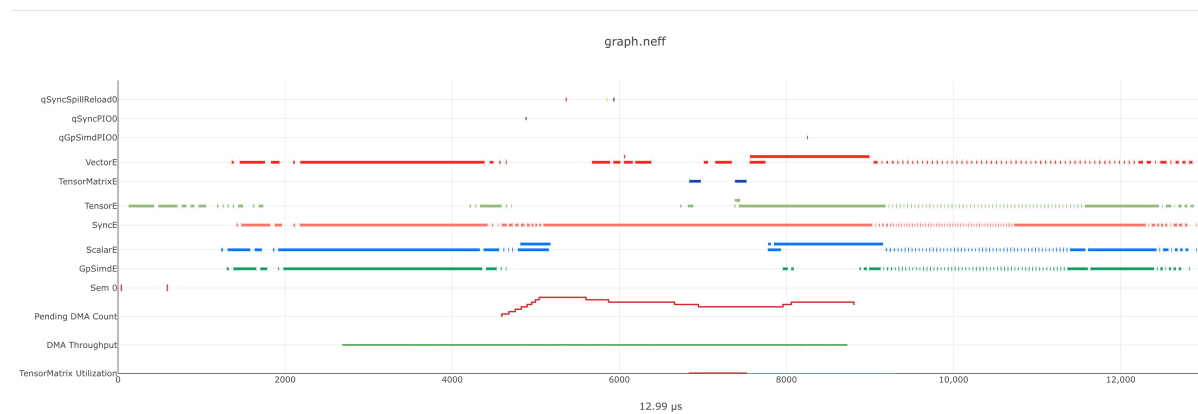
```

## Understanding a Neuron profile

The section provides a quick overview on what features and information are available through the Neuron web profile viewer.

For more information on terms used, please check out the [Neuron Glossary](#).

## Timeline



The execution timeline is plotted based on the elapsed nanoseconds since the start of execution.

Starting from the bottom, the **TensorMatrix Utilization** shows the efficiency of the **TensorEngine**, and the **Pending DMA Count** and **DMA Throughput** rows show the DMA activity. In general, we want these to be as high as possible, and in some cases may help give clues as to whether the workload is memory or compute bound.

Next are the individual **NeuronCore** engine executions. These rows show the start and end times for instructions executed by each engine, and clicking on one of these bars will show more detailed information, as well as any dependencies that were found. For models involving collective compute operations, you will additionally see rows labeled with **CC-core**, which are used to synchronize the CC operations.

Towards the top is the DMA activity. These can include the transfers of input and output tensors, intermediate tensors, and any additional spilling or loading to and from the on-chip SRAM memory.

## Features

The following are some useful features that may help with navigating a profile:

- Dragging your cursor across a portion of the timeline will zoom in to the selected window, providing a more in depth view of the execution during that time period.
- Hovering over a point will reveal a subset of information associated with it.
- Clicking a point will open a text box below the timeline with all the information associated with it.
- Right-clicking a point will drop a marker at a certain location. This marker will persist when zooming in and out.
  - All marker information can be found by clicking the **Annotations** button.
  - Markers can be saved and loaded by using a provided name for the marker set.
  - Individual markers can be renamed or deleted in this menu as well.
  - Time span between markers will automatically be shown, and users can change the marker name next to `diff` vs to calculate time between other markers.

Marker Set Name:

\* Note: Loading will replace the existing markers in the profiler and render the ones saved in the set.

|          |          |                 |          |
|----------|----------|-----------------|----------|
| "first"  | 0.068 ms | diff vs "first" | 0.0 ms   |
| "second" | 0.11 ms  | diff vs "first" | 0.040 ms |

- The “Search” tab can be used to find and highlight specific points in the profile related to the queried field(s).
- Click on the “Box Select” button in the top-right corner of the timeline and then click and drag on any region of the plot to select all events in that region and get summary statistics such as total duration and breakdowns of opcodes, transfer\_sizes, and more.
- The **Edit view settings** can be used to further customize the timeline view. Editing any settings will update the URL accordingly, which can be used to re-visit the current view at a later time.
  - For example, changing the **Instruction Grouping** dropdown option to “Layer” will re-color the timeline based on the associated framework layer name.
  - To speed up initial load times, the default will be a **Minimal View** which only shows the instructions executed and the model FLOPs utilization (MFU) over time. Changing between the minimal and full views can also be done through the **Reset to Full View** or **Reset to Minimal View** buttons.

Summary **View Settings** Layer Summary Selection Summary NEFF Header NEFF Nodes Model Info DMA Queues Info NC Memory Usage Info Terminology Help

**DMA color group**  
 DMAs will be colored by the selected grouping

**Instructions color group**  
 Instructions will be colored by the selected grouping

**Layer group depth**  
 current depth: 0 (use 0 for full name)

**Semaphore id**  
 Select the semaphore to be shown in the timeline

☒ Zoom based on timestamp ☐ Zoom based on markers

**Start timestamp**  
 Select the start timestamp for the plot

**End timestamp**  
 Select the start timestamp for the plot

Additionally, there are various summary tabs that can be clicked to provide more information on the model/NEFFs.

- **Layer Summary** shows timing information, FLOPs and instructions counts per layer.
- **Selection Summary** shows summarized information for all data points in the selected window when using the “Box Select” mode.
- **NEFF Header** shows details on the profiled NEFF, such as the number of NeuronCores required to execute.
- **NEFF Nodes** shows input, output, and weight tensor information, including name, size, and shape.
- **Model Info** shows a summary of the NTFF, such as the NeuronCore the model was executed on, number of notifications, and hardware execution time.
- **DMA Queues Info** shows more information on the queues used for data movement.
- **NC Memory Usage Info** shows a snapshot of the device memory usage breakdown before profiling was started.
- **Terminology** shows a description of metrics provided in the summary table.

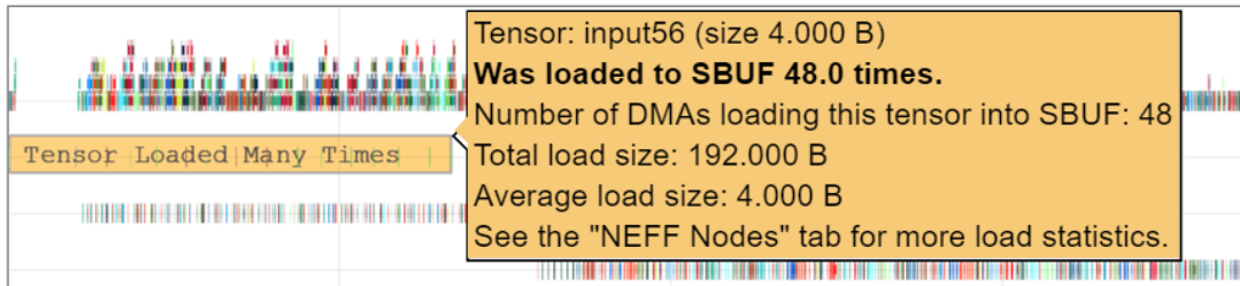
---

|                         |                               |                               |                                   |                             |                            |                            |                                 |                                      |                             |                      |
|-------------------------|-------------------------------|-------------------------------|-----------------------------------|-----------------------------|----------------------------|----------------------------|---------------------------------|--------------------------------------|-----------------------------|----------------------|
| <a href="#">Summary</a> | <a href="#">View Settings</a> | <a href="#">Layer Summary</a> | <a href="#">Selection Summary</a> | <a href="#">NEFF Header</a> | <a href="#">NEFF Nodes</a> | <a href="#">Model Info</a> | <a href="#">DMA Queues Info</a> | <a href="#">NC Memory Usage Info</a> | <a href="#">Terminology</a> | <a href="#">Help</a> |
|-------------------------|-------------------------------|-------------------------------|-----------------------------------|-----------------------------|----------------------------|----------------------------|---------------------------------|--------------------------------------|-----------------------------|----------------------|

---

## Performance Warnings

Furthermore, `neuron-profile` will automatically highlight some potential performance issues with warning annotations. For example if a tensor has been loaded more than 2 times a warning annotation (seen below as an orange box) will be drawn, encircling the dma instructions where the tensor was loaded many times. Hover on the annotation to see more details about loading the tensor. Another kind of warning annotation will highlight areas of high throttling. This provides the user a potential reason for slow down (thermal protection). Specific throttling details are shown when hovering the annotation.



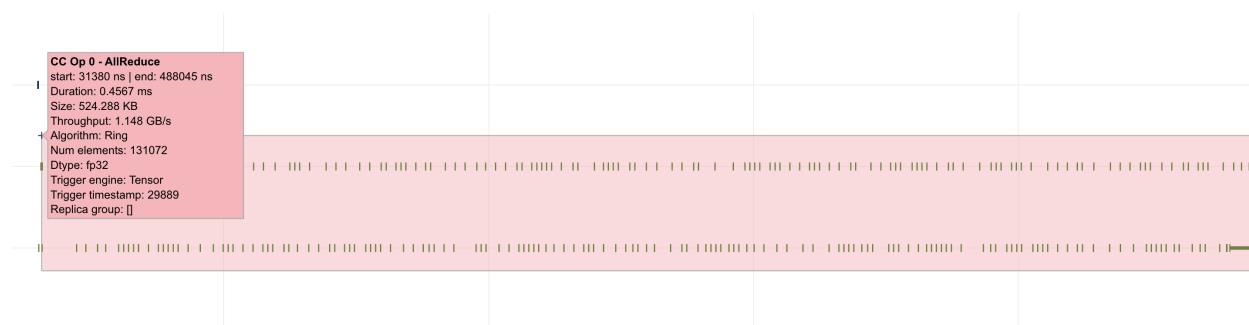
## Collectives

For models involving collective operations, the timeline will show a box around all data points related to each operation. Hovering the top left of the box will reveal more information associated with the operation.

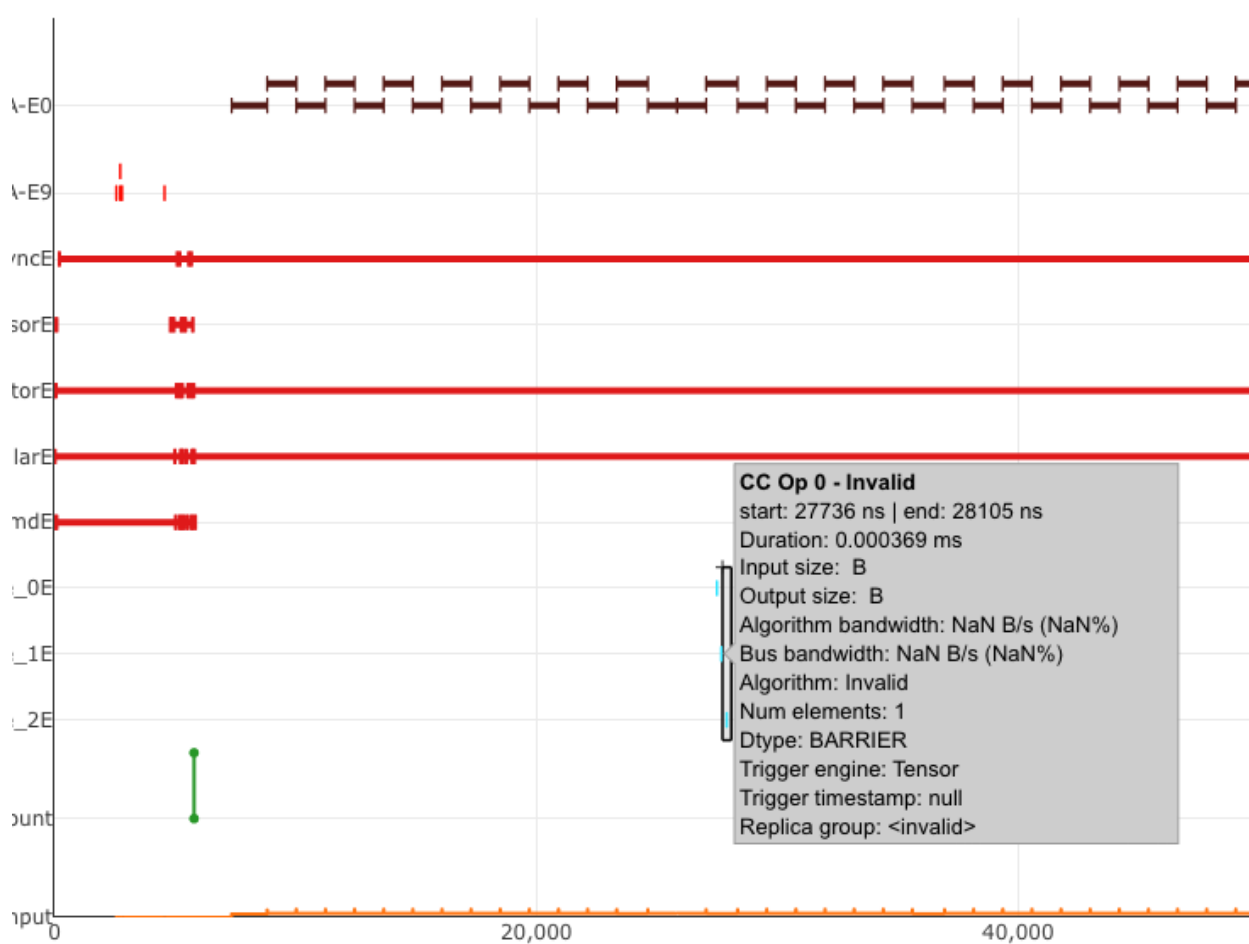
---

**Note:** this feature requires profiles to be captured with Neuron Runtime 2.20 or higher.

---



Additionally, for any on-device collectives synchronization barrier, a similar box will be display indicating a barrier instead of an actual collectives operation.



## Event Details

The information when a point is clicked is grouped by categories such as *Timing* or *IDs* for convenience. Each row will also include a tool tip on the right side, which can be hovered for an explanation on what the field represents. For instruction *Operands* specifically, clicking on the tooltip will reveal a breakdown of fields that compose an operand, as well as a generic example for reference. The examples may not apply directly to the currently viewed profile.

Event Details - TensorMatrixE

### Timing:

- Start: 7796431 ns ⓘ
- End: 7796623 ns ⓘ
- Duration: 192 ns ⓘ
- Evt wait time: 25 ns ⓘ

### Operation:

- Opcode: MATMUL ⓘ
- Operands: S[5] (Tensor)++@complete acc\_flags=3 src=bfloat16@0x2ff00[1,0,0][128,1,1] dst=0x2001800[1,0,0][128,1,1] 128\*128 debugHint=1 ⓘ
- Compiler opcode: MATMUL ⓘ
- Compiler operands: S[5] (Tensor)++@complete acc\_flags=3 src=bfloat16@0x2ff00[1,0,0][128,1,1] dst=0x2001800[1,0,0][128,1,1] 128\*128 ⓘ
- Instruction type: TRANSPOSE ⓘ
- Instruction dependency pcs:
  - GpSimd: 857
  - Scalar: 955 ⓘ

The HLO instruction that produced this instruction.

### HLO:

- Hlo name: %add.150 = add(%x.148, %y.149) ⓘ
- Hlo attrs: [{"op\_type":"","source\_file":"","source\_line":"0"}] ⓘ

### IDs:

- Pc: 21011 ⓘ
- Compiler pc: 20859 ⓘ
- Penguin id: sg0001:846 ⓘ
- Raw penguin id: 281474976711502 ⓘ
- Bir instruction name: l-846-3040\_sg0001 ⓘ
- Bir id: sg0000:67956 ⓘ
- Raw bir id: 67956 ⓘ
- Crc in: 51516 ⓘ
- Crc out: 65535 ⓘ

### Details:

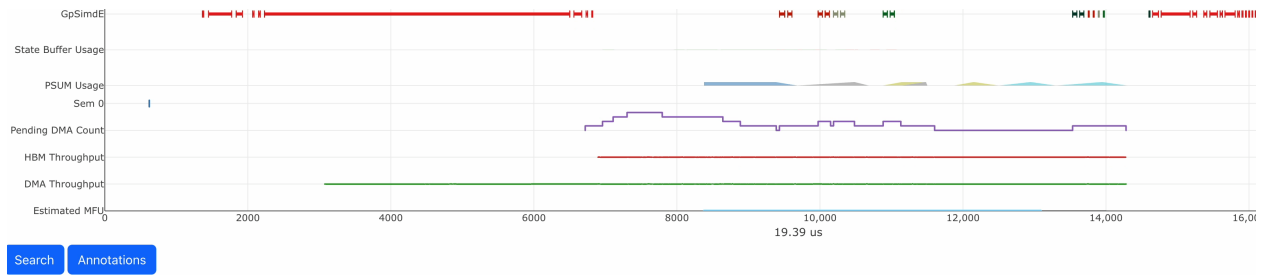
- Elements: 2097152

Close

## Framework Stack Trace

The Framework Stack Trace feature shows up in the Event Details when an instruction on the device profile is clicked. This can be used to map the device instructions back to framework level code in JAX or PyTorch to better understand what part of the application code resulted in a particular device instruction.





## Event Details - SyncE

## Timing:

- Start: 6717 ns
- End: 6766 ns
- Duration: 49 ns
- Evt wait time: 0 ns
- Trigger: 6000000000

## Operation:

- Opcode: WRITE
- Operands: data=0x1,0x1 size=4,4 dst=0xffff1022038,0xffff1002038
- Compiler opcode: PSEUDO\_DMA\_TRIGGER
- Compiler operands: qSyncIO0 block\_id=0

## HLO:

- Layer: aten\_\_addmm\_dot
- Hlo name: %dot = dot(%reshape, %transpose.3)
- Hlo attrs: [{"op\_type": "aten\_\_addmm", "source\_file": "/home/ubuntu/memory\_bound\_workload/torch\_workload/lib/python3.10/site-packages/torch/nn/modules/linear.py", "source\_line": "125"}]

## Framework Stack Trace:

- Stack Trace: [Click to see full stack trace](#)

Find the link to view the complete stack trace in the event details when you click on the instruction in the device profile. This will open up the stack trace viewer

## IDs:

- Id: 16328761139897035000
- Fully qualified subgraph: /sg00
- Pc: 26
- Compiler pc: 1

To enable tracking of the stack trace information, you need to set environment variables before compiling your NEFF:

```
export XLA_IR_DEBUG=1
export XLA_HLO_DEBUG=1
```

Once you have the NEFF, you can simply capture the profile as usual. While viewing the profile use the `--framework-source-root` to pass the path to framework source files. This is optional and is only needed if you want to view your code along side the profile.

```
$ neuron-profile view -n file.neff -s profile.ntff --framework-source-root /path/to/
↪ framework/source/files
```

## ▼ Complete Framework Stack Trace

[/home/ubuntu/memory\\_bound\\_workload/torch\\_workload/lib/python3.10/site-packages/torch/nn/modules/linear.py:125 forward](#)  
[/home/ubuntu/memory\\_bound\\_workload/torch\\_workload/lib/python3.10/site-packages/torch/nn/modules/module.py:1747 \\_call\\_impl](#)  
[/home/ubuntu/memory\\_bound\\_workload/torch\\_workload/lib/python3.10/site-packages/torch/nn/modules/module.py:1736 \\_wrapped\\_call\\_impl](#)  
[/home/ubuntu/memory\\_bound\\_workload/debug\\_info\\_test.py:27 forward](#)  
[/home/ubuntu/memory\\_bound\\_workload/torch\\_workload/lib/python3.10/site-packages/torch/nn/modules/module.py:1747 \\_call\\_impl](#)  
[/home/ubuntu/memory\\_bound\\_workload/torch\\_workload/lib/python3.10/site-packages/torch/nn/modules/module.py:1736 \\_wrapped\\_call\\_impl](#)  
[/home/ubuntu/memory\\_bound\\_workload/debug\\_info\\_test.py:50 main](#)  
[/home/ubuntu/memory\\_bound\\_workload/debug\\_info\\_test.py:57 <module>](#)

home/ubuntu/memory\_bound\_workload/debug\_info\_test.py

```

23         self.fc5 = nn.Linear(layers[5], output_size)
24     def forward(self, x):
25         x = F.relu(self.fc1(x))
26         x = F.relu(self.fc2(x))
27         x = self.fc3(x)
28         return F.log_softmax(x, dim=1)
29
30
31 def main():
32     # Fix the random number generator seeds for reproducibility
33     torch.manual_seed(0)
34     # XLA: Specify XLA device (defaults to a NeuronCore on Trn1 instance)
35     device = xm.xla_device()
36     # Start the profiler context-manager
37
38     # IMPORTANT: the model has to be transferred to XLA within
39     # the context manager, otherwise profiling won't work
40     model = MLP().to(device)
41
42     EPOCHS = 3
43     # start training loop
44     for epoch in range(EPOCHS):
45         train_x = torch.randn(1,10).to(device)
46
47         # Forward
48         output = model(train_x)
49
50     # XLA: collect ops and run them in XLA runtime
51     xm.mark_step()
52
53
54
55

```

[<](#) [Summary](#) | [View Settings](#) | [Layer Summary](#) | [Selection Summary](#) | [NEFF Header](#) | [NEFF Nodes](#) | [Model Info](#) | [DMA Queues Info](#) | [NC Memory Usage Info](#) | [Scratchpad Memory Usage](#) | [Com](#)

## Searching Profiles

Searching helps identify specific data points that may be worth investigating, such as all instructions related to a specific layer or operation. In the “Search” tab, select the corresponding field of interest and enter the value to search for. Multiple fields can be searched together. Please refer to the tooltip within the tab for more help on the query syntax. The search results will also include a summary of all data points found within the current time range.

Found 1448 results.
< Prev
Showing 1 / 1448
Next >
Highlight All

Category: instruction ▼

+ Add field

Search

Note: fields are joined by "AND"

Tip: enclose string with forward slashes to use regex, eg. /value/

X

Field: layer ▼

Search Result Summary

Duration: 0.69 **ms**

Start Time: 7.65 **ms**

End Time: 8.33 **ms**

Event Count: 1448

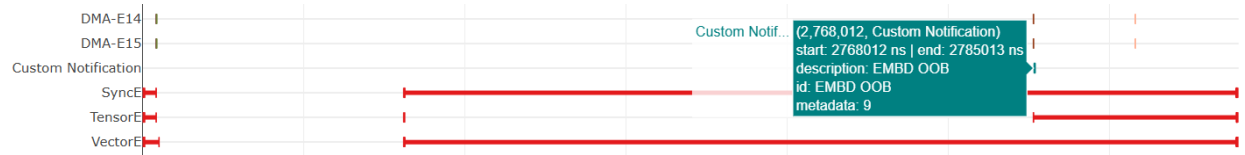
Event Duration Sum: 341.17 **us**

Event Duration Active: 315.28 **us** (45.84 %)

|  |   |
|--|---|
| GpSimdE (count = 9, active dur = 0.058%)         | ▼ |
| TensorE (count = 663, active dur = 7.88%)        | ▼ |
| TensorMatrixE (count = 668, active dur = 31.34%) | ▼ |
| VectorE (count = 108, active dur = 10.32%)       | ▼ |

## Hardware Errors

Invalid code can lead to errors on Neuron hardware. These errors will be displayed in Neuron Profile's Custom Notification timeline, as shown below. For example an Out of Bounds (OOB) error is displayed as:



Users can correlate the error to the time it occurred and view nearby events to help debug.

## View Scratchpad Usage With Memory Tracker

The Memory Tracker feature in Neuron Profiler provides detailed insights into scratchpad memory usage over time, showing how memory is allocated and utilized by different tensors during model execution. This is particularly useful for understanding memory bottlenecks and optimizing memory usage patterns.

To enable Memory Tracker, you need to set environment variables before compiling your NEFF:

```
export XLA_IR_DEBUG=1
export XLA_HLO_DEBUG=1
```

Then compile your model with these debug flags enabled. After compilation, capture the profile with the `--enable-dge-notifs` flag or set `NEURON_RT_ENABLE_DGE_NOTIFICATIONS=1`:

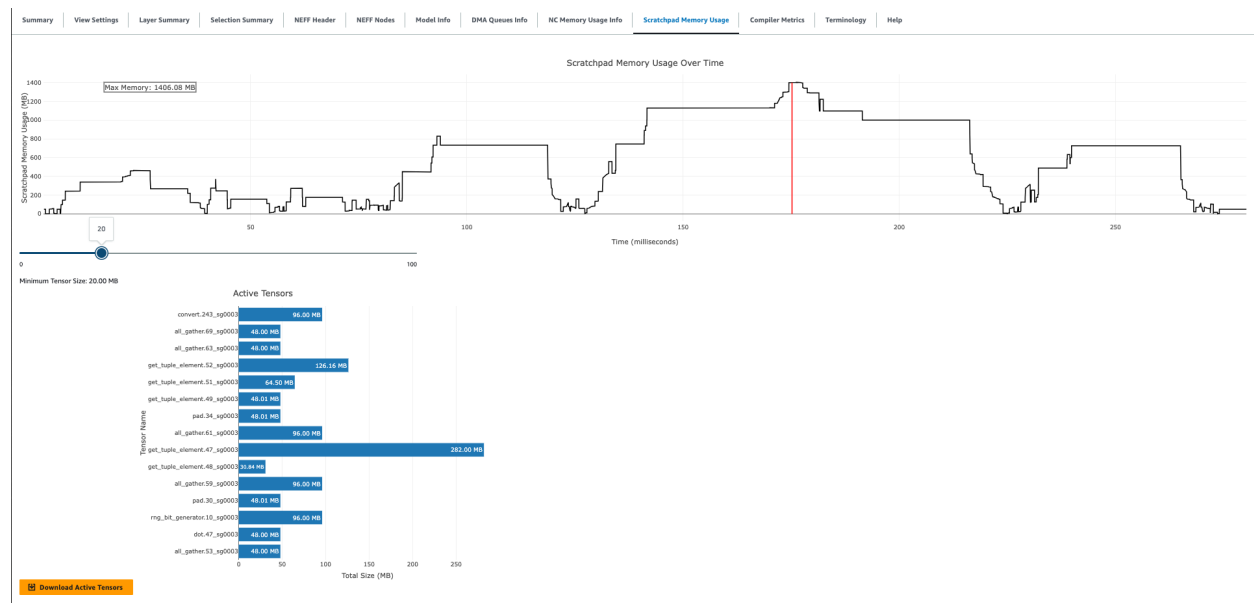
```
$ neuron-profile capture -n file.neff --enable-dge-notifs
```

Finally, view the profile with Memory Tracker enabled:

```
$ neuron-profile view -n file.neff -s profile.ntff --enable-memory-tracker
```

The Memory Tracker displays a timeline showing scratchpad memory usage over time, with a detailed breakdown of which tensors are consuming memory at any given point. This visualization helps identify:

- Peak scratchpad memory usage
- Memory allocation patterns
- Tensor-specific memory consumption
- Potential memory optimization opportunities



You can interact with the Memory Tracker timeline similar to other profile views - clicking on memory usage bars will show detailed information about the tensors using memory at that time, and you can zoom in to specific time ranges to get a more detailed view of memory allocation patterns.

## Viewing Profiles with Perfetto

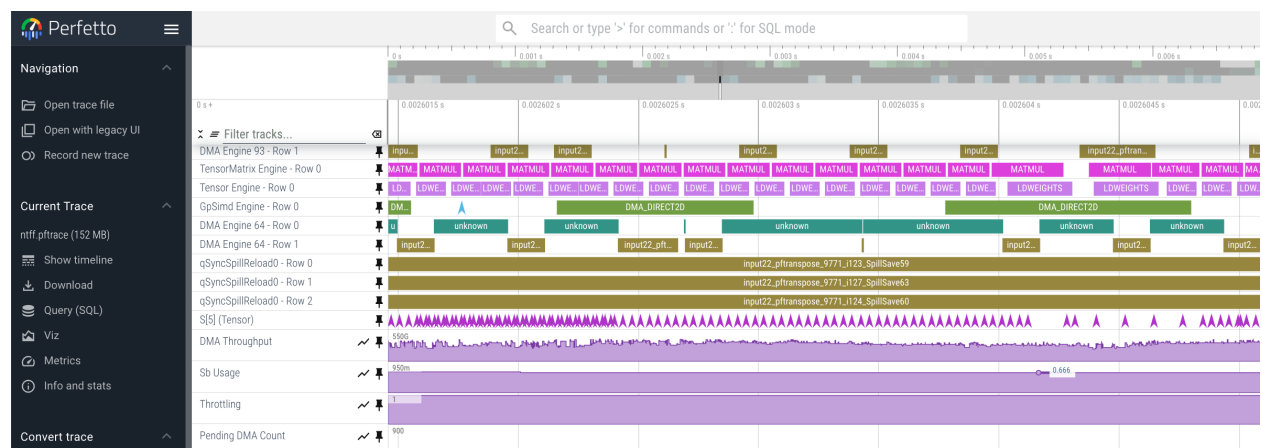
Perfetto is an open-source trace analysis toolkit with a powerful UI for visualizing and analyzing trace data. Users of Neuron Profiler have the option of viewing their profiles in the Perfetto UI.

To process your profile and generate a Perfetto trace file that can be viewed in the Perfetto UI run the following command:

```
$ neuron-profile view -n file.neff -s profile.ntff --output-format perfetto
```

This will generate a ntff.pftrace file. Go to <https://ui.perfetto.dev/> in your browser and open the ntff.pftrace file to view your profile in Perfetto.

**Note:** When loading trace files in the Perfetto UI, your data is processed locally and not uploaded to Perfetto's servers.



## Viewing Large Profiles In Perfetto

Your browser may run out of memory when viewing `ntff.pftrace` (Perfetto trace) files that are more than a few hundred MB. To get around this problem you can use the trace processor script by running the following command on your local system where you wish to view the profile

```
curl -LO https://get.perfetto.dev/trace_processor
chmod +x ./trace_processor
./trace_processor --httpd ntff.pftrace
```

Now go to <https://ui.perfetto.dev/> in your browser and in the dialog box that pops up click the “YES, use loaded trace” button.

For more information on using the trace processor script and viewing large traces, please refer to the Perfetto documentation at <https://perfetto.dev/docs/visualization/large-traces>.

## Showing Dependencies In Perfetto

By default Neuron Profiler does not process dependencies for profiles to be viewed in Perfetto because Perfetto renders the full dependency chain which can be visually overwhelming. To include dependencies that can be viewed when clicking instructions and DMAs in the Perfetto UI, use the `--show-perfetto-flows` flag when processing your profile.

```
$ neuron-profile view -n file.neff -s profile.ntff --output-format perfetto --show-
  ↳perfetto-flows
```

## CLI reference

### neuron-profile capture

**neuron-profile capture** [parameters] [inputs...]

Takes a given compiled NEFF, executes it, and collects the profile results. When no inputs are provided, all-zero inputs are used, which may result in inf or NaNs. It is recommended to use `--ignore-exec-errors`

- `-n, --neff` (string): the compiled NEFF to profile
- `-s, --session-file` (string): the file to store profile session information in
- `--ignore-exec-errors`: ignore errors during execution
- `inputs` (positional args): list of inputs in the form of `<NAME> <FILE_PATH>` separated by space. Eg `IN1 x.npy IN2 y.npy`

The following `neuron-profile capture` arguments are only relevant for multi-worker jobs

- `--collectives-profile-id` (string): worker id which will be profiled. Passing all profiles all workers. (default: all)
- `-r, --collectives-workers-per-node` (int): the number of workers on the current node. The global worker id (rank) of worker `n` on current node is `collectives-worker-start-id+n`
- `--collectives-worker-count` (int): total number of Neuron workers across all nodes for this collectives run.
- `--collectives-worker-start-id` (int): The rank offset for the first worker on the current node. For example, if node 0 has workers 0,1 and node 1 has workers 2,3 then `collectives-worker-start-id` for node 0 and 1 will be 0 and 2, respectively. (default: 0)

**neuron-profile view** [parameters]

- `-n, --neff-path` (string): the compiled NEFF file location
- `-s, --session-file` (string): the profile results NTFF file location

- `-d, --session-dir` (string): directory containing profile files for multi-worker runs
- `--output-format` (string): how the processed profile should be presented. The default `db` write processed data to the database. `summary-text` and `summary-json` print the summary data as a table or json, respectively, without writing to the database. The `perfetto` option writes processed data to Perfetto's native protobuf based tracing format, and can be visualized in the Perfetto UI. The `JSON` option writes processed data to human-readable JSON. (default: `db`)
- `--output-file` (string): file path to write results to, if applicable for the given output format
- `--db-endpoint` (string): the endpoint of InfluxDB (default: `http://localhost:8086`)
- `--db-org` (string): the org name of InfluxDB
- `--db-bucket` (string): name of the InfluxDB bucket where ingested profile data is stored. Also used in the URL for viewing the profile (Optional)
- `--port` (int): the port number of the http server (default: `3001`)
- `--force`: force overwrite an existing profile in the database
- `--terminology`: print a helpful table of terminology used by the profiler
- `--enable-memory-tracker`: Enable Memory Tracker to view scratchpad usage over time with a breakdown of usage per tensor. This requires having set `XLA_IR_DEBUG=1` and `XLA_HLO_DEBUG=1` before NEFF compilation and passing `--enable-dge-notifs` when capturing the profile.

## FAQ

### Difference between TensorE and TensorMatrixE Rows in Timeline

- TensorE includes instruction trace for LoadStationary (LoadWeight)
- TensorMatrixE includes instruction trace for MultiplyMoving (Matmul)
- Both instruction traces happen on the same TensorE engine, but we separate them into two rows to de-clutter the timeline due to the background load stationary feature (loading stationary matrix for the next matmul in parallel to current matmul). See more info in [NKI architecture guide](#).

## Troubleshooting

### InfluxDB not installed

```
$ neuron-profile view -n file.neff -s profile.ntff
ERRO[0001] To install influxdb, go to https://portal.influxdata.com/downloads/ and
↪ follow the instructions there
influxdb not setup correctly: exec: "influx": executable file not found in $PATH
```

```
$ neuron-profile view -n file.neff -s profile.ntff
ERRO[0000]
influxdb token not setup correctly: exit status 1
Try executing "systemctl start influxdb" and "influx setup"
```

Running `neuron-profile view` without InfluxDB installed will result in an error and a pointer to the InfluxDB installation instructions. Please follow the provided instructions and retry.

## Too many open files

```
influxdb2client E! Write error: internal error: unexpected error writing points to_
↪ database: [shard 10677] open /home/ubuntu/.influxdbv2/engine/data/7caae65aaa48380d/
↪ autogen/10677/index/0/MANIFEST: too many open files
```

InfluxDB will encounter “too many open files” and out of memory errors after a few hundred buckets have been created. Two ways to solve this are to delete unused buckets or increase the system file descriptor limit.

To increase the file descriptor limit, add the following lines to `/etc/security/limits.d/efa.conf` and `/etc/security/limits.conf`:

```
*          soft    nofile    1048576
*          hard    nofile    1048576
```

Add the following lines to `/etc/sysctl.conf`

```
fs.file-max = 197341270
vm.max_map_count=1048576
```

Commit changes by running `sudo sysctl -p`.

## When viewing UI “FATAL - Failed metadata query”

If you are SSH port forwarding the web UI from a remote machine to your local desktop you will need to port forward both the web UI (3001) and the database (8086) like so:

```
ssh -L 3001:localhost:3001 -L 8086:localhost:8086 remote_machine
```

## Visual Artifacts when viewing profiles

Some users have reported visual artifacts when viewing certain profiles in browsers other than Chrome. If you encounter this issue, please try using Chrome. For more details, refer to the GitHub issue: <https://github.com/aws-neuron/aws-neuron-sdk/issues/1033>

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.3.2 Neuron Profiler 2.0 (Beta) User Guide

### Table of contents

- *Overview*
- *Capturing profiles*
- *JAX User Experience*
  - *JAX Setup*
  - *JAX Profiler*

- *Custom Annotations in JAX*
  - *JAX Profiling using environment variable*
  - *Full JAX Example*
- *PyTorch User Experience*
  - *PyTorch Setup*
  - *PyTorch Profiler*
  - *PyTorch Profiling using Environment Variable*
  - *Full PyTorch Example*
- *Non-framework Specific User Experience*
  - *Profile Capture Environment Variables*
  - *Example Capturing Profile of Application Using Environment Variables*
  - *Example Capturing Profile of ncnn-test Using Environment Variables*
  - *CLI reference for System Profiles*
  - *Example of using System Profiles CLI*
  - *neuron-profile inspect Output*
- *EKS User Experience*
- *Processing and Viewing Profiles*
  - *Neuron Profiler UI*
  - *Neuron Profiler System Profile UI*
  - *Viewing Profiles with Perfetto*
  - *Perfetto Output View Options*
  - *Generating JSON Output From Profiles*
  - *Processing only system or device profiles*
- *Filtering System Profiles*
  - *Capture-Time Filtering*
  - *Processing-Time Filtering*
- *Troubleshooting*
  - *Incomplete JAX Profiles*



## Overview

Neuron Profiler 2.0 offers a user-friendly experience for capturing and analyzing application performance through both high-level system profiles and detailed device-level profiles. Users can profile their workloads using framework-specific APIs within their application code or by setting an environment variable before execution. This tool supports profiling for both single-node and distributed workloads, integrating with environments such as ParallelCluster and EKS. Once captured, profile results can be explored through multiple interfaces: the Neuron Profiler UI, the open-source trace viewer [Perfetto](#), or by exporting to a human-readable JSON format. This flexibility in data capture and visualization enables users to gain comprehensive insights into their application's performance across various scenarios and scales.

---

**Note:** Neuron Profiler 2.0 is a set of new features currently in beta that enhance and simplify the experience of capturing and viewing profiles. It is not a replacement of *Neuron Profiler*, which is the existing feature set specifically for capturing and viewing device profiles.

---

## Capturing profiles

Neuron Profiler 2.0 offers several flexible options for capturing profiles. Users can either set an environment variable `NEURON_RT_INSPECT_ENABLE` or use the PyTorch or JAX profiling APIs from their application code for fine-grained control over which sections of their code are profiled. PyTorch and JAX users who prefer not to modify their application code can still enable profiling by setting the environment variable before running their application.

## JAX User Experience

### JAX Setup

Follow the JAX Setup instructions to install the required JAX Neuron Plugin and the latest Neuron Driver, Runtime and Tools packages.

### JAX Profiler

The JAX context-managed profiling API allows you to profile blocks of code. This will capture a system profile including a Neuron Runtime API trace and Python trace for your application code in the captured block. This will also capture device profiles for any compiled graphs (NEFFs) executed on NeuronCores within this block. To use the profiler, import the `jax` package.

```
import jax
```

Profiling is enabled for all code enclosed in the context when using `with jax.profiler.trace(os.environ["NEURON_RT_INSPECT_OUTPUT_DIR"]):`

---

**Note:** It is important to pass the output directory `os.environ["NEURON_RT_INSPECT_OUTPUT_DIR"]` to `with jax.profiler.trace` and run `export NEURON_RT_INSPECT_OUTPUT_DIR=<your output directory>` before enabling profiling. This ensures all captured profile data is saved to the correct output directory.

---

## Custom Annotations in JAX

To add custom annotations to blocks of code in your profile, you can use `jax.profiler.TraceAnnotation`. Annotation names can be created at runtime, such as in the [example here](#) using with `jax.profiler.TraceAnnotation("my_label"+str(i))`. For more information on TraceAnnotations, see the official [JAX documentation](#).

## JAX Profiling using environment variable

Instead of using the `jax.profiler` context manager, you can enable profiling for your entire application using an environment variable. This is desirable if you want to capture a profile without modifying your application code. To enable profiling with the environment variable `NEURON_RT_INSPECT_ENABLE=1` and `NEURON_RT_INSPECT_OUTPUT_DIR=.` /output before running your application.

For example:

```
# make sure to remove call to with jax.profiler.trace from python script
NEURON_RT_INSPECT_ENABLE=1 NEURON_RT_INSPECT_OUTPUT_DIR=./output python jax_script.py
```

When using the `NEURON_RT_INSPECT_ENABLE` environment variable instead of `jax.profiler`, system profiles will not contain a framework and application code trace, only Neuron Runtime API trace.

Do not set the `NEURON_RT_INSPECT_ENABLE` environment variable and use the `jax.profiler` within your application code at the same time. Use one or the other.

For more profiling options that can be set through environment variables, see the section [Profile Capture Environment Variables](#).

## Full JAX Example

Create a file `jax_script.py` which performs repeated matrix multiplications distributed across Neuron devices.

```
from functools import partial
import os
import jax
import jax.numpy as jnp

from jax.sharding import Mesh, NamedSharding, PartitionSpec as P
from jax.experimental.shard_map import shard_map
from time import sleep

os.environ["XLA_FLAGS"] = "--xla_dump_hlo_snapshots --xla_dump_to=./dump"

jax.config.update("jax_default_prng_impl", "rbg")

mesh = Mesh(jax.devices(), ('i',))

def device_put(x, pspec):
    return jax.device_put(x, NamedSharding(mesh, pspec))

lhs_spec = P('i', None)
lhs = device_put(jax.random.normal(jax.random.key(0), (128, 128)), lhs_spec)

rhs_spec = P('i', None)
```

(continues on next page)

(continued from previous page)

```

rhs = device_put(jax.random.normal(jax.random.key(1), (128, 16)), rhs_spec)

@jax.jit
@partial(shard_map, mesh=mesh, in_specs=(lhs_spec, rhs_spec), out_specs=rhs_spec)
def matmul_allgather(lhs_block, rhs_block):
    rhs = jax.lax.all_gather(rhs_block, 'i', tiled=True)
    return lhs_block @ rhs

with jax.profiler.trace(os.environ["NEURON_RT_INSPECT_OUTPUT_DIR"]):
    out = matmul_allgather(lhs, rhs)
    for i in range(10):
        with jax.profiler.TraceAnnotation("my_label"+str(i)):
            out = matmul_allgather(lhs, rhs)
            sleep(0.001)

expected = lhs @ rhs
with jax.default_device(jax.devices('cpu')[0]):
    equal = jnp.allclose(jax.device_get(out), jax.device_get(expected), atol=1e-3,
    ↪rtol=1e-3)
    print("Tensors are the same") if equal else print("Tensors are different")

```

Set your profile output directory and run the script:

```

export NEURON_RT_INSPECT_OUTPUT_DIR=./output
python jax_script.py

```

## PyTorch User Experience

### PyTorch Setup

Follow the [PyTorch Setup](#) instructions to install the required PyTorch Neuron packages as well as the latest Neuron Driver, Runtime and Tools.

### PyTorch Profiler

The PyTorch context-managed profiling API allows you to profile blocks of code. This will capture a system profile including a Neuron Runtime API trace and Python trace for your application code in the captured block. This will also capture device profiles for any compiled graphs executed on NeuronCores within this block. To use the profiler, import it in your application:

```

from torch_neuronx.experimental import profiler

```

Then profile a block of code using:

```

with torch_neuronx.experimental.profiler.profile(
    port=9012,
    profile_type='system',
    target='neuron_profile_perfetto',
    output_dir=os.environ['NEURON_RT_INSPECT_OUTPUT_DIR'],
    ms_duration=30000) as profiler:

```

After modifying your code to call the profiler, run your application as you normally would but set the environment variable `NEURON_RT_INSPECT_OUTPUT_DIR` to specify the output directory.

```
NEURON_RT_INSPECT_OUTPUT_DIR=./output python application.py
```

---

**Note:** it is essential to set `output_dir=os.environ['NEURON_RT_INSPECT_OUTPUT_DIR']` when starting the profiler from your application code. This ensures that all profile data sources dump to the same output directory.

---

## PyTorch Profiling using Environment Variable

Instead of using the `torch_neuronx.experimental.profiler.profile` context manager, you can enable profiling for your entire application using environment variable. This is desirable if you want to capture a profile without modifying your application code. To enable profiling with environment variable `NEURON_RT_INSPECT_ENABLE=1` and `NEURON_RT_INSPECT_OUTPUT_DIR=./output` before running your application.

For example

```
# make sure to remove call to with torch_neuronx.experimental.profiler.profile from
python script
NEURON_RT_INSPECT_ENABLE=1 NEURON_RT_INSPECT_OUTPUT_DIR=./output python pytorch_script.py
```

When using the `NEURON_RT_INSPECT_ENABLE` environment variable instead of `torch_neuronx.experimental.profiler.profile` system profiles will not contain a framework and application code trace, only Neuron Runtime API trace.

Do not set the `NEURON_RT_INSPECT_ENABLE` environment variable and use the `torch_neuronx.experimental.profiler.profile` within your application code at the same time. Use one or the other.

For more profiling options that can be set through environment variables, see the section *Profile Capture Environment Variables*.

## Full PyTorch Example

Create a file `train_torchrun_context.py` with the following contents

```
import os

import torch
import torch.nn as nn
import torch.nn.functional as F

# XLA imports
import torch_xla
import torch_xla.core.xla_model as xm
import torch_xla.debug.profiler as xp

import torch_neuronx
from torch_neuronx.experimental import profiler

os.environ["NEURON_CC_FLAGS"] = "--cache_dir=./compiler_cache"
```

(continues on next page)

(continued from previous page)

```

# Global constants
EPOCHS = 2

# Declare 3-layer MLP Model
class MLP(nn.Module):
    def __init__(self, input_size=10, output_size=2, layers=[5, 5]):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return F.log_softmax(x, dim=1)

def main():
    # Fix the random number generator seeds for reproducibility
    torch.manual_seed(0)

    # XLA: Specify XLA device (defaults to a NeuronCore on Trn1 instance)
    device = xm.xla_device()

    # Start the profiler context-manager
    with torch_neuronx.experimental.profiler.profile(
        port=9012,
        profile_type='system',
        target='neuron_profile_perfetto',
        output_dir=os.environ['NEURON_RT_INSPECT_OUTPUT_DIR'],
        ms_duration=30000) as profiler:

        # IMPORTANT: the model has to be transferred to XLA within
        # the context manager, otherwise profiling won't work
        model = MLP().to(device)
        optimizer = torch.optim.SGD(model.parameters(), lr=0.01)
        loss_fn = torch.nn.NLLLoss()

        # start training loop
        print('-----Training -----')
        model.train()
        for epoch in range(EPOCHS):
            optimizer.zero_grad()
            train_x = torch.randn(1, 10).to(device)
            train_label = torch.tensor([1]).to(device)

            # forward
            loss = loss_fn(model(train_x), train_label)

            # back
            loss.backward()
            optimizer.step()

```

(continues on next page)

(continued from previous page)

```

    # XLA: collect ops and run them in XLA runtime
    xm.mark_step()

    print('-----End Training -----')

if __name__ == '__main__':
    main()

```

Run this workload with the following command:

```
NEURON_RT_INSPECT_OUTPUT_DIR="output" python simple_demo.py
```

## Non-framework Specific User Experience

You can also control profiling with environment variables. This is useful when you can't easily change your application code, such as when running an executable which calls the Neuron Runtime or in a containerized environment where the application code is built into the container image.

### Profile Capture Environment Variables

- `NEURON_RT_INSPECT_ENABLE`: Set to 1 to enable system and device profiles. For control over which profile types are captured use `NEURON_RT_INSPECT_SYSTEM_PROFILE` and `NEURON_RT_INSPECT_DEVICE_PROFILE`.
- `NEURON_RT_INSPECT_OUTPUT_DIR`: The directory where captured profile data will be saved to. Defaults to `./output`.
- `NEURON_RT_INSPECT_SYSTEM_PROFILE`: Set to 0 to disable the capture of system profiles. Defaults to 1 when `NEURON_RT_INSPECT_ENABLE` is set to 1.
- `NEURON_RT_INSPECT_DEVICE_PROFILE`: Set to 0 to disable the capture of device profiles. Defaults to 0 when `NEURON_RT_INSPECT_ENABLE` is set to 1.

### Example Capturing Profile of Application Using Environment Variables

Instead of using the PyTorch or JAX profilers you can profile your Python application (or any application calling the Neuron Runtime API) using environment variables.

```
NEURON_RT_INSPECT_ENABLE=1 NEURON_RT_INSPECT_OUTPUT_DIR=./output python app.py
```

See *Profile Capture Environment Variables* for other profiling options that can be set via environment variable.

## Example Capturing Profile of nccom-test Using Environment Variables

Profiling can be enabled using environment variables. For simplicity, we have a quick way to generate a Neuron workload through using `nccom-test`. `nccom-test` is a benchmarking tool which is already available with Neuron AML.

```
export NEURON_RT_INSPECT_ENABLE=1
export NEURON_RT_INSPECT_OUTPUT_DIR=./output
nccom-test allr allg -b 512kb -e 512kb -r 32 -n 10 -d fp32 -w 1 -f 512
```

**Note:** If you have problems with `nccom-test` add the `--debug` flag. If using a `trn1.2xlarge` instance, change `-r 32` to `-r 2` to use fewer neuron cores.

To understand the profiling output see this section: [Inspect Output](#)

## CLI reference for System Profiles

In addition to controlling profiling with environment variables, you can use the `neuron-profile inspect` command line interface for profiling applications. This provides the same functionality as environment variables but helps you avoid typos, invalid arguments, and provides a useful `--help` command to explain available options.

```
Usage:
neuron-profile [OPTIONS] inspect [inspect-OPTIONS] [userscript...]

Application Options:
-v, --version                Show version and exit

Help Options:
-h, --help                  Show this help message

[inspect command options]
  -o, --output-dir=          Output directory for the captured profile data,
  ↪ including system and device profiles (default: ./output)
  -n, --num-trace-events=    Maximum number of trace events to capture when
  ↪ profiling. Once hitting this limit, no new events are recorded
  --capture-system-profiles  Disable capture of system profile data. Can reduce
  ↪ output size.
  --capture-device-profiles  Disable capture of device profile data. Can reduce
  ↪ output size.

[inspect command arguments]
userscript:                  Run command/script that launches a Neuron workload. E.
  ↪ g. 'python app.py' or './runscript.sh'
```

## Example of using System Profiles CLI

User can provide any type of their own script to generate a Neuron workload such as Pytorch to the System Profiles CLI. For simplicity, we have a quick way to generate a Neuron workload through using `nccom-test`. `nccom-test` is a benchmarking tool which is already available with Neuron AMI and `aws-neuronx-tools` package.

```
ubuntu@ip-172-31-63-210:~$ neuron-profile inspect -o inspect-output-nccom-test nccom-
↪test allg -b 512kb -e 512kb -r 32 -n 10 -d fp32 -w 1 -f 512
INFO[0000] Running command "nccom-test allg -b 512kb -e 512kb -r 32 -n 10 -d fp32 -w 1 -
↪f 512" with profiling enabled
      size(B)    count(elems)    type    time:avg(us)    algbw(GB/s)    busbw(GB/s)
      524288      131072      fp32         24.15         21.71         21.03
Avg bus bandwidth: 21.0339GB/s
```

**Note:** If you have problems with `nccom-test` add the `-debug` flag. If using a `trn1.2xlarge` instance, change `-r 32` to `-r 2` to use fewer neuron cores.

### neuron-profile inspect Output

The above command shows a Neuron workload execution is being traced and output to `inspect-output-nccom-test` directory. You will see the output directory contains a single NEFF file and a device profile (NTFF) for all Neuron Cores which executed that NEFF. You will also see `ntrace.pb` and `trace_info.pb` files storing the system profile data. Below showing what the outputs will look like:

```
ubuntu@ip-172-31-63-210:~$ tree inspect-output-nccom-test
inspect-output-nccom-test
├── i-012590440bb9fd263_pid_98399
│   ├── 14382885777943380728_instid_0_vnc_0.ntff
│   ├── 14382885777943380728_instid_0_vnc_1.ntff
│   ├── 14382885777943380728_instid_0_vnc_10.ntff
│   └── 14382885777943380728_instid_0_vnc_11.ntff
├── ...
├── 14382885777943380728_instid_0_vnc_8.ntff
├── 14382885777943380728_instid_0_vnc_9.ntff
├── cpu_util.pb
├── host_mem.pb
├── neff_14382885777943380728.neff
├── ntrace.pb
└── trace_info.pb

2 directories, 74 files
```

To view a summary of the captured profile data run the command

```
neuron-profile view -d inspect-output-nccom-test --output-format summary-text
```



## EKS User Experience

Capturing a profile on EKS is most easily done through setting of environment variables as described in the section *Non-framework specific User Experience*. By using environment variables, users do not need to change application code in their container image or modify their run commands.

Update the deployment yaml to include the `NEURON_RT_INSPECT_ENABLE` and `NEURON_RT_INSPECT_OUTPUT_DIR` environment variables. For distributed workloads, it's important that `NEURON_RT_INSPECT_OUTPUT_DIR` points to a directory on a shared volume which all workers have access to.

```
apiVersion: v1
kind: Pod
metadata:
  name: trn1-mlp
spec:
  restartPolicy: Never
  schedulerName: default-scheduler
  nodeSelector:
    beta.kubernetes.io/instance-type: trn1.32xlarge
  containers:
    - name: trn1-mlp
      env:
        - name: NEURON_RT_INSPECT_ENABLE
          value: "1"
        - name: NEURON_RT_INSPECT_OUTPUT_DIR
          value: "/shared/output"
      command: ['torchrun']
      args:
        - '--nnodes=1'
        - '--nproc_per_node=32'
        - 'train_torchrun.py'
      image: ${ACCOUNT_ID}.dkr.ecr.${REGION}.amazonaws.com/${REPO}:mlp
      imagePullPolicy: IfNotPresent
      resources:
        limits:
          aws.amazon.com/neuron: 16
```

**Note:** EKS users running PyTorch and JAX applications are still free to change their application code and use the PyTorch or JAX Python profiling APIs if they want finer-grained control over profiling. However, using the environment variables conveniently allows profiling without modifying the container image or application code.

## Processing and Viewing Profiles

Users have three output options for interacting with their captured profiles

- Neuron Profiler UI - Neuron's custom UI which allows easily drilling down to detailed device profiles from high level system profiles
- Perfetto - Allows sharing profiles as a single file and viewing your profiles in the Perfetto UI at <https://ui.perfetto.dev/>
- JSON - human-readable text output that enables simple scripting

## Neuron Profiler UI

To view a profile in the Neuron Profiler UI run the following command to process a profile and launch the UI

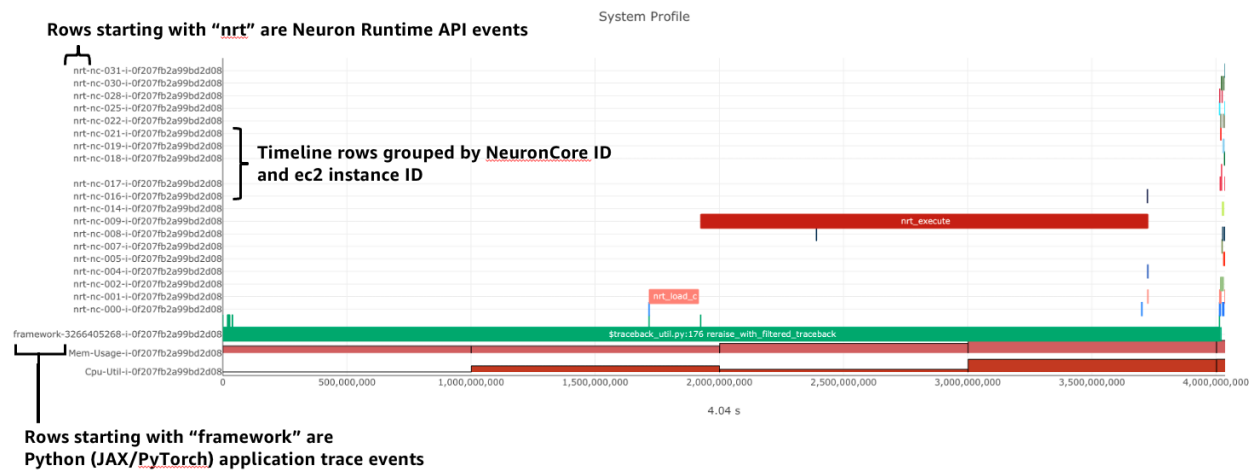
```
neuron-profile view -d ./output
```

To view profiles with the Neuron Profiler UI running locally you will need to have InfluxDB installed on your system. To install and setup InfluxDB follow the [directions in the official Neuron Profile documentation](#).

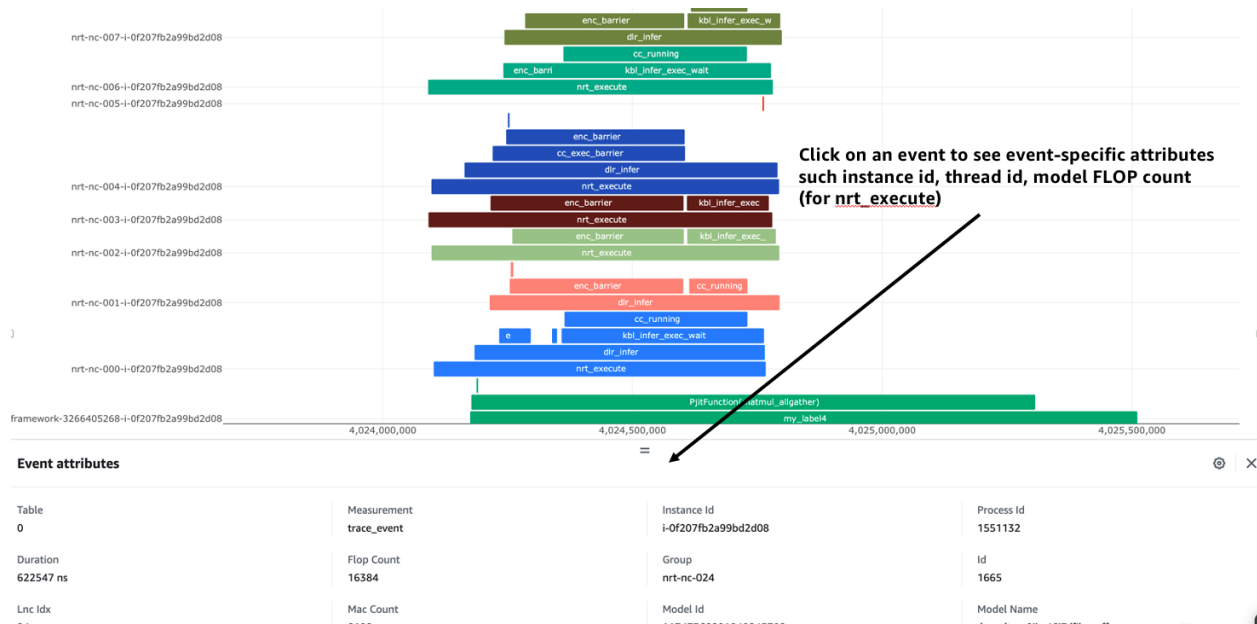
## Neuron Profiler System Profile UI

The system profile timeline shows a trace of Neuron Runtime API calls, ML framework function calls, CPU utilization, and memory usage on each of the instances in your workload. The Neuron Runtime API trace is grouped by NeuronCore IDX and ec2 instance ID. For example, all events in the row labeled nrt-nc-003-i-0f207fb2a99bd2d08 are associated with NeuronCore 3 and instance i-0f207fb2a99bd2d08.

Framework function traces are grouped by thread id and ec2 instance id. For example, all events in the row framework-3266405268-i-0f207fb2a99bd2d08 are framework or application function calls made on thread 3266405268 running on instance i-0f207fb2a99bd2d08.

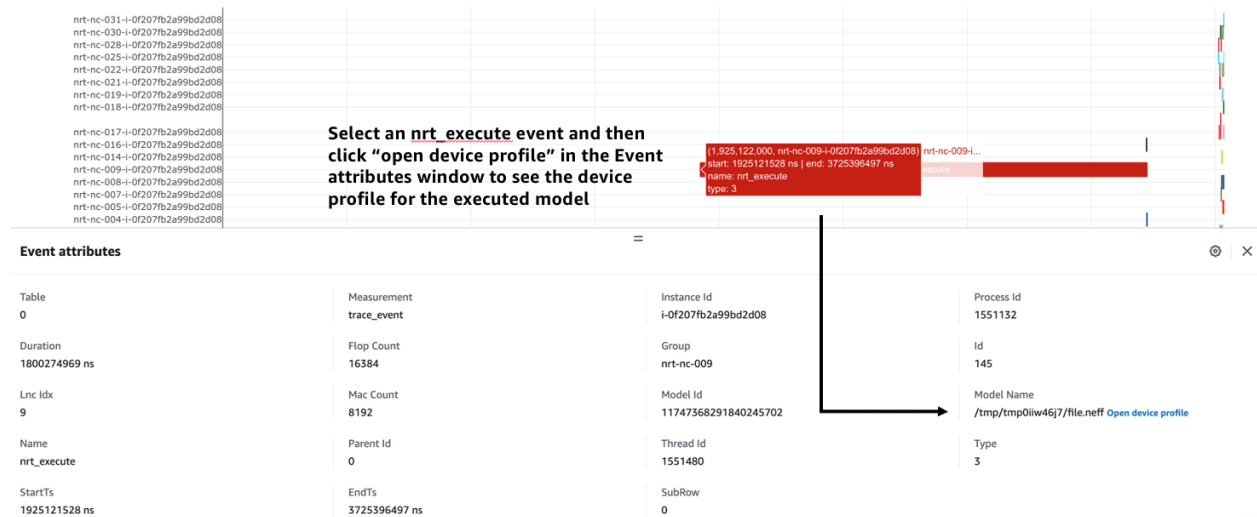


Clicking on trace events in the timeline shows a "Event attributes" view with a list of attributes associated with that event. For example, clicking on an nrt\_execute event (the Neuron Runtime API call for executing a compiled model on a NeuronCore) will show events such as Flop count (the number of floating point operations for a single execution of the model), the model name, and the NeuronCore idx and ec2 instance id associated with the function call.

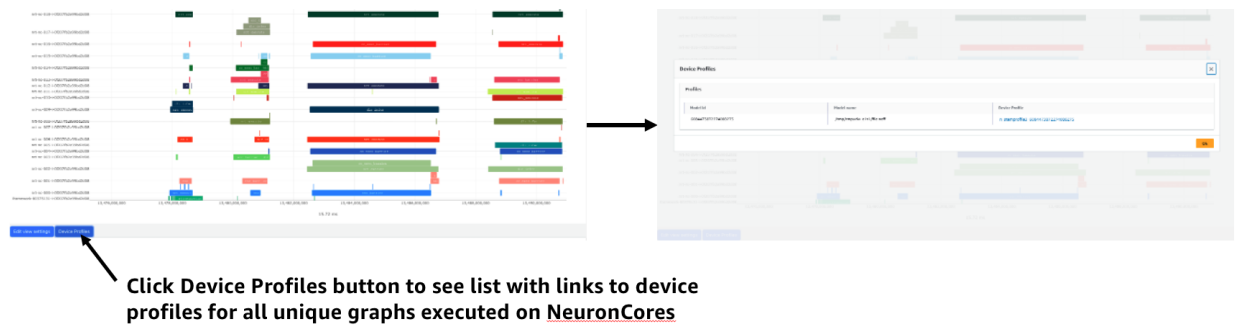


Neuron Profiler 2.0 allows users to drill-down from a system timeline to a device profile timeline in order to see a detailed view of hardware activity during the execution of a graph. To do this, select an `nrt_execute` event in the timeline and in the “Event attributes” view select the “Open device profile” button under the Model Name attribute. This will open a new window with a device profile. For help understanding a device profile see the section documentation section “Understanding a Neuron Profile”

## Drill-down to device profiles



To see a list of all device profiles that were captured during your workload press the “Device Profiles” button at the bottom of the timeline. From this list you can see all unique compiled graphs (NEFFs) that were executed on NeuronCores during your workload. For each graph there is a link to a device profile that will show a detailed view of hardware activity on the NeuronCore during execution of this graph.



## Viewing Profiles with Perfetto

Perfetto is an open-source trace analysis toolkit with a powerful UI for visualizing and analyzing trace data. Users of Neuron Profiler have the option of viewing their profiles in the Perfetto UI.

The `--output-format perfetto` option writes processed data to Perfetto's native protobuf-based tracing format which can be visualized in the Perfetto UI at <https://ui.perfetto.dev/>.

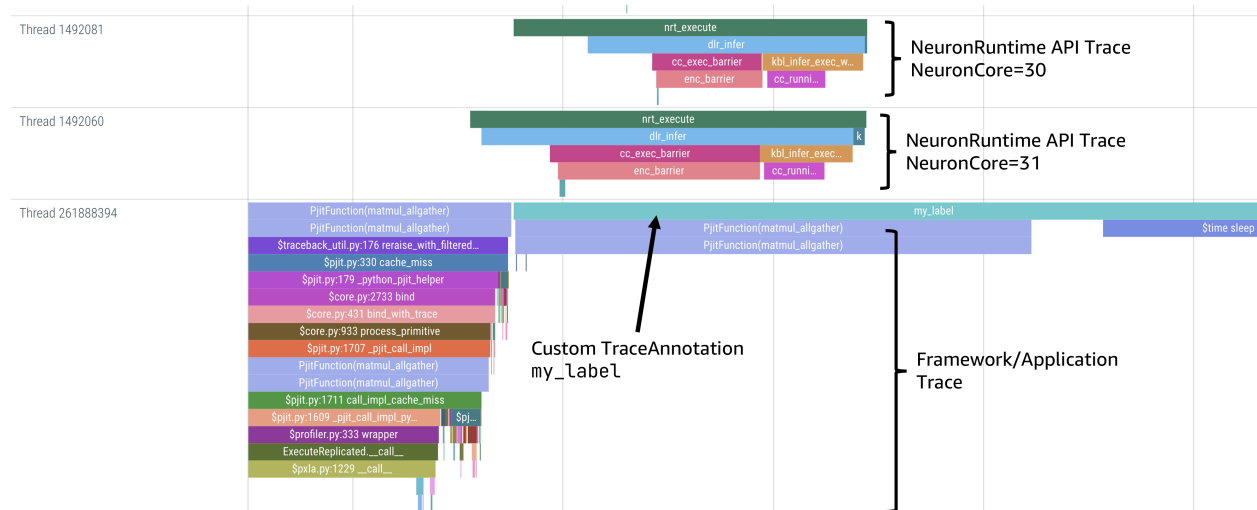
Example:

```
neuron-profile view -d ./output --output-format perfetto
```

This will generate a `system_profile.pftrace` file for the system profile and a `device_profile_model_<model_id>.pftrace` file for each unique compiled model that was executed on a Neuron Device.

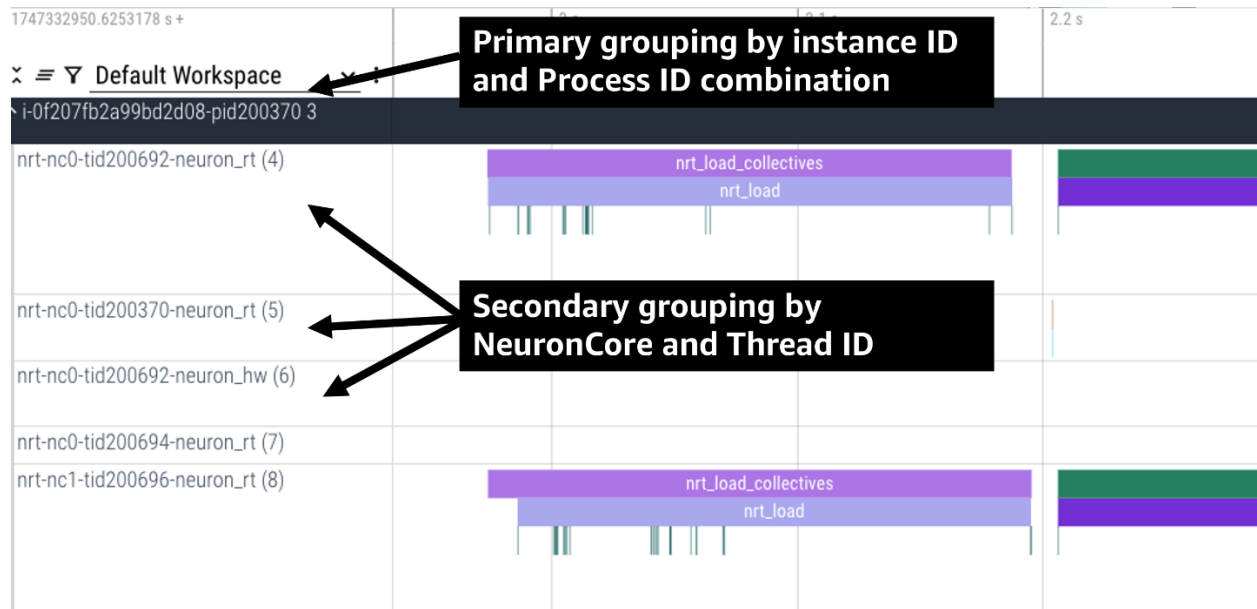
To view the system profile, go to <https://ui.perfetto.dev/> and open the `system_profile.pftrace` file.

**Note:** When loading trace files in the Perfetto UI, your data is processed locally and not uploaded to Perfetto's servers.



To view a device profile go to <https://ui.perfetto.dev/> and open the `device_profile_model_<model_id>.pftrace` file. This will show a detailed view of hardware activity on the NeuronCore during execution of this graph.

**Note:** Your browser may run out of memory when viewing \*.pfttrace (Perfetto trace) files that are more than a few hundred MB. See the section [Viewing Large Profiles in Perfetto](#) for directions on how to view large traces using the trace processor.



## Generating JSON Output From Profiles

The `--output-format json` option writes processed profile data to human-readable JSON that can be used for scripting and manual inspection.

```
neuron-profile view -d ./output --output-format json
```

This will generate a `system_profile.json` file containing the system profile data and a `device_profile_model_<model_id>.json` file for each unique compiled model that was executed on a Neuron Device.

The `system_profile.json` JSON contains the following data types:

- `trace_events`: Neuron Runtime API trace events and Framework/Application trace events containing timestamps, durations, names, and the ec2 instance-id to differentiate between events from different compute nodes in a distributed workload.

```
{
  "Neuron_Runtime_API_Event": {
    "duration": 27094,
    "group": "nrt-nc-000",
    "id": 1,
    "instance_id": "i-0f207fb2a99bd2d08",
    "lnc_idx": "0",
    "name": "nrt_tensor_write",
    "parent_id": 0,
    "process_id": "1627711",
    "size": "4",
    "tensor_id": "4900392441224765051",
    "tensor_name": "_unknown_",
    "thread_id": 1627711,
    "timestamp": 1729888371056597613,
    "type": 11
  }
}
```

(continues on next page)

(continued from previous page)

```

    },
    "Framework_Event": {
      "duration": 3758079,
      "group": "framework-80375131",
      "instance_id": "i-0f207fb2a99bd2d08",
      "name": "PjitFunction(matmul_allgather)",
      "process_id": "701",
      "thread_id": 80375131,
      "timestamp": 1729888382798557372,
      "type": 99999
    }
  }
}

```

- mem\_usage: sampled host memory usage

```

{
  "duration": 1,
  "instance_id": "i-0f207fb2a99bd2d08",
  "percent_usage": 9.728179797845964,
  "timestamp": 1729888369286687792,
  "usage": 51805806592
}

```

- cpu\_util: sampled CPU utilization. Results are provided per core and per ec2 instance involved in a distributed workload

```

{
  "cpu_id": "47",
  "duration": 1,
  "instance_id": "i-0f207fb2a99bd2d08",
  "timestamp": 1729888371287337243,
  "util": 2.3255813
},

```

### Processing only system or device profiles

To reduce processing times it is possible to skip processing of system or device profiles. Sometimes users may only be interested in one or want to start with a limited set of profiling data before exploring the full profile.

To skip processing of device profiles use the `--ignore-device-profile` option. To skip processing of system profiles use the `--ignore-system-profile` option. These options can be used with the `--output-format` values `db` (default), `perfetto`, or `json`.

For example:

```
neuron-profile view -d ./output --ignore-device-profile --output-format perfetto
```

## Filtering System Profiles

This guide explains how to filter system trace events to optimize memory usage, reduce output size, and speed up trace processing. **Capture-time filtering** reduces memory usage and trace file size by only collecting specific events, but filtered data cannot be recovered later. **Processing-time filtering** preserves the complete trace and allows flexible analysis with different filters, but requires more memory and storage during capture.

### Capture-Time Filtering

Configure filters before trace capture using environment variables or API functions. You can use NeuronCore filters to only capture events for specific NeuronCores (for example only events associated with NeuronCore 0 or all the NeuronCores on a specific NeuronDevice). You can use event type filters to only capture specific events (for example model execute or collectives events). It is possible to combine both NeuronCore and event type filters.

### Filtering by NeuronCore

If capture is enabled for a NeuronCore then a ring buffer will be allocated in host memory for storing those core's events. Thus filtering by NeuronCore decreases host memory usage during capture.

### Default Behavior

By default, all visible NeuronCores are enabled for capture.

### Using Environment Variables

```
# Filter to capture events only from NeuronCore 0
export NEURON_RT_INSPECT_EVENT_FILTER_NC=0

# Filter to capture events from NeuronCores 0, 2, and 4
export NEURON_RT_INSPECT_EVENT_FILTER_NC=0,2,4

# Filter to capture events from a range of NeuronCores (0 through 3)
export NEURON_RT_INSPECT_EVENT_FILTER_NC=0-3

# Reset to default behavior
unset NEURON_RT_INSPECT_EVENT_FILTER_NC # Back to capturing all visible cores
```

### Using API Functions

```
#include <nrt/nrt_sys_trace.h>

// Allocate and configure trace options
nrt_sys_trace_config_t *config;
nrt_sys_trace_config_allocate(&config);
nrt_sys_trace_config_set_defaults(config);

// Enable capture only for specific NeuronCores
```

(continues on next page)



(continued from previous page)

```

// Disable all cores since by default they are all enabled
int num_cores = 128;
for (int i=0; i<num_cores; i++) {
    nrt_sys_trace_config_set_capture_enabled_for_nc(config, i, false); // disable NC i
}

// Then enable specific cores
nrt_sys_trace_config_set_capture_enabled_for_nc(config, 0, true); // Enable NC 0
nrt_sys_trace_config_set_capture_enabled_for_nc(config, 2, true); // Enable NC 2

// Start tracing with the configuration
nrt_sys_trace_start(config);

// Your application code here...

// Stop tracing and cleanup
nrt_sys_trace_stop();
nrt_sys_trace_config_free(config);

```

## Filtering by Event Type

### Default Behavior

By default, all event types are enabled for capture.

### Getting Available Event Types

You can discover all available event types using the `nrt_sys_trace_get_event_types` API.

```

#include <nrt/nrt_sys_trace.h>

// Get all available event types
const char **event_types = nullptr;
size_t count = 0;
NRT_STATUS status = nrt_sys_trace_get_event_types(&event_types, &count);

if (status == NRT_SUCCESS) {
    printf("Available event types:\n");
    for (size_t i = 0; i < count; ++i) {
        printf("  %s\n", event_types[i]);
    }

    // Free the event types array
    for (size_t i = 0; i < count; ++i) {
        free((void*)event_types[i]);
    }
    free((void*)event_types);
}

```

## Using Environment Variables

The `NEURON_RT_INSPECT_EVENT_FILTER_TYPE` environment variable supports:

- **Default:** If not set, all event types are captured
- **Specific event types:** Use exact event names from `nrt_sys_trace_get_event_types()`
- **Event categories:** Use hardware or software to filter by category
- **Exclusion:** Use `^` prefix to exclude specific events from a category

```
# Filter to capture only specific event types
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=model_load,nrt_execute,runtime_execute

# Filter to capture all hardware events
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=hardware

# Filter to capture all software events
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=software

# Filter to capture all hardware events EXCEPT cc_exec
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=hardware,^cc_exec

# Filter to capture all software events EXCEPT model_load
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=software,^model_load

# Mix categories and specific events
export NEURON_RT_INSPECT_EVENT_FILTER_TYPE=hardware,tensor_read,tensor_write

# Reset to default behavior
unset NEURON_RT_INSPECT_EVENT_FILTER_TYPE # Back to capturing all event types
```

The hardware group contains events that are executed on the ML accelerator. These are `nc_exec_running`, `cc_running`, `cc_exec_barrier`, `numerical_err`, `nrt_model_switch`, `timestamp_sync_point`, `hw_notify`. The software group contains all other events.

## Using API Functions

Use the `nrt_sys_trace_config_set_capture_enabled_for_event_type` API to filter by event type.

```
#include <nrt/nrt_sys_trace.h>

// Configure trace options
nrt_sys_trace_config_t *config;
nrt_sys_trace_config_allocate(&config);
nrt_sys_trace_config_set_defaults(config); // By default, all event types are enabled

// Disable specific event types (others remain enabled)
nrt_sys_trace_config_set_capture_enabled_for_event_type(config, "device_exec", false);

// Or disable all first, then enable only specific ones
const char **all_event_types = nullptr;
size_t all_count = 0;
```

(continues on next page)

(continued from previous page)

```

nrt_sys_trace_get_event_types(&all_event_types, &all_count);

// Disable all event types first
for (size_t i = 0; i < all_count; ++i) {
    nrt_sys_trace_config_set_capture_enabled_for_event_type(config, all_event_types[i],
    ↪ false);
}

// Enable only specific event types
nrt_sys_trace_config_set_capture_enabled_for_event_type(config, "model_load", true);
nrt_sys_trace_config_set_capture_enabled_for_event_type(config, "nrt_execute", true);

// Verify which event types are enabled
const char **enabled_types = nullptr;
size_t enabled_count = 0;
nrt_sys_trace_config_get_enabled_event_types(config, &enabled_types, &enabled_count);
printf("Enabled event types: %zu\n", enabled_count);
for (size_t i = 0; i < enabled_count; ++i) {
    printf("  %s\n", enabled_types[i]);
}

// Clean up memory (caller is responsible)
for (size_t i = 0; i < enabled_count; ++i) {
    free((void*)enabled_types[i]);
}
free((void*)enabled_types);

for (size_t i = 0; i < all_count; ++i) {
    free((void*)all_event_types[i]);
}
free((void*)all_event_types);

// Start tracing
nrt_sys_trace_start(config);

// Your application code here...

// Cleanup
nrt_sys_trace_stop();
nrt_sys_trace_config_free(config);

```

## Tips

1. **Memory Optimization:** Use NeuronCore filtering to avoid allocating ring buffers for unused cores and decrease host memory usage. Use both event type or NeuronCore to decrease output trace sizes.
2. **Event Type Discovery:** Use `nrt_sys_trace_get_event_types()` to discover available event types
3. **Category Filtering:** Use hardware/software categories for broad filtering
4. **Exclusion Filtering:** Use ^ prefix to exclude specific events from categories
5. **Combine Filters:** Use both NeuronCore and event type filters together for maximum optimization

## Processing-Time Filtering

Apply filters when viewing or processing already captured profiles. This approach allows you to analyze the same trace data in different ways without recapturing. The filters can be used for any neuron-profile output format including `--output-format json` and `--output-format perfetto`.

### Filtering by NeuronCore

Use the `--system-trace-filter-neuron-core` to only process events for specific NeuronCores. The IDs are local to the instance and not global IDs.

If the `--system-trace-filter-neuron-core` argument is not set then events from all NeuronCores will be included in the processed trace.

```
# Filter by single neuron core
neuron-profile view -d ./output --system-trace-filter-neuron-core "0" --output-format-
↪ perfetto

# Filter by multiple neuron cores
neuron-profile view -d ./output --system-trace-filter-neuron-core "0,1,2,3" --output-
↪ format perfetto
```

### Filtering by Event Type

Use the `--system-trace-filter-event-type` to only process specific trace events types.

If the `--system-trace-filter-event-type` argument is not set then all event types will be included in the processed trace.

```
# Filter by single event type
neuron-profile view -d ./output --system-trace-filter-event-type "nrt_execute" --output-
↪ format perfetto

# Filter by multiple event types
neuron-profile view -d ./output --system-trace-filter-event-type "nrt_execute,nrt_load" -
↪ -output-format perfetto
```

### Filtering by Instance ID

Use the `--system-trace-filter-instance-id` to only process events for specific ec2 instances.

If the `--system-trace-filter-instance-id` argument is not set then events from all instances will be included in the processed trace.

```
# Filter by single instance
neuron-profile view -d ./output --system-trace-filter-instance-id "i-abc123" --output-
↪ format perfetto

# Filter by multiple instances (comma-separated)
neuron-profile view -d ./output --system-trace-filter-instance-id "i-abc123,i-def456,i-
↪ ghi789" --output-format perfetto
```

## Troubleshooting

### Incomplete JAX Profiles

If your JAX profile has fewer events than expected or lacks the Runtime API trace, check whether `jax.profiler.stop_trace` is being called inside a `with jax.profiler.trace` context block. This can prematurely stop tracing. Use `jax.profiler.stop_trace` only when profiling was started with `jax.profiler.start_trace`, not when using the context-managed `with jax.profiler.trace` API.

Also when using `jax.profiler` within your script ensure that the environment variable `NEURON_RT_INSPECT_ENABLE` is not set to 1. Additionally, ensure that `NEURON_RT_INSPECT_OUTPUT_DIR` is set to the correct output directory and this is the output directory passed to `with jax.profiler.trace`.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.4 Third-party solutions

AWS Neuron integrates with multiple third-party partner solutions that allow you to run deep learning workloads on Amazon EC2 instances powered by AWS Trainium and AWS Inferentia chips. The following list gives an overview of third-party solutions that work with AWS Neuron.

### Table of contents

- [Weights & Bias](#)
- [Datadog](#)

### 6.4.1 Weights & Bias

Weights & Biases is a machine learning platform for developers to build better models faster. Use W&B's lightweight, interoperable tools to quickly track experiments, version and iterate on datasets, evaluate model performance, reproduce models, visualize results and spot regressions, and share findings with colleagues.

[Weights & Bias documentation](#)

### 6.4.2 Datadog

Datadog, an observability and security platform, provides real-time monitoring for cloud infrastructure and ML operations. Datadog is excited to launch its AWS Neuron integration, which pulls metrics collected by Neuron SDK's Neuron Monitor tool into Datadog, enabling users to track the performance of their Trainium and Inferentia-based instances. By providing real-time visibility into model performance and hardware usage, Datadog helps customers ensure efficient training and inference, optimized resource utilization, and the prevention of service slowdowns.

[Datadog documentation](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 6.5 Other Tools

*This document is relevant for: Inf1*

### 6.5.1 Neuron Check Model

#### Overview

Neuron Check Model tool provides user with basic information about the compiled and uncompiled model's operations without the use of TensorBoard-Neuron. For additional visibility into the models, please see [Neuron Plugin for TensorBoard \(Inf1\)](#).

Neuron Check Model tool scans the user's uncompiled model and provides a table of the operations within the uncompiled model. By default, the table shows each operation type and number of instances of that type within model, and whether the type is supported in Neuron. If `--show_names` option is specified, the table shows each operation by name and whether the type of that operation is supported in Neuron.

If the model is already compiled, the tool also provides the table of operations as for uncompiled model. The table include the Neuron subgraph type and number of instances of that type, along with operations that have not been compiled to Neuron. Additionally, the tool displays a message showing the minimum number of NeuronCores required to run the model, followed by another table which shows the list of Neuron subgraphs by name and the number of pipelined NeuronCores used by each subgraph. More information about NeuronCore pipeline can be found in [NeuronCore Pipeline](#). If `--expand_subgraph` option is specified, the operations within each subgraph are printed below the subgraph information.

Neuron Check Model tool is currently available for TensorFlow and MXNet. To check PT model, please use `torch.neuron.analyze_model` function as shown in PyTorch-Neuron Getting Started tutorial [/src/examples/pytorch/resnet50.ipynb](#)

#### TensorFlow-Neuron Check Model

The following example shows how to run TensorFlow-Neuron Check Model tool with TensorFlow ResNet50 tutorial.

1. Start with the TensorFlow ResNet50 tutorial at `/src/examples/tensorflow/tensorflow_resnet50/resnet50.ipynb` and do the first three steps of the tutorial. Please stay in the Python environment that you setup during the tutorial.
2. Install needed `tensorflow_hub` package and download the tool:

```
pip install tensorflow_hub
wget https://raw.githubusercontent.com/aws/aws-neuron-sdk/master/src/neuron-gatherinfo/
  ↪ tf_neuron_check_model.py
python tf_neuron_check_model.py -h
```

```
usage: tf_neuron_check_model.py [-h] [--show_names] [--expand_subgraph]
                                model_path

positional arguments:
  model_path            a TensorFlow SavedModel directory (currently supporting
                        TensorFlow v1 SaveModel only).

optional arguments:
  -h, --help            show this help message and exit
  --show_names          list operation by name instead of summarizing by type
```

(continues on next page)

(continued from previous page)

```
(caution: this option will generate many lines of output
for a large model).
--expand_subgraph show subgraph operations.
```

3. After step 3 of the TensorFlow ResNet50 tutorial, you can check the uncompiled model to see Neuron supported operations (currently supporting TensorFlow v1 SaveModel only):

```
$ python tf_neuron_check_model.py ws_resnet50/resnet50/

* The following table shows the supported and unsupported operations within this
uncompiled model.
* Each line shows an operation type, the number of instances of that type within model,
* and whether the type is supported in Neuron.
* Some operation types are excluded from table because they are no-operations or
training-related operations:
['Placeholder', 'PlaceholderWithDefault', 'NoOp', 'Const', 'Identity', 'IdentityN',
'VarHandleOp',
'VarIsInitializedOp', 'AssignVariableOp', 'ReadVariableOp', 'StringJoin',
'ShardedFilename', 'SaveV2',
'MergeV2Checkpoints', 'RestoreV2']
```

| Op Type          | Num Instances | Neuron Supported ? |
|------------------|---------------|--------------------|
| Pad              | 2             | Yes                |
| RandomUniform    | 54            | Yes                |
| Sub              | 54            | Yes                |
| Mul              | 54            | Yes                |
| Add              | 54            | Yes                |
| Conv2D           | 53            | Yes                |
| BiasAdd          | 54            | Yes                |
| FusedBatchNormV3 | 53            | Yes                |
| Relu             | 49            | Yes                |
| MaxPool          | 1             | Yes                |
| AddV2            | 16            | Yes                |
| Fill             | 56            | Yes                |
| Mean             | 1             | Yes                |
| MatMul           | 1             | Yes                |
| Softmax          | 1             | Yes                |
| Pack             | 1             | Yes                |

```

* Total inference operations: 504
* Total Neuron supported inference operations: 504
* Percent of total inference operations supported by Neuron: 100.0
```

4. You can also check the compiled model to see the number of pipeline NeuronCores for each subgraph:

```
$ python tf_neuron_check_model.py ws_resnet50/resnet50_neuron/

* Found 1 Neuron subgraph(s) (NeuronOp(s)) in this compiled model.
* Use this tool on the original uncompiled model to see Neuron supported operations.
* The following table shows all operations, including Neuron subgraphs.
* Each line shows an operation type, the number of instances of that type within model,
```

(continues on next page)

(continued from previous page)

```
* and whether the type is supported in Neuron.
* Some operation types are excluded from table because they are no-operations or
↳ training-related operations:
['Placeholder', 'PlaceholderWithDefault', 'NoOp', 'Const', 'Identity', 'IdentityN',
↳ 'VarHandleOp',
'VarIsInitializedOp', 'AssignVariableOp', 'ReadVariableOp', 'StringJoin',
↳ 'ShardedFilename', 'SaveV2',
'MergeV2Checkpoints', 'RestoreV2']
```

| Op Type  | Num Instances | Neuron Supported ? |
|----------|---------------|--------------------|
| NeuronOp | 1             | Yes                |

```
* Please run this model on Inf1 instance with at least 1 NeuronCore(s).
* The following list show each Neuron subgraph with number of pipelined NeuronCores used
↳ by subgraph
* (and subgraph operations if --expand_subgraph is used):
```

| Subgraph Name  | Num   |
|--|-------|
| ↳ Pipelined NeuronCores  |       |
| -----  | ----- |
| ↳ -----  |       |
| conv5_block3_3_bn/FusedBatchNormV3/ReadVariableOp/neuron_op_d6f098c01c780733 | 1     |

5. When showing subgraph information, you can use `--expand_subgraph` to show operation types in each subgraph:

```
$ python tf_neuron_check_model.py ws_resnet50/resnet50_neuron/ --expand_subgraph
```

(output truncated to show subgraph information only)

| Subgraph Name  | Num   |
|--|-------|
| ↳ Pipelined NeuronCores  |       |
| -----  | ----- |
| ↳ -----  |       |
| conv5_block3_3_bn/FusedBatchNormV3/ReadVariableOp/neuron_op_d6f098c01c780733 | 1     |

| Op Type        | Num Instances |
|----------------|---------------|
| MatMul         | 1             |
| Relu           | 49            |
| Add            | 16            |
| FusedBatchNorm | 53            |
| BiasAdd        | 54            |
| Conv2D         | 53            |
| Pad            | 2             |
| Mean           | 1             |
| MaxPool        | 1             |
| Softmax        | 1             |

6. Use `--show_names` to see full operation names (caution: this option will generate many lines of output for a large model):

```
$ python tf_neuron_check_model.py ws_resnet50/resnet50_neuron/ --show_names
```

(continues on next page)



(continued from previous page)

```
* Found 1 Neuron subgraph(s) (NeuronOp(s)) in this compiled model.
* Use this tool on the original uncompiled model to see Neuron supported operations.
* The following table shows all operations, including Neuron subgraphs.
* Each line shows an operation name and whether the type of that operation is supported
  ↳ in Neuron.
* Some operation types are excluded from table because they are no-operations or
  ↳ training-related operations:
  ['Placeholder', 'PlaceholderWithDefault', 'NoOp', 'Const', 'Identity', 'IdentityN',
  ↳ 'VarHandleOp',
  'VarIsInitializedOp', 'AssignVariableOp', 'ReadVariableOp', 'StringJoin',
  ↳ 'ShardedFilename', 'SaveV2',
  'MergeV2Checkpoints', 'RestoreV2']
```

| Op Name  | Op Type  |
|--|----------|
| ↳ Neuron Supported ?   |          |
| -----  | -----    |
| ↳ -----  |          |
| conv5_block3_3_bn/FusedBatchNormV3/ReadVariableOp/neuron_op_d6f098c01c780733 | NeuronOp |
| ↳ Yes  |          |

```
* Please run this model on Inf1 instance with at least 1 NeuronCore(s).
* The following list show each Neuron subgraph with number of pipelined NeuronCores used
  ↳ by subgraph
* (and subgraph operations if --expand_subgraph is used):
```

| Subgraph Name  | Num   |
|--|-------|
| ↳ Pipelined NeuronCores  |       |
| -----  | ----- |
| ↳ -----  |       |
| conv5_block3_3_bn/FusedBatchNormV3/ReadVariableOp/neuron_op_d6f098c01c780733 | 1     |

## MXNet-Neuron Check Model

The following example shows how to run MXNet-Neuron Check Model tool with MXNet ResNet50 tutorial.

1. Start with the MXNet ResNet50 tutorial at `/src/examples/mxnet/resnet50/resnet50.ipynb` and do the first three steps of the tutorial. Please stay in the Python environment that you setup during the tutorial.

2. Download the tool:

```
wget https://raw.githubusercontent.com/aws/aws-neuron-sdk/master/src/neuron-gatherinfo/
  ↳ mx_neuron_check_model.py
python mx_neuron_check_model.py -h
```

```
usage: mx_neuron_check_model.py [-h] [--show_names] [--expand_subgraph]
                                model_path
```

positional arguments:

model\_path path prefix to MXNet model (the part before `-symbol.json`)

optional arguments:

(continues on next page)

(continued from previous page)

```
-h, --help          show this help message and exit
--show_names        list operation by name instead of summarizing by type
                    (caution: this option will generate many lines of output
                    for a large model).
--expand_subgraph   show subgraph operations.
```

3. After step 3 of MXNet ResNet50 tutorial, you can check the uncompiled model to see Neuron supported operations:

```
$ python mx_neuron_check_model.py resnet-50

* The following table shows the supported and unsupported operations within this
  ↳ uncompiled model.
* Each line shows an operation type, the number of instances of that type within model,
* and whether the type is supported in Neuron.
* Some operation types are excluded from table because they are no-operations or
  ↳ training-related operations:
  ['null']
```

| Op Type        | Num Instances | Neuron Supported ? |
|----------------|---------------|--------------------|
| BatchNorm      | 51            | Yes                |
| Convolution    | 53            | Yes                |
| Activation     | 50            | Yes                |
| Pooling        | 2             | Yes                |
| elemwise_add   | 16            | Yes                |
| Flatten        | 1             | Yes                |
| FullyConnected | 1             | Yes                |
| SoftmaxOutput  | 1             | No                 |

```

* Total inference operations: 175
* Total Neuron supported inference operations: 174
* Percent of total inference operations supported by Neuron: 99.4
```

4. You can also check the compiled model to see the number of pipeline NeuronCores for each subgraph:

```
$ python mx_neuron_check_model.py resnet-50_compiled

* Found 1 Neuron subgraph(s) (_neuron_subgraph_op(s)) in this compiled model.
* Use this tool on the original uncompiled model to see Neuron supported operations.
* The following table shows all operations, including Neuron subgraphs.
* Each line shows an operation type, the number of instances of that type within model,
* and whether the type is supported in Neuron.
* Some operation types are excluded from table because they are no-operations or
  ↳ training-related operations:
  ['null']
```

| Op Type             | Num Instances | Neuron Supported ? |
|---------------------|---------------|--------------------|
| _neuron_subgraph_op | 1             | Yes                |
| SoftmaxOutput       | 1             | No                 |

(continues on next page)

(continued from previous page)

```
* Please run this model on Inf1 instance with at least 1 NeuronCore(s).
* The following list show each Neuron subgraph with number of pipelined NeuronCores used,
↳by subgraph
* (and subgraph operations if --expand_subgraph is used):
```

| Subgraph Name        | Num Pipelined NeuronCores |
|----------------------|---------------------------|
| -----                | -----                     |
| _neuron_subgraph_op0 | 1                         |

5. When showing subgraph information, you can use `--expand_subgraph` to show operation types in each subgraph:

```
$ python mx_neuron_check_model.py resnet-50_compiled --expand_subgraph
```

(output truncated to show subgraph information only)

| Subgraph Name        | Num Pipelined NeuronCores |
|----------------------|---------------------------|
| -----                | -----                     |
| _neuron_subgraph_op0 | 1                         |
| Op Type              | Num Instances             |
| -----                | -----                     |
| BatchNorm            | 51                        |
| Convolution          | 53                        |
| Activation           | 50                        |
| Pooling              | 2                         |
| elemwise_add         | 16                        |
| Flatten              | 1                         |
| FullyConnected       | 1                         |

6. Use `--show_names` to see full operation names (caution: this option will generate many lines of output for a large model):

```
$ python mx_neuron_check_model.py resnet-50_compiled --show_names
```

```
* Found 1 Neuron subgraph(s) (_neuron_subgraph_op(s)) in this compiled model.
* Use this tool on the original uncompiled model to see Neuron supported operations.
* The following table shows all operations, including Neuron subgraphs.
* Each line shows an operation name and whether the type of that operation is supported,
↳in Neuron.
* Some operation types are excluded from table because they are no-operations or
↳training-related operations:
['null']
```

| Op Name              | Op Type             | Neuron Supported ? |
|----------------------|---------------------|--------------------|
| -----                | -----               | -----              |
| _neuron_subgraph_op0 | _neuron_subgraph_op | Yes                |
| softmax              | SoftmaxOutput       | No                 |

```
* Please run this model on Inf1 instance with at least 1 NeuronCore(s).
* The following list show each Neuron subgraph with number of pipelined NeuronCores used,
↳by subgraph
* (and subgraph operations if --expand_subgraph is used):
```

(continues on next page)

(continued from previous page)

| Subgraph Name        | Num Pipelined NeuronCores |
|----------------------|---------------------------|
| -----                | -----                     |
| _neuron_subgraph_op0 | 1                         |

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## 6.5.2 Using Neuron GatherInfo Tool to collect debug and support information

### Overview

The Neuron GatherInfo tool `neuron-gatherinfo.py` can assist in automating the collection and packaging of information from Neuron SDK tools that is useful to both user and AWS for issue resolution. The tool gathers log files and other system information. If being used to supply that info to AWS, the tool will redact proprietary and confidential information. The GatherInfo tool is supplied in source code form - available here: [Neuron Gatherinfo](#)

The tool enables developers to gather compiler and inference/runtime logs. Additionally, the common usage is from within one of the supported ML frameworks that have been integrated with Neuron, and information can be captured from those compile/runtime environments using the frameworks.

### Steps Overview:

1. Obtain a copy of `neuron-gatherinfo.py` from [Neuron Gatherinfo](#)
2. Install into a location in your \$PATH or into a location from where you can launch the script
3. Use with compile and/or runtime environments

### Neuron-CC information gathering

#### Step 1: Re-run the compile steps for your workload with increased verbosity or debug levels

- For TensorFlow-Neuron, change the Python code as shown. Note that 'compiler-workdir' is expected to be an empty directory to prevent files from other runs from interfering with the information gathering. The call to the compile function has to be augmented with the **verbose** and the **\*\*compiler\_workdir\*\*** arguments. In addition, please capture the stdout messages into a file (for example, by redirecting the stdout to a file)

```
tfn.saved_model.compile(model_dir, compiled_model_dir, compiler_args=['--verbose', '2',
↪ '--pipeline', 'compile', 'SaveTemps'], compiler_workdir='./compiler-workdir')
```

- For Neuron Apache MXNet, add compiler arguments as shown below and run the compilation process from an empty workdir:

```
import mxnet as mx
import os

from packaging import version
mxnet_version = version.parse(mx.__version__)
if mxnet_version >= version.parse("1.8"):
    import mx_neuron as neuron
```

(continues on next page)

(continued from previous page)

```

else:
    from mxnet.contrib import neuron

...
os.environ['SUBGRAPH_INFO'] = '1'
compile_args = { '--verbose' : 2, '--pipeline' : 'compile', 'flags' : ['SaveTemps'] }
csym, cargs, cauxs = neuron.compile(sym, args, auxs, inputs=inputs, **compile_args)

```

## Step 2: Run neuron-gatherinfo.py to gather information to share

The output result will be a tar.gz file.

## Neuron Runtime information gathering

### Step 1: EXECUTE inference steps for your workload with increased verbosity or debug levels

In the case of runtime information, the tool **neuron-dump.py** is used by **neuron-gatherinfo.py** to gather that information. Make sure that you have the neuron tools package (aws-neuron-tools) installed.

## Step 2: Run neuron-gatherinfo.py to gather information to share

The output result will be a tar.gz file.

## Tool Usage Reference

Run neuron-gatherinfo.py using the “—help” option:

```

bash $ ~/bin/neuron-gatherinfo.py --help
usage: neuron-gatherinfo.py [-h] [--additionalfileordir ADDFLDIR] [-c CCDIR]
                          [-i] [-f FILTERFILE] [-m] -o OUTDIR [-r RTDIR] -s
                          STDOUT [-v]

Usage: /home/user/bin/neuron-gatherinfo.py [options]
This program is used to gather information from this system for analysis
and debugging

optional arguments:
  -h, --help                show this help message and exit
  --additionalfileordir ADDFLDIR
                          Additional file or directory that the user wants to
                          provide in the archive. The user can sanitize this
                          file or directory before sharing
  -c CCDIR, --compileroutdir CCDIR
                          Location of the neuron-cc generated files
  -i, --include              By default, only the lines containing (grep) patterns
                          like 'nrtd|neuron|kernel:' from the syslog are copied.
                          Other lines are excluded. Using this option allows the

```

(continues on next page)

(continued from previous page)

```

timestamp section of other lines to be included. The
rest of the contents of the line itself are elided.
Providing the timestamp section may provide time
continuity while viewing the copied syslog file
-f FILTERFILE, --filter FILTERFILE
-m, --modeldata      By using this option, the entire compiler work
                      directory's contents will be included (excluding the
                      .pb files, unless an additional option is used). This
                      would include model information, etc. The files that
                      are included, by default, are these: graph_def.neuron-
                      cc.log, all_metrics.csv, hh-tr-operand-
                      tensortensor.json
-o OUTDIR, --out OUTDIR
                      The output directory where all the files and other
                      information will be stored. The output will be stored
                      as an archive as well as the actual directory where
                      all the contents are copied. This will allow a simple
                      audit of the files, if necessary. *** N O T E ***:
                      Make sure that this directory has enough space to hold
                      the files and resulting archive
-r RTDIR, --runtimeoutdir RTDIR
                      Location of the neuron runtime generated files
-s STDOUT, --stdout STDOUT
                      The file where the stdout of the compiler run was
                      saved
-v, --verbose         Verbose mode displays commands executed and any
                      additional information which may be useful in
                      debugging the tool itself

```

## Examples

### Example 1: no ML model information gathered (default behavior)

In this case, the tool will archive just the default information gathering:

```

bash $ sudo ~/bin/neuron-gatherinfo.py -o compile-and-run-info-for-debugging-no-model-
↳info -i --verbose -s stdout-from-compile_resnet50.out -c compiler-workdir

Running cmd: lscpu and capturing output in file: /home/user/tutorials-3/compile-and-run-
↳info-for-debugging-no-model-info/neuron-gatherinfo/report-lscpu.txt
Running cmd: lshw and capturing output in file: /home/user/tutorials-3/compile-and-run-
↳info-for-debugging-no-model-info/neuron-gatherinfo/report-lshw.txt
Running cmd: lspci | grep -i Amazon and capturing output in file: /home/user/tutorials-3/
↳compile-and-run-info-for-debugging-no-model-info/neuron-gatherinfo/report-lspci.txt
Running cmd: neuron-cc --version and capturing output in file: /home/user/tutorials-3/
↳compile-and-run-info-for-debugging-no-model-info/neuron-gatherinfo/report-neuron-cc.txt
Running cmd: neuron-ls and capturing output in file: /home/user/tutorials-3/compile-and-
↳run-info-for-debugging-no-model-info/neuron-gatherinfo/report-neuron-ls.txt
<SNIP>
*****
Archive created at:

```

(continues on next page)

(continued from previous page)

```

/home/user/tutorials-3/compile-and-run-info-for-debugging-no-model-info/neuron-
↳gatherinfo.tar.gz
From directory:
/home/user/tutorials-3/compile-and-run-info-for-debugging-no-model-info/neuron-
↳gatherinfo
*****

```

## Example 2 : model ML information gathered using the “—modeldata” option

In this case, the tool will archive the compiler work directory in addition to the default information gathering

```

bash $ sudo ~/bin/neuron-gatherinfo.py -o compile-and-run-info-for-debugging -i --
↳verbose -s stdout-from-compile_resnet50.out -c compiler-workdir --modeldata

<SNIP>
Running cmd: lscpu and capturing output in file: /home/user/tutorials-3/compile-and-run-
↳info-for-debugging/neuron-gatherinfo/report-lscpu.txt
Running cmd: lshw and capturing output in file: /home/user/tutorials-3/compile-and-run-
↳info-for-debugging/neuron-gatherinfo/report-lshw.txt
Running cmd: lspci | grep -i Amazon and capturing output in file: /home/user/tutorials-3/
↳compile-and-run-info-for-debugging/neuron-gatherinfo/report-lspci.txt
Running cmd: neuron-cc --version and capturing output in file: /home/user/tutorials-3/
↳compile-and-run-info-for-debugging-no-model-info/neuron-gatherinfo/report-neuron-cc.txt
Running cmd: neuron-ls and capturing output in file: /home/user/tutorials-3/compile-and-
↳run-info-for-debugging-no-model-info/neuron-gatherinfo/report-neuron-ls.txt
<SNIP>

*****
Archive created at:
/home/user/tutorials-3/compile-and-run-info-for-debugging/neuron-gatherinfo.tar.
↳gz
From directory:
/home/user/tutorials-3/compile-and-run-info-for-debugging/neuron-gatherinfo
*****

*****
Based on your command line option, we're also packaging these files:

graph_def.neuron-cc.log
all_metrics.csv
hh-tr-operand-tensortensor.json

And this directory: /home/user/tutorials-3/compiler-workdir

*****

```

*This document is relevant for: Inf1*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 6.5.3 NeuronPerf Overview

NeuronPerf is a lightweight Python library that can help you easily benchmark your models with Neuron hardware.

NeuronPerf supports Neuron releases for PyTorch, Tensorflow, and MXNet. It is used internally by the Neuron team to generate performance benchmarking numbers.

When interacting with NeuronPerf, you will typically import the base package along with one of the submodule wrappers, for example:

```
import neuronperf
import neuronperf.torch
```

You may then benchmark and/or compile one or more models with NeuronPerf. For example,

```
reports = neuronperf.torch.benchmark(model, inputs, ...)
```

The compile and benchmark methods must be accessed through one of the supported framework submodules.

#### Benchmarking

All NeuronPerf benchmark calls require a minimum of two arguments:

1. A filename
2. Inputs

The filename may refer to:

1. A Neuron-compiled model (e.g. `my_model.pt`)
2. A Model Index.

A Model Index is useful for benchmarking more than one model in a single session.

#### Compiling

NeuronPerf also provides a standard interface to all Neuron frameworks through the `compile` API.

```
model_index = neuronperf.torch.compile(model, inputs, ...)
```

This is completely optional. You may use the standard compilation guides for supported frameworks.

#### Next Steps

Take a look at the simple `neuronperf_examples`, `neuronperf_benchmark_guide`, `neuronperf_compile_guide`, and `neuronperf_api`.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



## 6.5.4 Neuron Calculator

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 6.5.5 Neuron Plugin for TensorBoard (Trn1)

### Table of Contents

- [Overview](#)
- [Enable profiling on Trn1](#)
- [Launch TensorBoard](#)
- [View results in TensorBoard](#)
- [Neuron Trace View](#)
- [Neuron Operator View](#)
- [Neuron Operator Timeline View](#)
- [Troubleshooting](#)
  - [TensorBoard launch fails](#)

### Overview

This guide is for developers who want to better understand how their model is executed using Neuron SDK through TensorBoard.

The Neuron plugin for TensorBoard provides metrics to the performance of machine learning tasks accelerated using the Neuron SDK. It is compatible with TensorBoard versions 1.15 and higher. It provides visualizations and profiling results for graphs executed on NeuronCores.

---

**Note:** The following information is compatible with Neuron SDK for Trn1. For a walkthrough on Inf1, please check out the guide [Neuron Plugin for TensorBoard \(Inf1\)](#).

---

### Enable profiling on Trn1

---

**Note:** Profiling is currently only supported with PyTorch Neuron (`torch-neuronx`).

---

Please refer to the following guides:

- **PyTorch-Neuron**
  - [torch-neuronx-profiling-with-tb](#)

## Launch TensorBoard

In this step, we will process the Neuron profile data and launch TensorBoard.

1. Install the Neuron plugin for Tensorboard on your EC2 instance.

```
python -m pip config set global.extra-index-url "https://pip.repos.neuron.amazonaws.com"
pip install tensorboard-plugin-neuronx
```

**Note:** If using TensorBoard >= 2.5, please use the `--load_fast=false` option when launching. `tensorboard --logdir results --load_fast=false`

2. After you see the following message, TensorBoard is ready to use. By default, TensorBoard will be launched at `localhost:6006`.

```
...
Serving TensorBoard on localhost; to expose to the network, use a proxy or pass --bind_
all
TensorBoard 2.4.1 at http://localhost:6006/ (Press CTRL+C to quit)
```

## View results in TensorBoard

In this step, we will view the Neuron plugin for TensorBoard from a browser on your local development machine.

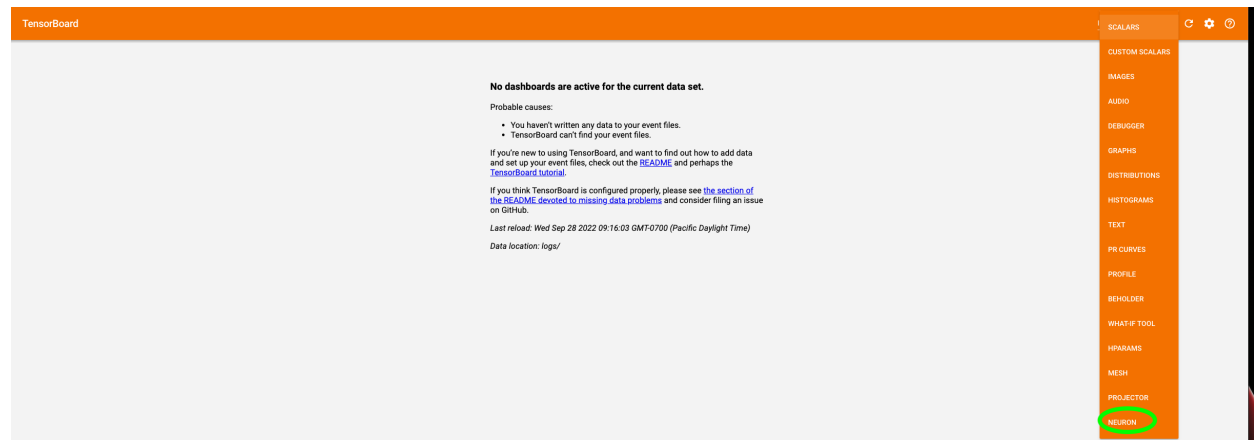
1. Connect to the EC2 instance where TensorBoard is running while enabling port forwarding. In this example, we assume TensorBoard has been launched using the default address `localhost:6006`.

```
# if Ubuntu-based AMI
ssh -i <PEM key file> ubuntu@<instance DNS> -L 6006:localhost:6006

# if AL2-based AMI
ssh -i <PEM key file> ec2-user@<instance DNS> -L 6006:localhost:6006
```

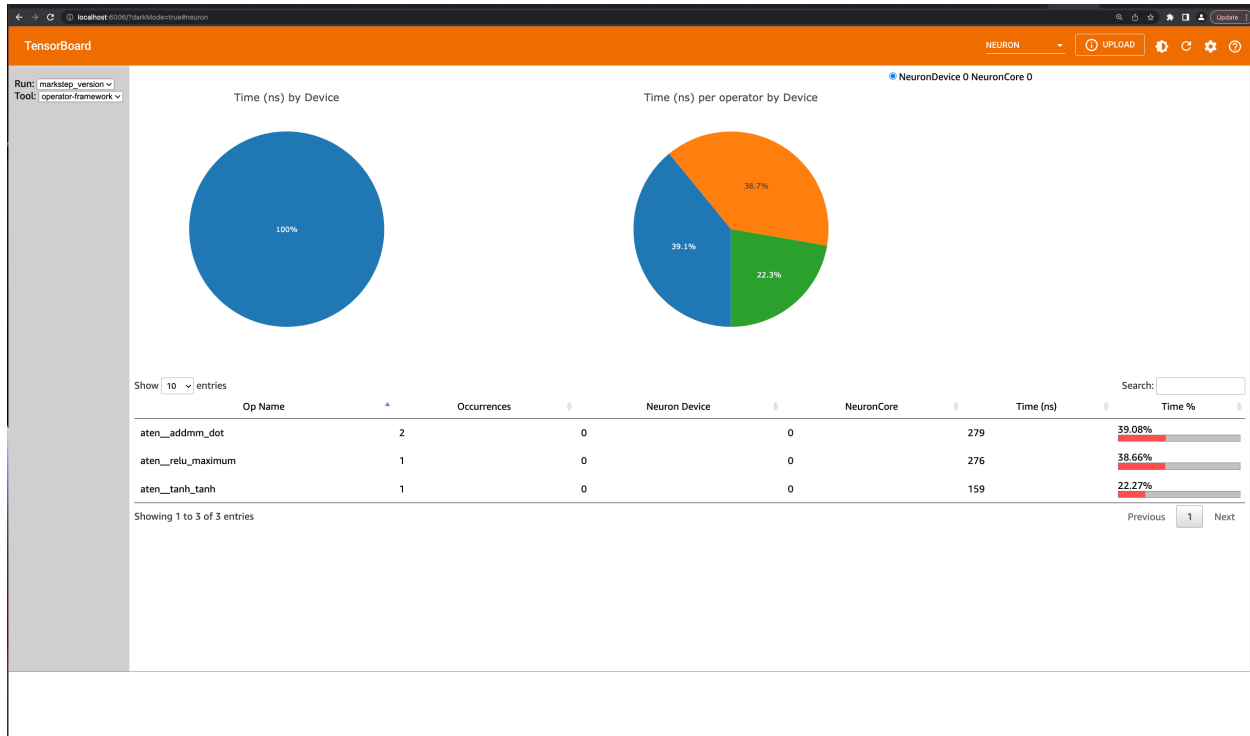
2. In a browser, visit .

3. In the top navigation bar, switch from **Graphs** to **Neuron**. If it does not show up, please wait a while and refresh the page while the plugin loads. If the issue persists, check the Inactive dropdown list on the right and check for Neuron.





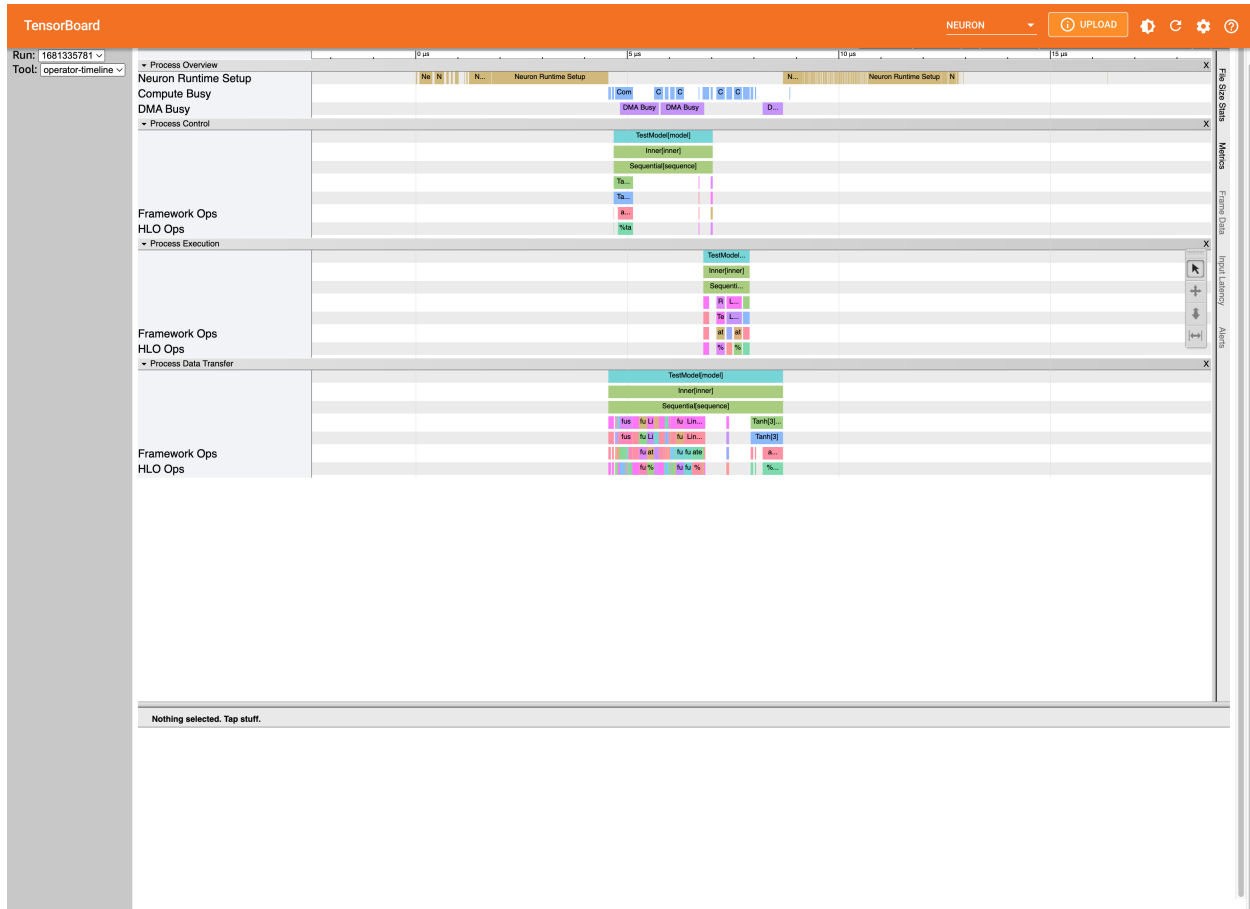
## Neuron Operator View



The operator view can show timing information for both the framework operators and HLO operators by selecting the `operator-framework` and `operator-hlo` tools respectively. The pie charts show breakdowns of the time taken by device, as well as per operator on a single device. The table below lists out the operators and can be sorted by clicking on the column headers. For fused operations, hover over the ? to see which operators are being executed.

For a quick glance at the most time consuming operators, click the `Time %` column in the table to sort by the relative time spent on this type of operation compared to the rest of the model.

## Neuron Operator Timeline View



The operator timeline view is a detailed look into a single execution with Neuron. A high level overview at the top breaks down the execution into categories, including Neuron Runtime setup time, as well as NeuronCore compute engine and DMA engine busyness. Activity on the compute and DMA engines are further categorized as compute, control, and data transfer intervals which are shown as separate processes, with each showing a hierarchical view of the framework operators and their corresponding HLO operation. The fused operations can be a result of compiler optimizations or are operations that are running in parallel on the device. Each bar can be clicked to show information regarding which operators are overlapped.

This view can give better insight into how operators translate to Neuron, as well as how certain Neuron compiler options may improve performance.

## Troubleshooting

### TensorBoard launch fails

```
ImportError: cannot import name 'Mapping' from 'collections'
```

This is an issue with Python 3.10 and a dependency of an old tensorboard version. To workaround this error, please run `pip install --upgrade tensorboard`. For more information, see <https://github.com/tensorflow/tensorboard/pull/5490>.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1*

## 6.5.6 Neuron Plugin for TensorBoard (Inf1)

### Table of Contents

- *Overview*
- *Compile the neural network*
- *Enable profiling*
- *Launch TensorBoard*
- *View results in TensorBoard*
- *Visualize graphs executed on Neuron*
  - *Show how the graph was partition to run on NeuronCores*
  - *Inspect which operators consumes the most time*
  - *Check out Neuron support operators for each framework*
  - *Filter view by device*
  - *Expand/collapse subgraphs and view operator details*
- *Viewing the Neuron profile data*
  - *See performance summary*
  - *Get a breakdown of time spent per NeuronCore*
  - *Get a breakdown of time spent per operator*

### Overview

This guide is for developers who want to better understand how their model is executed using Neuron SDK through TensorBoard.

The Neuron plugin for TensorBoard provides metrics to the performance of machine learning tasks accelerated using the Neuron SDK. It is compatible with TensorBoard versions 1.15 and higher. It provides visualizations and profiling results for graphs executed on NeuronCores.

---

**Note:** The following information is compatible with Neuron SDK for Inf1. For a walkthrough on the latest version, please check out the guide *Neuron Plugin for TensorBoard (Trn1)*.

---

---

**Note:** Graph visualization is currently only supported for TensorFlow-Neuron. Support for MXNet-Neuron and PyTorch-Neuron visualization will be added in a future release.

---

## Compile the neural network

3. Refer to the following guides on how to compile a graph using Neuron SDK.

- **TensorFlow-Neuron**
  - /src/examples/tensorflow/tensorflow\_resnet50/resnet50.ipynb
- **PyTorch-Neuron:**
  - “Compile model for Neuron” in PyTorch-Neuron Resnet50 Tutorial
- **MXNet-Neuron:**
  - /src/examples/mxnet/resnet50/resnet50.ipynb

## Enable profiling

In this step, we enable Neuron profile data collection and collect results from executing an inference.

4.1. To start profiling the neural network and collect inference traces, create a directory where profile data will be dumped and set the NEURON\_PROFILE environment variable. In this example, we will assume this directory is \$HOME/profile

```
mkdir -p $HOME/profile
export NEURON_PROFILE=$HOME/profile
```

4.2. Ensure Neuron Tools are executable by setting the PATH environment variable.

```
export PATH=/opt/aws/neuron/bin:$PATH
```

4.3. Execute inference!

---

**Note:** Please run the inference script outside of Jupyter notebook. Profiling in Jupyter notebook is not supported at this time.

---



---

**Note:** Please ensure the inference script executes only one inference, as profiling results are currently only supported for a single inference.

---

For more info on how to execute inference, refer to the following guides:

- **TensorFlow-Neuron**
  - /src/examples/tensorflow/tensorflow\_resnet50/resnet50.ipynb
- **PyTorch-Neuron**
  - “Run inference on Single Core” in /src/examples/pytorch/resnet50.ipynb
- **MXNet-Neuron**
  - /src/examples/mxnet/resnet50/resnet50.ipynb

4.4. Check if profiling results were successfully saved. In the directory pointed to by NEURON\_PROFILE environment variable set in Step 4.1, there should be at least two files, one with the .neff extension and one with the .ntff extension. For TensorFlow-Neuron users, the graph file (.pb) will also be in this directory.

```
ls $NEURON_PROFILE
```

## Launch TensorBoard

In this step, we will process the Neuron profile data and launch TensorBoard.

### 5.1. Install the Neuron plugin for Tensorboard.

If you are using the DLAMI TensorFlow-Neuron Conda environment, please run the following to update TensorBoard before installing the Neuron plugin.

```
pip install "tensorboard<=2.4.0" --force-reinstall
```

Modify Pip repository configurations to point to the Neuron repository:

```
tee $VIRTUAL_ENV/pip.conf > /dev/null <<EOF
[global]
extra-index-url = https://pip.repos.neuron.amazonaws.com
EOF
```

```
pip install tensorboard-plugin-neuron
```

5.2. After collecting the raw profile data, we need to post-process it to create the log files used by the Neuron plugin. This can be done when launching TensorBoard by passing an extra flag `--run_neuron_profiler`. Using this flag will create the directory specified by `--logdir` and populate it with Neuron plugin data. Please note that the `NEURON_PROFILE` environment variable set in Step 4.1 must still point to the same directory as before.

```
tensorboard --logdir results --run_neuron_profiler
```

---

**Note:** If using TensorBoard `>= 2.5`, please use the `--load_fast=false` option when launching. `tensorboard --logdir results --run_neuron_profiler --load_fast=false`

---

5.3. After you see the following message, TensorBoard is ready to use. By default, TensorBoard will be launched at `localhost:6006` on the Deployment Instance.

```
...
Running neuron-profile
Serving TensorBoard on localhost; to expose to the network, use a proxy or pass --bind_
↪all
TensorBoard 2.4.1 at http://localhost:6006/ (Press CTRL+C to quit)
```

## View results in TensorBoard

In this step, we will view the Neuron plugin for TensorBoard from a browser on your local development machine.

6.1. Connect to the Deployment Instance while enabling port forwarding. In this example, we assume TensorBoard has been launched using the default address `localhost:6006` on the Deployment Instance.

```
# if Ubuntu-based AMI
ssh -i <PEM key file> ubuntu@<instance DNS> -L 6006:localhost:6006
```

(continues on next page)



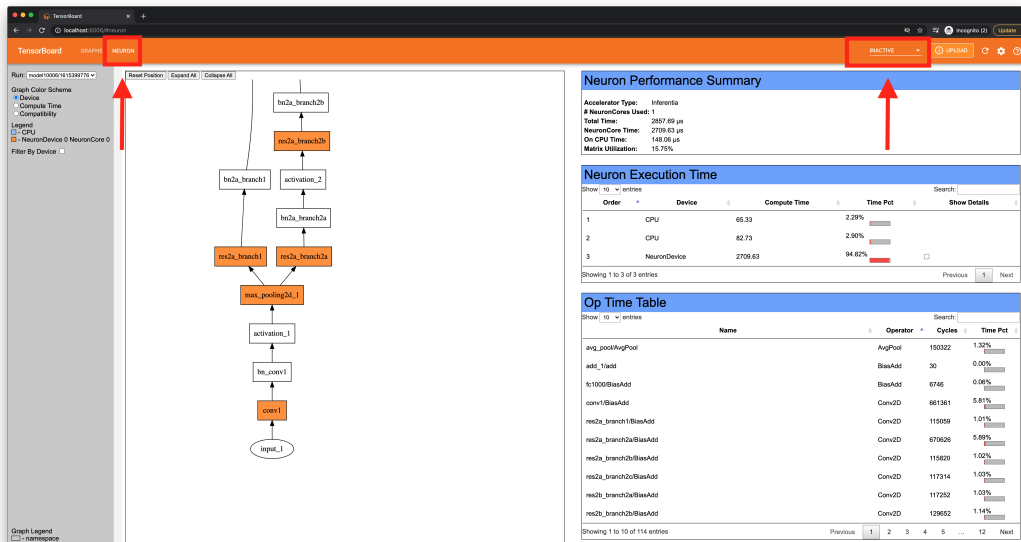
(continued from previous page)

```
# if AL2-based AMI
```

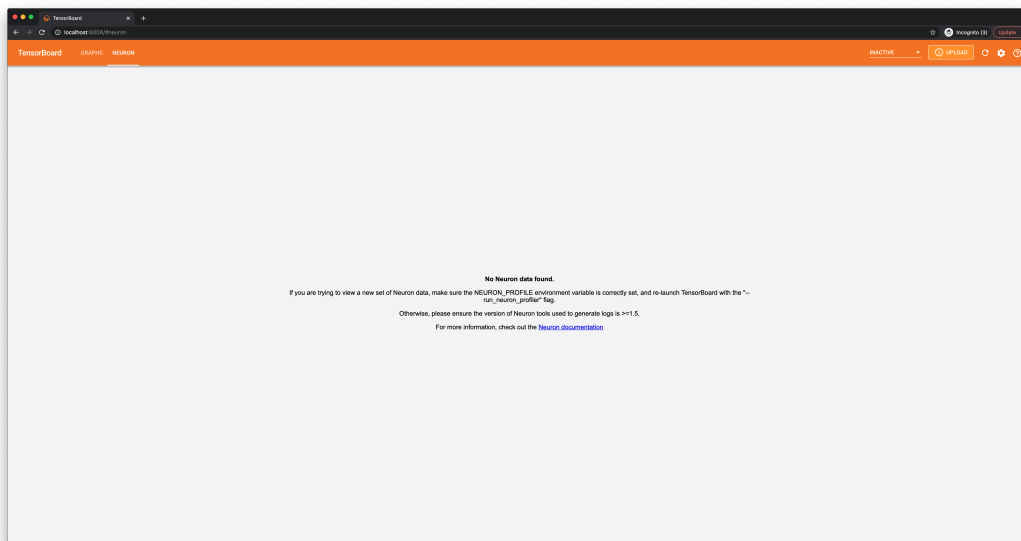
```
ssh -i <PEM key file> ec2-user@<instance DNS> -L 6006:localhost:6006
```

6.2. In a browser, visit .

6.3. In the top navigation bar, switch from Graphs to Neuron. If it does not show up, please wait a while and refresh the page while the plugin loads. If the issue persists, check the Inactive dropdown list on the right and check for Neuron.



6.4. If TensorBoard failed to find the generated logs, you will see the following message:

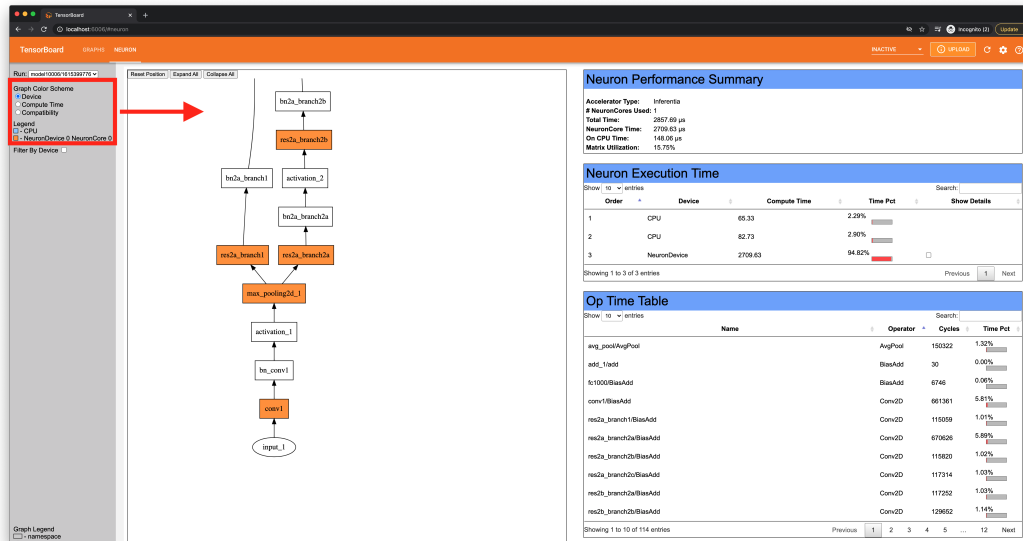


In this case, please check the console output on the Deployment Instance where TensorBoard was launched for any warnings or error messages, and make sure the version of the `aws-neuron-tools` package is compatible.

## Visualize graphs executed on Neuron

### Show how the graph was partitioned to run on NeuronCores

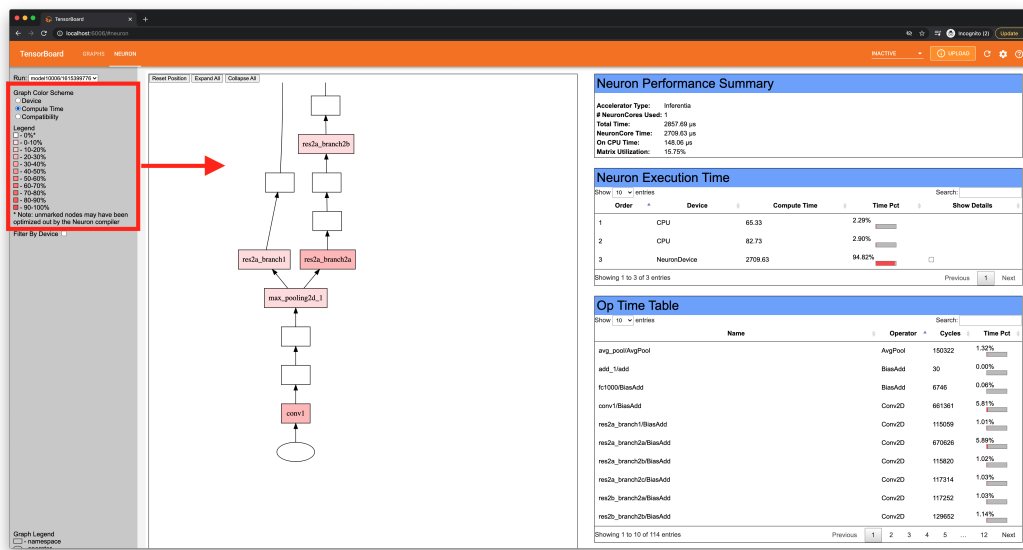
To view how the graph was partitioned to run on NeuronCores, select “Device” under “Graph Color Schemes” in the left navigation bar.



Each operator will be colored according to the device used. In this example, light blue indicates an operator was executed on CPU, and orange indicates the operator was executed on NeuronCores. Operators that are white may have been optimized by the Neuron compiler and fused into another operation.

### Inspect which operators consumes the most time

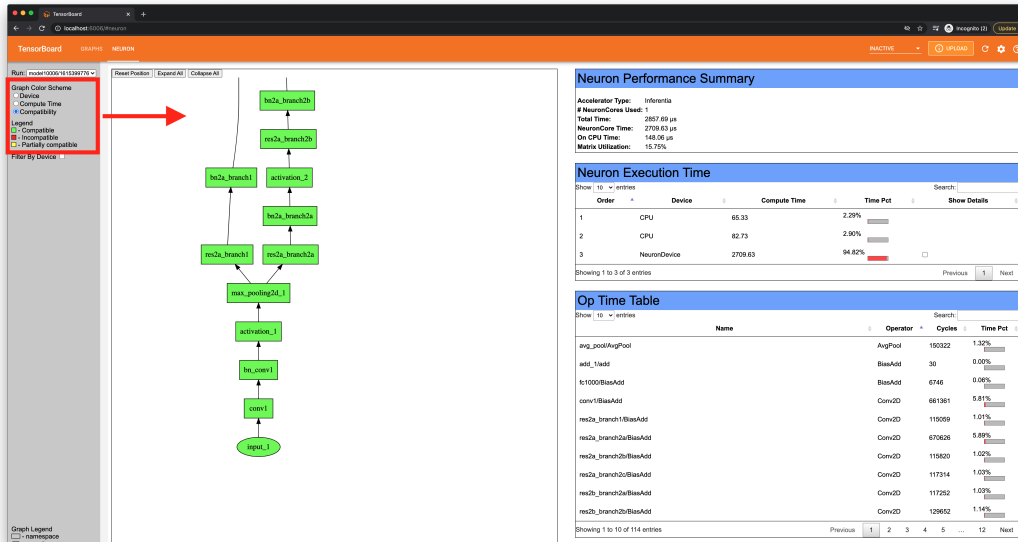
You can also view how long each operator took by changing to the “Compute time” color scheme.



This view will show time taken by each layer and will be colored according to how much relative time the layer took to compute. A lighter shade of red means that a relatively small portion of compute time was spent in this layer, while a darker red shows that more compute time was used.

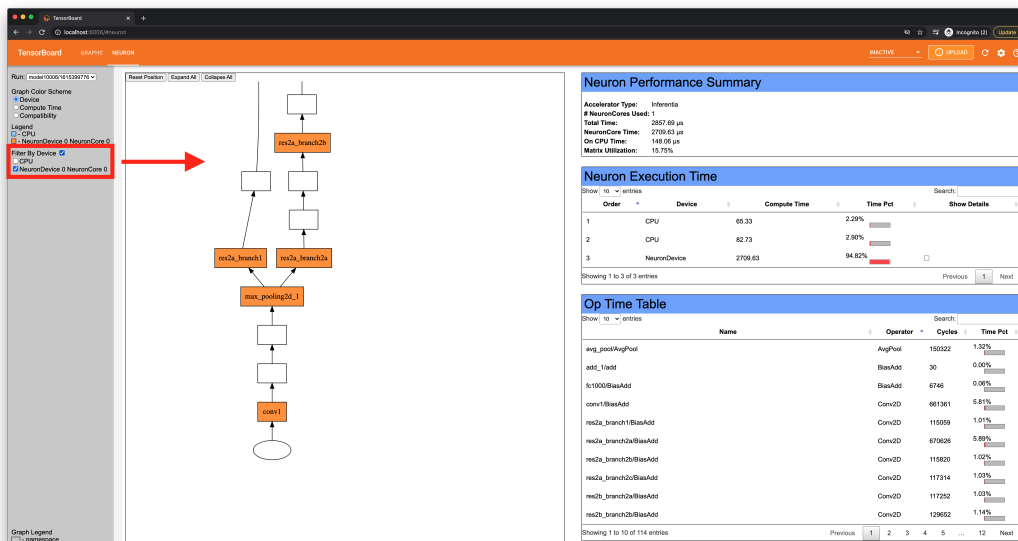
## Check out Neuron support operators for each framework

The “Compatibility” color scheme allows you to better understand what operators are currently supported by the Neuron compiler - green for compatible ops, red for incompatible ops, and yellow for subgraphs that contain both compatible and incompatible ops.

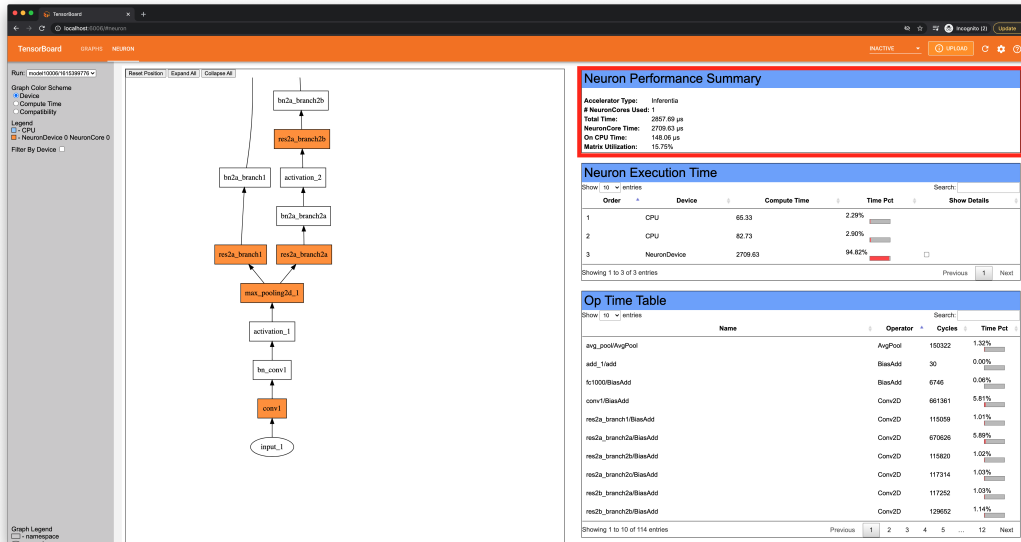


## Filter view by device

Additionally, you can choose to filter by CPU and NeuronCores, which will only color ops that match the selected device(s).

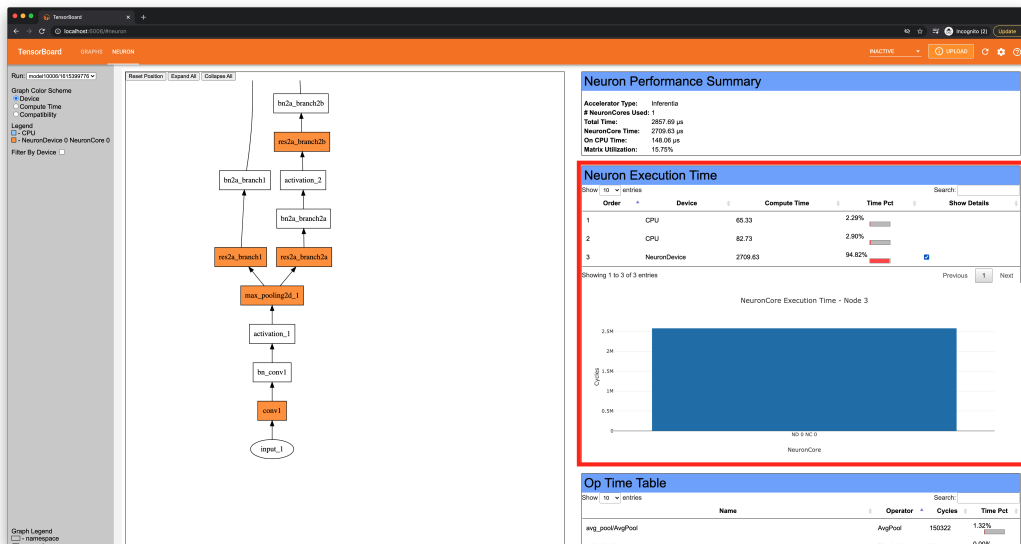






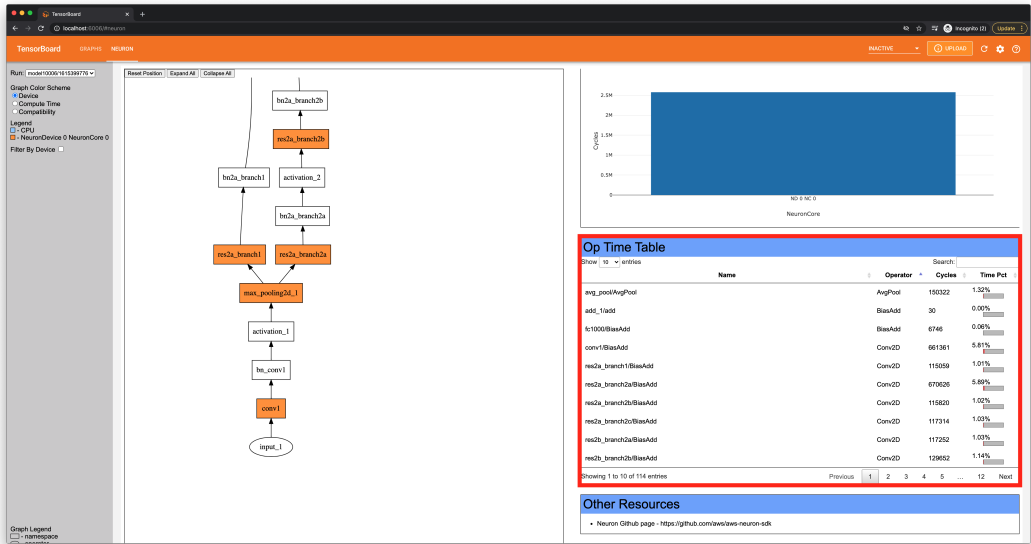
## Get a breakdown of time spent per NeuronCore

Next, the “Neuron Execution” will give more details on how a graph was partitioned for Neuron. Each entry in the table will show the order it was executed in, what type of device was used, the compute time (in microseconds), and the percentage of total time spent. To dive deeper into subgraphs, you can check the “Show Details” box to display the breakdown per NeuronCore.



Get a breakdown of time spent per operator

The “Op Time Table” section shows the cycle count per operator, much like the “Compute time” coloring for graph visualization. This table can be sorted by clicking the column names, and searched using the provided text box in the top right corner. Due to Neuron compiler optimizations, some of the compute may not be associated with any specific operator and will be categorized as unknown. Additionally, time spent moving data to and from NeuronCores will fall under (ND\_ENGINE\_LOAD).



This document is relevant for: Inf1  
This document is relevant for: Inf2, Trn1, Trn2

6.5.7 Track Training Progress in TensorBoard using PyTorch Neuron

Table of Contents

- Multi-layer perceptron MNIST model
- Output TensorBoard logs
- View loss in TensorBoard

This tutorial explains how to track training progress in TensorBoard while running a multi-layer perceptron MNIST model on Trainium using PyTorch Neuron.

## Multi-layer perceptron MNIST model

This tutorial is based on the MNIST example for PyTorch Neuron on Trainium. For the full tutorial, please see [Multi-Layer Perceptron Training Tutorial](#).

### Output TensorBoard logs

To generate TensorBoard logs, we first modify the training script to use the SummaryWriter:

```
from torch.utils.tensorboard import SummaryWriter
writer = SummaryWriter('./output')
```

In the training loop, we can then use the add\_scalar API to log the loss per step.

```
writer.add_scalar("step loss", loss, idx)
```

At the end of the script, add `writer.flush()` to ensure all logs are written.

Save the following code as `train_tb.py` and run it as `python3 train_tb.py` on a Trn1 instance. The generated logs can be found in the `./output` directory that was passed to SummaryWriter.

```
import os
import time
import torch
import torch.nn as nn
import torch.nn.functional as F

from torchvision.datasets import mnist
from torch.optim import SGD
from torch.utils.data import DataLoader
from torchvision.transforms import ToTensor

# XLA imports
import torch_xla.core.xla_model as xm

from torch.utils.tensorboard import SummaryWriter

# Declare 3-layer MLP for MNIST dataset
class MLP(nn.Module):
    def __init__(self, input_size = 28 * 28, output_size = 10, layers = [120, 84]):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return F.log_softmax(x, dim=1)

# Load MNIST train dataset
train_dataset = mnist.MNIST(root='./MNIST_DATA_train', \
                             train=True, download=True, transform=ToTensor())
```

(continues on next page)

(continued from previous page)

```

def main():
    # Prepare data loader
    train_loader = DataLoader(train_dataset, batch_size=32)

    # Fix the random number generator seeds for reproducibility
    torch.manual_seed(0)

    # XLA: Specify XLA device (defaults to a NeuronCore on Trn1 instance)
    device = 'xla'

    # Move model to device and declare optimizer and loss function
    model = MLP().to(device)
    optimizer = torch.optim.SGD(model.parameters(), lr=0.01)
    loss_fn = torch.nn.NLLLoss()

    # Use SummaryWriter to generate logs for TensorBoard
    writer = SummaryWriter('./output')

    # Run the training loop
    print('-----Training -----')
    model.train()
    start = time.time()
    for idx, (train_x, train_label) in enumerate(train_loader):
        optimizer.zero_grad()
        train_x = train_x.view(train_x.size(0), -1)
        train_x = train_x.to(device)
        train_label = train_label.to(device)
        output = model(train_x)
        loss = loss_fn(output, train_label)
        writer.add_scalar("step loss", loss, idx) # add the step loss to the TensorBoard_
↪logs
        loss.backward()
        optimizer.step()
        xm.mark_step() # XLA: collect ops and run them in XLA runtime
        if idx < 2: # skip warmup iterations
            start = time.time()

    # Compute statistics
    interval = idx - 2 # skip warmup iterations
    throughput = interval / (time.time() - start)
    print("Train throughput (iter/sec): {}".format(throughput))
    print("Final loss is {:.4f}".format(loss.detach().to('cpu')))

    # Ensure TensorBoard logs are all written
    writer.flush()

    # Save checkpoint for evaluation
    os.makedirs("checkpoints", exist_ok=True)
    checkpoint = {'state_dict': model.state_dict()}
    # XLA: use xm.save instead of torch.save to ensure states are moved back to cpu
    # This can prevent "XRT memory handle not found" at end of test.py execution

```

(continues on next page)



(continued from previous page)

```

xm.save(checkpoint, 'checkpoints/checkpoint.pt')

print('-----End Training -----')

if __name__ == '__main__':
    main()

```

## View loss in TensorBoard

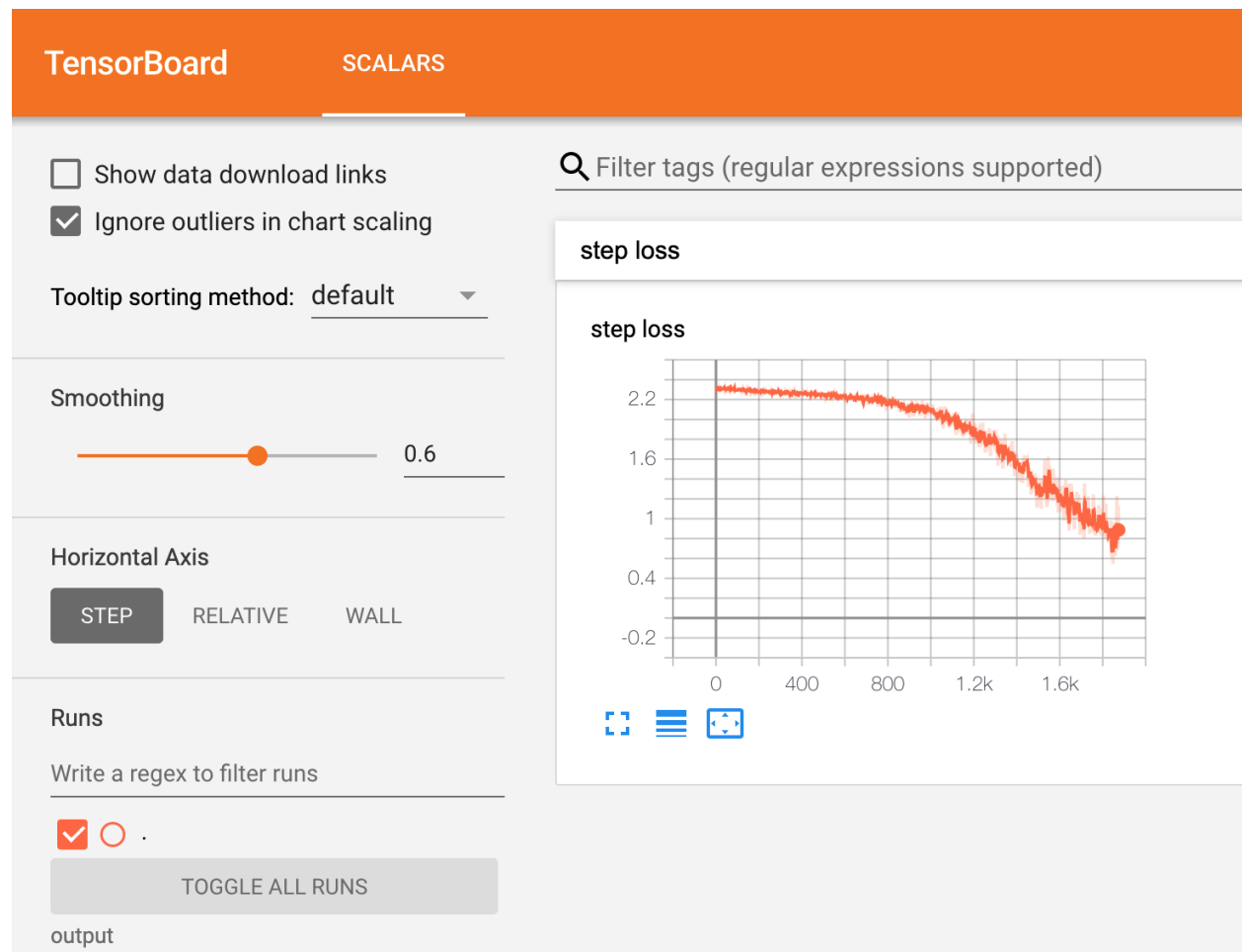
In order to view your training metrics, install TensorBoard in your Python environment:

```
pip install tensorboard
```

Then, launch TensorBoard with the ./output directory

```
tensorboard --logdir ./output
```

Once running, open a new SSH connection to the instance and port-forward TCP port 6006 (ex: -L 6006:127.0.0.1:6006). Once the tunnel is established, TensorBoard can then be accessed via web browser at the following URL: <http://localhost:6006>. Please note that you will not be able to access TensorBoard if you disconnect your port-forwarding SSH session to the Trainium instance.



In TensorBoard, you can now see the loss per step plotted. When capturing loss for multiple runs, you can plot them together on the same graph to compare runs. Be sure to change the output directory for different runs, for example `./output/run1` for the first, `./output/run2` for the second, etc.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 6.5.8 Neuron Plugin for TensorBoard Release Notes

### Table of Contents

- *Known Issues and Limitations - Updated 11/29/2022*
- *Neuron Plugin for TensorBoard release [2.6.7.0]*
- *Neuron Plugin for TensorBoard release [2.6.1.0]*
- *Neuron Plugin for TensorBoard release [2.5.39.0]*
- *Neuron Plugin for TensorBoard release [2.5.37.0]*
- *Neuron Plugin for TensorBoard release [2.5.26.0]*
- *Neuron Plugin for TensorBoard release [2.5.25.0]*
- *Neuron Plugin for TensorBoard release [2.5.0.0]*
- *Neuron Plugin for TensorBoard release [2.4.0.0]*
- *Neuron Plugin for TensorBoard release [2.3.0.0]*
- *Neuron Plugin for TensorBoard release [2.2.0.0]*
- *[2.1.2.0]*
- *[2.1.0.0]*
- *[2.0.29.0]*
- *[1.15.0.1.2.6.0]*
- *[1.15.0.1.1.1.0]*
- *[1.15.0.1.0.615.0]*
- *[1.15.0.1.0.600.0]*
- *[1.15.0.1.0.570.0]*
- *[1.15.0.1.0.513.0]*
- *[1.15.0.1.0.491.0]*
- *[1.15.0.1.0.466.0]*
- *[1.15.0.1.0.392.0]*
- *[1.15.0.1.0.366.0]*
- *[1.15.0.1.0.315.0]*
- *[1.15.0.1.0.306.0]*
- *[1.15.0.1.0.280.0]*

### Known Issues and Limitations - Updated 11/29/2022

The following are not limitations in the Neuron plugin, but may affect your ability to use TensorBoard.

- The Neuron plugin for Trn1 (`tensorboard-plugin-neuronx`) is not compatible with the Neuron plugin for Inf1 (`tensorboard-plugin-neuron`). Please ensure you only have only the correct package installed.

### Neuron Plugin for TensorBoard release [2.6.7.0]

Date: 04/01/2024

#### Summary

- Minor updates.

### Neuron Plugin for TensorBoard release [2.6.1.0]

Date: 12/21/2023

#### Summary

- Now uses local third-party dependencies instead of relying on a CDN.

### Neuron Plugin for TensorBoard release [2.5.39.0]

Date: 7/19/2023

#### Summary

- Minor updates.

### Neuron Plugin for TensorBoard release [2.5.37.0]

Date: 6/14/2023

#### Summary

- Minor updates.

### Neuron Plugin for TensorBoard release [2.5.26.0]

Date: 05/01/2023

#### Summary

- Neuron operator timeline view now includes Neuron Runtime setup/teardown time and a collapsed execution of NC engines and DMA - see Tensorboard tutorial for updated views.
- Improved execution categorization to include “control” instructions

### Neuron Plugin for TensorBoard release [2.5.25.0]

Date: 03/28/2023

#### Summary

- Supports INF2 and TRN1.

### Neuron Plugin for TensorBoard release [2.5.0.0]

Date: 12/09/2022

#### Summary

- Added support for PyTorch Neuron on Trn1 (`torch-neuronx`) with new views! Includes a trace view, an operator view, and an operator timeline view. For more info, check out the documentation [Neuron Plugin for TensorBoard \(Trn1\)](#).

---

#### Important:

- You must update to the latest Neuron Tools (`aws-neuronx-tools` version 2.6 or newer) and install `tensorboard-plugin-neuronx` for proper functionality of the Neuron plugin on Trn1.
  - For Inf1, please continue to use `tensorboard-plugin-neuron`. Refer to the getting started guide on Inf1 [Neuron Plugin for TensorBoard \(Inf1\)](#).
- 

### Neuron Plugin for TensorBoard release [2.4.0.0]

Date: 04/29/2022

## Summary

- Minor updates.

## Neuron Plugin for TensorBoard release [2.3.0.0]

Date: 03/25/2022

## Summary

- Minor updates.

## Neuron Plugin for TensorBoard release [2.2.0.0]

Date: 10/27/2021

## New in this release

- Neuron Plugin for TensorBoard now support applications built with Neuron Runtime 2.x (`libnrt.so`).

---

### Important:

- You must update to the latest Neuron Driver (`aws-neuron-dkms` version 2.1 or newer) for proper functionality of the new runtime library.
  - Read *Introducing Neuron Runtime 2.x (`libnrt.so`)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.
  - Read *Migrate your application to Neuron Runtime 2.x (`libnrt.so`)* for detailed information of how to migrate your application.
- 

## [2.1.2.0]

Date: 8/12/2021

## Summary

- Adds support for Neuron Tensorflow 2.5+

### [2.1.0.0]

Date: 5/28/2021

#### Summary

- No major changes or fixes. Released with other Neuron packages.

### [2.0.29.0]

Date: 4/30/2021

#### Summary

- First release Neuron plugin for TensorBoard. Check out it out here: [Neuron Plugin for TensorBoard \(Inf1\)](#).
  - The Neuron plugin is now compatible with TensorBoard 2.0 and higher, in addition to TensorBoard 1.15
  - Provides a centralized place to better understand execution using Neuron SDK.
  - Continues support visualization for TensorFlow graphs, with support for PyTorch and MXNet coming in future releases.
- Neuron plugin for TensorBoard is supported for Neuron tools  $\geq 1.5$ , which is first introduced in Neuron v1.13.0 release
- TensorBoard-Neuron is deprecated, and only supported for Neuron tools  $\leq 1.4.12.0$ . The final version, 1.4.12.0 is part of Neuron v1.12.2 release.

### [1.15.0.1.2.6.0]

Date: 2/24/2021

#### Summary

- Fix for CVE-2021-3177.

### [1.15.0.1.1.1.0]

Date: 12/23/2020

#### Summary

- Minor internal improvements.

**[1.15.0.1.0.615.0]**

Date: 11/17/2020

**Summary**

- Fix issue with viewing chrome trace in Neuron profile plugin in Chrome 80+.

**Resolved Issues**

- Updated dependencies to polyfill missing APIs used by chrome trace in newer browser versions.

**[1.15.0.1.0.600.0]**

Date: 09/22/2020

**Summary**

- Minor internal improvements.

**[1.15.0.1.0.570.0]**

Date: 08/08/2020

**Summary**

- Minor internal improvements.

**[1.15.0.1.0.513.0]**

Date: 07/16/2020

**Summary**

- Minor internal improvements.

**[1.15.0.1.0.491.0]**

Date 6/11/2020

### Summary

Fix issue where utilization was missing in the op-profile view.

### Resolved Issues

- The op-profile view in the Neuron Profile plugin now correctly shows the overall NeuronCore utilization.

[1.15.0.1.0.466.0]

Date 5/11/2020

### Summary

Fix potential installation issue when installing both tensorboard and tensorboard-neuron.

### Resolved Issues

- Added tensorboard as a dependency in tensorboard-neuron. This prevents the issue of overwriting tensorboard-neuron features when tensorboard is installed after tensorboard-neuron.

### Other Notes

[1.15.0.1.0.392.0]

Date 3/26/2020

### Summary

Added ability to view CPU node latency in the Graphs plugin and the Neuron Profile plugins.

### Major New Features

- Added an aggregate view in addition to the current Neuron subgraph view for both the Graphs plugin and the Neuron Profile plugin.
- When visualizing a graph executed on a Neuron device, CPU node latencies are available when coloring the graph by “Compute time” using the “neuron\_profile” tag.
- The Neuron Profile plugin now has an overview page to compare time spent on Neuron device versus on CPU.



## Other Notes

- Requires Neuron-RTD config option “enable\_node\_profiling” to be set to “true”

### [1.15.0.1.0.366.0]

Date 02/27/2020

## Summary

Reduced load times and fixed crashes when loading large models for visualization.

## Resolved Issues

- Enable large attribute filtering by default
- Reduced load time for graphs with attributes larger than 1 KB
- Fixed a fail to load graphs with many large attributes totaling more than 1 GB in size

### [1.15.0.1.0.315.0]

Date 12/20/2019

## Summary

No major changes or fixes. Released with other Neuron packages.

### [1.15.0.1.0.306.0]

Date 12/1/2019

## Summary

## Major New Features

## Resolved Issues

## Known Issues & Limits

Same as prior release

### Other Notes

[1.15.0.1.0.280.0]

Date 11/29/2019

### Summary

Initial release packaged with DLAMI.

### Major New Features

N/A, initial release.

See user guide here: <https://github.com/aws/aws-neuron-sdk/blob/master/docs/neuron-tools/getting-started-tensorboard-neuron.md>

### Resolved Issues

N/A - first release

### Known Issues & Limits

- Must install TensorBoard-Neuron by itself, or after regular TensorBoard is installed. If regular Tensorboard is installed after TensorBoard-Neuron, it may overwrite some needed files.
- Utilization missing in Op Profile due to missing FLOPs calculation (see overview page instead)
- Neuron Profile plugin may not immediately show up on launch (try reloading the page)
- Graphs with NeuronOps may take a long time to load due to attribute size
- Instructions that cannot be matched to a framework layer/operator name show as "" (blank)
- CPU Usage section in chrome-trace is not applicable
- Debugger currently supports TensorFlow only
- Visualization requires a TensorFlow-compatible graph

### Other Notes

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## COMPILER

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 7.1 Neuron Compiler

The Neuron Compiler accepts Machine Learning models in various formats (TensorFlow, MXNet, PyTorch, XLA HLO) and optimizes them to run on Neuron devices.

The Neuron compiler is invoked within the ML framework, where ML models are sent to the compiler by the Neuron Framework plugin. The resulting compiler artifact is called a NEFF file (Neuron Executable File Format) that in turn is loaded by the Neuron runtime to the Neuron device.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### 7.1.1 NeuronX Compiler for Trn1 & Inf2

*This document is relevant for: Inf2, Trn1, Trn2*

#### API Reference Guide

*This document is relevant for: Inf2, Trn1, Trn2*

#### Neuron Compiler CLI Reference Guide (`neuronx-cc`)

This document describes the command line interface of the Neuron Compiler.

This reference is not relevant for applications that run the Neuron Compiler from within a machine learning framework (*PyTorch-Neuron* for example) since these options are passed from the framework directly to the compiler. Using the compiler command line may be desirable for applications that do not use a framework or customize existing frameworks. It is also possible to specify compiler options within the framework which will forward these options to the compiler using `NEURON_CC_FLAGS`.

## Usage

*Optional parameters are shown in square brackets.*

## Neuron Compiler Command-Line Interface

**neuronx-cc** <command> [parameters]

### Available Commands:

- compile
- list-operators

Common parameters for the Neuron CLI:

- **--help:** Display a usage message of compiler options.  
Use `neuronx-cc <command> --help` for information on a specific command.

**neuronx-cc** compile [parameters]

Compile a model for use on the AWS Machine Learning Accelerator.

```
neuronx-cc compile <model_files>
--framework <framework_name>
--target <instance_family>
[--model-type <model>]
[--auto-cast <cast_mode>]
[--auto-cast-type <data_type>]
[--distribution-strategy <distribution_type>]
[--logical-nc-config <shard_degree>], or [-lnc <shard_degree>]
[--optlevel <opt_level>], or [-O <opt_level>]
[--enable-mixed-precision-accumulation]
[--enable-saturate-infinity]
[--enable-fast-context-switch]
[--enable-fast-loading-neuron-binaries]
[--logfile <filename>]
[--output <filename>]
[--verbose <level>]
```

### Compile Parameters:

- **<model\_files>: Input containing model specification.**  
The number of arguments required varies between frameworks:
  - **XLA:** A local filename of a HLO file (hlo.pb) generated via XLA. See [hlo.proto](#) for the .proto description and [inspect-compiled-programs](#) for more information on how to generate such files.
- **--framework:** Framework used to generate training model.  
Valid values:
  - XLA
- **--target:** Name of the Neuron instance family on which the compiled model will be run.  
Valid values:
  - inf2
  - trn1
  - trn1n

- `trn2`
- `--model-type`: Permit the compiler to attempt model-specific optimizations based upon type of model being compiled. (Default: `generic`)

Valid values:

- `generic`: Perform optimizations applicable to all types of inference and training models.
- `transformer`: Perform optimizations specific to [Transformer](#) models.
- `unet-inference`: Perform optimizations specific to certain [U-Net](#) model architectures when performing inference. U-Net models often have certain structures that result in excessive performance-impacting data transfers; this option allows the compiler to apply additional memory optimizations to prevent these data transfers and also allows the compiler to map larger normalization operators which would otherwise not successfully execute.
- `--auto-cast`: Controls how the compiler makes tradeoffs between performance and accuracy for FP32 operations. (Default: `matmult`)

Valid values:

- `matmult`: Only cast FP32 operations that use the Neuron matrix-multiplication engine.
- `all`: Cast all FP32 operations to achieve highest performance. This option can potentially lower precision/accuracy.
- `none`: Leave all data types as defined in the model. Do not apply auto-casting data type optimizations.

A more complete discussion on how to use this option and its arguments is in [Mixed Precision and Performance-accuracy Tuning for Training](#).

---

**Note:** If the `--auto-cast` option is specified, the `--auto-cast-type` compiler flag can be optionally set to define which lower-precision data type the compiler should use.

---

- `--auto-cast-type`: When auto-cast mode is enabled, cast the FP32 operators to the lower-precision data type specified by this option. (Default: `bf16`)

Valid values:

- `bf16`: Cast the FP32 operations selected via the `--auto-cast` option to BF16 to achieve highest performance and preserve dynamic range.
- `fp16`: Cast the FP32 operations selected via the `--auto-cast` option to FP16 to achieve improved performance relative to FP32 and increased precision relative to BF16.
- `tf32`: Cast the FP32 operations selected via the `--auto-cast` option to TensorFloat-32.
- `fp8_e4m3`: Cast the FP32 operations selected via the `--auto-cast` option to a signed 8-bit floating point represented as a 4-bit exponent and 3-bit mantissa.

---

**Note:** If multiple competing options are specified then the option right-most on the command line will supersede previous options.

---

- `--distribution-strategy`: Permit the compiler to attempt model-specific optimizations based upon type of model being compiled. (Default: `generic`)

Valid values:

- `llm-training`: Enable the compiler to perform optimizations applicable to large language model (LLMs) training runs that shard parameters, gradients, and optimizer states across data-parallel workers. This is equivalent to the previously documented option argument value of `NEMO`, which will be deprecated in a future release.
- `--logical-nc-config`: Instructs the compiler to shard the input graph across physical NeuronCore accelerators. Possible numeric values are {1, 2}. (Only available on `trn2`; Default: 2)

Valid values:

- 1: instructs the compiler to shard the input graph across 1 physical NeuronCore, i.e., do not perform any input graph sharding.

- 2: [default on trn2] instructs the compiler to shard the input graph across 2 physical NeuronCores.
- **--optlevel**: Specify the level of optimization the compiler should perform. Possible numeric values are {1, 2, 3}. (Default: 2)

Valid values:

- 1: enables the core performance optimizations in the compiler, while also minimizing compile time.
- 2: [default] provides the best balance between model performance and compile time.
- 3: may provide additional model execution performance but may incur longer compile times and higher host memory usage during model compilation.

---

**Note:** This option supercedes, and deprecates, the **-enable-experimental-01** option introduced in an earlier release.

---

- **--enable-mixed-precision-accumulation**: Perform intermediate calculations of accumulation operators (such as softmax and layernorm) in FP32 and cast the result to the model-designated datatype. This improves the operator's resulting accuracy.
- **--enable-saturate-infinity**: Convert +/- infinity values to MAX/MIN\_FLOAT for compiler-introduced matrix-multiply transpose computations that have a high risk of generating Not-a-Number (NaN) values. There is a potential performance impact during model execution when this conversion is enabled. (Only needed on trn1; while the trn2 compiler will accept this flag for compatibility reasons, it has no effect on the compilation.)
- **--enable-fast-context-switch: Optimize for faster model switching rather than execution latency.**  
This option will defer loading some weight constants until the start of model execution. This results in overall faster system performance when your application switches between models frequently on the same Neuron Core (or set of cores).
- **--enable-fast-loading-neuron-binaries: Save the compilation output file in an uncompressed format.**  
This creates executable files which are larger in size but faster for the Neuron Runtime to load into memory during model execution.
- **--logfile**: Filename where compiler writes log messages. (Default: "log-neuron-cc.txt").
- **--output**: Filename where compilation output (NEFF archive) will be recorded. (Default: "file.neff")
- **--verbose**: Specify the level of output produced by the compiler. (Default: warning)

Valid values:

- **info**: Informational messages regarding the progress of model compilation (written to stdout).
- **warning**: Diagnostic messages that report model code that is not inherently erroneous but may be risky or suggest there may have been an error (written to stderr).
- **error**: The compiler detected a condition causing it not complete the compilation successfully (written to stderr).
- **critical**: The compiler encountered an unrecoverable error terminates immediately (written to stderr).
- **debug**: Extensive information regarding the compiler's internal execution phases (written to stdout).

**Example:**

Compiling an XLA HLO:

```
neuronx-cc compile bert-model.hlo --framework XLA --target trn1 --model-type_
↪transformer --output bert.neff
```

**neuronx-cc list-operators [parameters]**

Returns a newline ('\\n') separated list of operators supported by the Neuron Compiler.

```
neuronx-cc list-operators
--framework <value>
```

*List-Operators Parameters:*

- `--framework`: Framework in which the operators were registered.

Valid values:

- XLA: Operator names will be formatted according to the value used by XLA compiler in XlaBuilder.

*Example:*

```
neuronx-cc list-operators -framework XLA
...
```

*Exit Statuses:*

- **0**: Compilation succeeded
- **<>0**: An error occurred during compilation.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Developer Guide

*This document is relevant for:* Inf2, Trn1, Trn2

## Mixed Precision and Performance-accuracy Tuning (neuronx-cc)

### Table of contents

- *Overview*
- *Neuron Hardware*
- *Performance-accuracy tradeoffs*
- *What is the difference between Data Types?*
- *Should I downcast operations to smaller Data Types?*

## Overview

The Neuron Compiler supports machine learning models with FP32, TF32, FP16 and BF16 (Bfloat16) tensors and operators. The Neuron hardware supports a mix of 32, 16, and 8 bit datatypes. This guide explains how to apply the available auto-cast methods and their performance / accuracy trade-offs when compiling a model with Neuron.

---

**Note:** Neuron Compiler support for INT8 is planned for a future Neuron SDK release. See [Neuron Compiler: Enable Neuron INT8 support](#) for details.

---

## Neuron Hardware

The Neuron v2 hardware supports matrix multiplication using FP16, BF16, TF32, and FP32 on its matrix multiply (“matmult”) engine, and accumulations using FP32. Operators such as activations or vector operations are supported using FP32, TF32, FP16, and BF16. Supporting FP16 and BF16 allows Neuron to have significantly higher performance than executing everything as FP32.

## Performance-accuracy tradeoffs

**By default**, the Neuron Compiler will **automatically cast FP32 matrix multiplication operations to BF16**. The remaining operations are performed in the data type specified by the model. The Neuron Compiler provides CLI options that direct the compiler to cast to other data types, thereby giving the ability to choose an accuracy-to-performance tradeoff in model execution. Deciding what CLI settings to use will be application specific and may require some experimentation. See *Neuron Compiler CLI Reference Guide* for details.

## What is the difference between Data Types?

The NeuronCore v2 support multiple data types (see *NeuronCore v2 Data Types*). Each data type provides benefits and drawbacks due to its dynamic range and numeric precision.

| Type | Minimum   | Maximum  | Strength   | Weakness  |
|------|-----------|----------|--|---|
| FP16 | -65504    | 65504    | Numeric Precision, High granularity, Mid-range numbers | Low range, medium precision                       |
| BF16 | -3.40E+38 | 3.40E+38 | Dynamic Range, Extremely small/large numbers           | Low precision                                     |
| TF32 | -3.40E+38 | 3.40E+38 | Dynamic Range, Extremely small/large numbers           | Medium precision                                  |
| FP32 | -3.40E+38 | 3.40E+38 | N/A  | Larger model size, potentially slower computation |

- FP16 provides a high density of representable values that are neither extremely small or extremely large. The density of representable values within the range is approximately an order of magnitude greater than BF16.
  - Conversion from FP32 to FP16 will perform well when values are relatively small but non-extreme (either very small or very large).
  - Conversion from FP32 to FP16 will perform badly if the original FP32 values are outside of the range of FP16. This will produce inf/-inf values and may result in NaN depending on the operation.
- BF16 provides a wider range of representable values which includes both very small and very large values. However, the overall density of representable values is usually lower than FP16 for more non-extreme values. The range is nearly identical to the range of FP32 but because the number of bits is halved, this means the individual values are sparse.
  - Conversion from FP32 to BF16 will perform well when the values are well-distributed throughout the range. Since BF16 covers the entire FP32 range, this means each original value can map to a relatively close downcast value.
  - Conversion from FP32 to BF16 will perform badly when fine granularity is needed. Since BF16 granularity is sacrificed for greater range it will almost always map worse to values that are within the FP16 range.



## Should I downcast operations to smaller Data Types?

This choice here is driven entirely by accuracy vs performance tradeoff. Casting operations to smaller 16-bit data types will provide a significant performance benefit but may end up sacrificing accuracy.

The compiler uses BF16 casting **by default** for matrix multiplication operations. The speedup from casting operations gives a significant performance boost and the range of representable values in BF16 allows for more safety compared to FP16 when the possible numeric range of input values is unknown.

The Neuron Compiler's `--auto-cast` and `--auto-cast-type` CLI options are used to direct the compiler to perform alternate casting operations. See the detailed list of the options in [Neuron v2 Compiler CLI Reference Guide](#).

It is recommended that you start with compiling the model to achieve high performance (default), you can then test the accuracy of the application and, if needed, try the next higher precision casting option until the desired accuracy and performance are achieved.

The option combinations to consider in a typical flow are:

| Compiler autocast   | Options Effect   | Performance   | Accuracy   |
|---|--|---|--|
| --auto-cast all<br>--auto-cast-type bf16<br>--auto-cast matmult<br>--auto-cast-type bf16<br>(default) | Best performance at the expense of precision                             | Performance <i>decreases</i> as you move down the table | Accuracy <i>increases</i> as you move down the table |
| --auto-cast all<br>--auto-cast-type fp16<br>--auto-cast matmult<br>--auto-cast-type fp16              | Best performance at the expense of dynamic range                         |   |  |
| --auto-cast all<br>--auto-cast-type tf32<br>--auto-cast matmult<br>--auto-cast-type tf32              | Balance of performance, dynamic range, and precision                     |   |  |
| --auto-cast none  | Disables all auto-casting, using the data types defined within the model |   |  |

Note that compiler has to preserve the input/output (i/o) tensor types requested by Framework, therefore no casting is done on the i/o tensors. Additional speedup can be obtained by casting them in the Framework prior to compilation.

To learn how to configure the compiler options from within your application's framework, please see:

- [Developer Guide for Training with PyTorch Neuron](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Misc (neuronx-cc)

*This document is relevant for: Inf2, Trn1, Trn2*

## Neuron Compiler FAQ (neuronx-cc)

### Table of contents

- *Where can I compile to Neuron?*
- *What is the difference between neuron-cc and neuronx-cc?*
- *Should I use neuron-cc or neuronx-cc?*
- *My current neural network is based on FP32, how can I use it with Neuron?*
- *Which operators does Neuron support?*
- *Any operators that Neuron Compiler doesn't support?*
- *Will I need to recompile again if I updated runtime/driver version?*
- *I have a NEFF binary, how can I tell which compiler version generated it?*
- *How long does it take to compile?*
- *Why is my model producing different results compared to CPU/GPU?*
- *Do you support model <insert model type>?*

## Where can I compile to Neuron?

The one-time compilation step from the standard framework-level model to NEFF binary may be performed on any EC2 instance or even on-premises.

We recommend using a high-performance compute server of choice (C5 or z1d instance types), for the fastest compile times and ease of use with a prebuilt [DLAMI](#). Developers can also install Neuron in their own environments; this approach may work well for example when building a large fleet for inference, allowing the model creation, training and compilation to be done in the training fleet, with the NEFF files being distributed by a configuration management application to the inference fleet.

## What is the difference between neuron-cc and neuronx-cc?

- neuron-cc is the Neuron Compiler with TVM front-end, neuron-cc supports only neuroncores-v1-arch.
- neuronx-cc is the Neuron Compiler with XLA front-end, neuronx-cc currently supports neuroncores-v2-arch, neuronx-cc support of neuroncores-v1-arch is currently a [Roadmap Item](#).

### Should I use `neuron-cc` or `neuronx-cc`?

See *What is the difference between `neuron-cc` and `neuronx-cc`?*

### My current neural network is based on FP32, how can I use it with Neuron?

Developers who want to train their models in FP32 for best accuracy can compile and deploy them with Neuron. The Neuron compiler automatically converts FP32 to internally supported datatypes, such as FP16 or BF16. You can find more details about FP32 data type support and performance and accuracy tuning in *Mixed Precision and Performance-accuracy Tuning (neuronx-cc)* or *Mixed precision and performance-accuracy tuning (neuron-cc)*. The Neuron compiler preserves the application interface - FP32 inputs and outputs. Transferring such large tensors may become a bottleneck for your application. Therefore, you can improve execution time by casting the inputs and outputs to FP16 or BF16 in the ML framework prior to compilation.

### Which operators does Neuron support?

You can use the `neuronx-cc list-operators` command on the cli to list the operators. See *Neuron Compiler CLI Reference Guide (neuronx-cc)*.

To request support for new operators, open an issue on our [GitHub forum](#).

### Any operators that Neuron Compiler doesn't support?

Models with control-flow and dynamic shapes are not supported now. You will need to partition the model using the framework prior to compilation.

---

**Note:** Starting with `neuroncores-v2-arch` Neuron supports control-flow and dynamic shapes.

Stay tuned and follow the *Neuron Roadmap*.

---

### Will I need to recompile again if I updated runtime/driver version?

The compiler and runtime are committed to maintaining compatibility for major version releases with each other. The versioning is defined as major.minor, with compatibility for all versions with the same major number. If the versions mismatch, an error notification is logged and the load will fail. This will then require the model to be recompiled.

### I have a NEFF binary, how can I tell which compiler version generated it?

\*\* We will bring a utility out to help with this soon.

### How long does it take to compile?

It depends on the model and its size and complexity, but this generally takes a few minutes.

### Why is my model producing different results compared to CPU/GPU?

neuroncores-v2-arch supports multiple casting modes for floating point numbers, each with associated implications for performance and accuracy. The default casting mode is a pragmatic balance between performance and accuracy, however on some models it may result in loss of precision.

See the `--auto-cast` and `--auto-cast-type` options in *Neuron Compiler CLI Reference Guide (neuronx-cc)* for details on how to adjust the casting mode.

### Do you support model *<insert model type>*?

neuronx-cc has explicit support for select model families using the `--model-type` option, though many other model types are supported. You can also inspect supported operators using the `list-operators` sub-command. See the *Neuron Compiler CLI Reference Guide (neuronx-cc)* for details. More generally, support for new operators and models is continually being added. See our *Roadmap* for details.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### Neuron Compiler (neuronx-cc) release notes

#### Neuron Compiler [2.17.194.0]

Date: 04/03/2025

- Minor bug fixes and performance enhancements for both the trn1 and trn2 platforms.

#### Neuron Compiler [2.16.372.0]

Date: 01/14/2025

- Minor bug fixes and performance enhancements for the trn2 platform.

#### Neuron Compiler [2.16.345.0]

Date: 12/20/2024

- Minor bug fixes and performance enhancements for the trn2 platform.

## Neuron Compiler

Date: 12/03/2024

- This release introduces the `trn2` option argument to the compiler `--target` option to specify that the compiler should generate code for a `trn2` instance family. Example usage: `neuronx-cc compile --target=trn2 ...`
- This release introduces the `--logical-nc-config` or `-lnc` compiler command line option in support of the Logical NeuronCore Configuration feature available in Trainium2 instances. The compiler's default is `LNC=2`. **Note: Use of this option is available only for Trainium2 instances.**

## Neuron Compiler [2.15.128.0]

Date: 09/16/2024

- This release introduces memory optimization that will reduce the generated compiler artifacts size (i.e., NEFFs) and the models' memory footprint. It is possible that some models may experience unexpected performance degradation. If this occurs, these optimizations can be disabled using the `-disable-dge` compiler command line option or the framework-level option `additional_compile_opt="--disable-dge"`

## Neuron Compiler [2.14.213.0]

Date: 07/03/2024

- Minor bug fixes and performance enhancements.
- Improved flash attention kernel performance.

## Neuron Compiler [2.13.72.0]

Date: 04/25/2024

- Minor bug fixes and enhancements.

## Neuron Compiler [2.13.68.0]

Date: 04/10/2024

- This release fixes hang issues related to Triton Inference Server.

## Neuron Compiler [2.13.66.0]

Date: 04/01/2024

- This release introduces a new `--enable-mixed-precision-accumulation` compiler option. This option instructs the compiler to perform intermediate calculations of reduction operators (such as the dot or reduce operators) in FP32 regardless of the operation's defined datatype. The final result of the operator will be cast from FP32 to the model-designated datatype (e.g., BF16). This helps to improve the operator's resulting accuracy.

**Neuron Compiler [2.12.68.0]**

Date: 01/18/2024

- Patch release with bug fixes.

**Neuron Compiler [2.12.54.0]**

Date: 12/21/2023

- The compiler now generates instructions to check if a model references an embedding table with an illegal index. The check is made at model execution time. If an attempted invalid table index is encountered, the model execution will continue and the user will see an error similar to:

WARNING: Received notification generated at runtime: failed to run scatter/gather (indirect memory copy with branch\_label\_id = xx), due to out-of-bound access.

When this occurs, users are encouraged to review the model's gather/scatter input values to determine if there is a coding error.

**Neuron Compiler [2.11.0.35]**

Date: 11/17/2023

- This release addresses performance related issues when training through `neuronx-nemo-megatron` library.

**Neuron Compiler [2.11.0.34]**

Date: 10/26/2023

- This release introduces the option-argument `llm-training` to the existing `--distribution_strategy` compiler option. This option-argument allows the compiler to make specific optimizations related to training distributed models. This new option-argument is equivalent to the previously introduced `nemo` option-argument, which will be deprecated in a future release.

**Neuron Compiler [2.10.0.35]**

Date: 09/26/2023

- This release addresses a compilation regression for certain configurations of Llama and Llama-2 inference models when it fails compilation with this error "IndirectLoad/Save requires contiguous indirect access per partition" .

There is still a known issue for some configurations of the model with the error "Too many instructions after unroll for function sg0000" . To mitigate this, recompile using the `--optlevel 1 (-O1)` option. A complete fix will be coming in the future release which will not require this option

## Neuron Compiler [2.10.0.34]

Date: 09/15/2023

- This release introduces a new `--optlevel (-O)` compiler option. This option allows the user to balance between compile-time and optimizations performed. Three levels are supported. Level `--optlevel 1 (-O1)` aims to minimize compile-time and allow for a more rapid model development cycle. Model execution time may be reduced. Level `--optlevel 3 (-O3)` performs whole-model optimization. This level will deliver the best performance however there will be longer compile-times and the compiler will use more host DRAM, potentially requiring a larger instance to compile the model. The default is `--optlevel 2 (-O2)` which provides a balance between model performance and compile time.

The previous `-enable-experimental-O1` flag introduced in the 02/08/2023 Neuron Compiler [2.4.0.21] release is now deprecated. Using this flag will generate a message similar to:

WARNING: Option `--enable-experimental-O1` is deprecated and will be removed in a future release." Use `--optlevel 1 (-O1)` instead.

## Neuron Compiler [2.9.0.16]

Date: 08/28/2023

- This release fixes an issue where any initial seed passed into the Random Number Generator operator was not honored. The `RngBitGenerator` operator now correctly accepts and uses setting the seed. Note that the current RNG implementation only supports 32-bit seeds.

## Neuron Compiler [2.8.0.25]

Date: 07/19/2023

- This release introduces a new optional `--distribution_strategy` compiler option. This option informs the compiler what type of distributed APIs are used to shard the model and allows the compiler to make API-specific optimizations. Currently following option-arguments are supported: `nemo`.

## Neuron Compiler [2.7.0.40]

Date: 06/14/2023

- This release introduces a new `--enable-saturate-infinity` compiler option. A computation that can generate +/- infinity is at a high risk of generating Not-a-Number (NaN) values when the infinity value is used in subsequent computations. This option helps avoid this by converting `+Inf/-Inf` values to `MAX/MIN_FLOAT` before operations that could produce NaN values for `+Inf/-Inf` inputs on the target architecture. While this option helps to avoid NaN values, there is a potential performance degradation that occurs during model execution when this conversion is enabled.

**Neuron Compiler [2.6.0.19]**

Date: 05/01/2023

- This release introduces a new `model-type` option argument: `unet-inference`. This option instructs the compiler to perform model-specific optimizations that produce executable models with improved performance on the specified target instance.
- Added support for the HLO operator `BitcastConvertType` and also added support for `TopK` (sampling mode) operator.

**Neuron Compiler [2.5.0.28]**

Date: 03/28/2023

- This release introduces the `trn1n` option argument to the compiler `target` option to specify that it should generate code for a `trn1n` instance type. Example usage: `neuronx-cc compile --target=trn1n ...`
- The compiler's usage message now includes the `inf2` option argument.
- A new 8-bit floating point data type, `fp8_e4m3`, is now supported and can be specified using the `auto-cast-type` option. This instructs the compiler to convert the FP32 operations selected via the `--auto-cast` option to a signed FP8 size with 4-bit exponent and 3-bit mantissa. Care must be taken to ensure that the down-casted values are representable within the 8-bit data range.

**Neuron Compiler [2.4.0.21]**

Date: 02/24/2023

- This release introduces the `inf2` option argument to the compiler `target` option to specify that it should generate code for an `inf2` instance type. Example usage: `neuronx-cc compile --target=inf2 ...`. The `inf2` option argument does not appear in the compiler's usage message. It will be added in the next release.

**Neuron Compiler [2.4.0.21]**

Date: 02/08/2023

- Added support for the following HLO operators: `SelectAndScatter`.
- Beta: `--enable-experimental-01` flag: This option reduces the compile-time with a negligible impact on model execution performance. It allows the compiler to execute compiler passes in parallel to perform the compilation. By default the compiler uses 8 processes. This can be changed via the CLI option `--num-parallel-jobs`. This option is expected to become the default in a future SDK release.

**Neuron Compiler [2.3.0.4]**

Date: 12/09/2022

- Added support for the following HLO operators: `rev` (reverse).
- The `pow()` function can now handle both integer and floating-point exponents.
- Optimization enhancements and bug fixes to improve model execution performance.



## Neuron Compiler [2.2.0.73]

Date: 10/27/2022

- Adding support for the following HLO operators: `LogicalNot`, `atan2` and `DynamicUpdateSlice` (for constant index).

## Neuron Compiler [2.1.0.76]

Date: 10/5/2022

The Neuron Compiler is an Ahead-of-Time compiler that accelerates models for execution on NeuronCores. This release supports compiling models for training on a Trn1 instance using Pytorch Neuron. Users typically access the compiler via the Framework to perform model compilation, although it can also be run as a command line tool (*neuronx-cc*).

The Neuron Compiler supports compiling models for mixed precision calculations. The trn1 hardware supports matrix multiplication using FP16, BF16, and FP32 on its Matrix Multiplication Engine, and accumulations using FP32. Operators such as activations or vector operations are supported using FP16, BF16, and FP32. Tensor transpose can be accomplished in FP16, BF16, FP32, or TF32 datatypes. By default, scalar and vector operations on FP32 values will be done in FP32, while matrix multiplications are cast to BF16 and transpose operations are cast to FP32. This default casting will generate the highest performance for a FP32 trained model.

By default, the compiler will target maximum performance by automatically casting the model to mixed precision. It also provides an option (`--auto-cast`) that allows the user to make tradeoffs between higher performance and optimal accuracy. The decision on what option argument to use with the `--auto-cast` option will be application specific. Compiler CLI options can be passed to the compiler via the framework.

## Known issues

- The Random Number Generator operation can be passed an initial seed value, however setting the seed is not supported in this release.
- The exponent value of the `pow()` function must be a compile-time integer constant.
- The compiler treats INT64 datatypes as INT32 by truncating the high-order bits. If possible, cast these values to 32 bits .
- Model compilation time is proportional to the model size and operators used. For some larger NLP models it may be upwards of 30 minutes.

## Supported Operators

The following XLA operators are supported by the Neuron Compiler. Future releases will broaden model support by providing additional XLA operators defined in [https://www.tensorflow.org/xla/operation\\_semantics](https://www.tensorflow.org/xla/operation_semantics).

The list of supported operators can also be retrieved from the command line using *neuronx-cc list-operators*.

| Supported XLA Operators | Notes |
|-------------------------|-------|
| Abs                     |       |
| Add                     |       |
| Allgather               |       |
| Allreduce               |       |

continues on next page

Table 7.1 – continued from previous page

| Supported XLA Operators | Notes   |
|-------------------------|---|
| Atan2                   |   |
| Batchnorm               |   |
| Batchnormgrad           |   |
| Batchnorminference      |   |
| BitcastConvertType      |   |
| Broadcast               |   |
| BroadcastInDim          |   |
| Ceil                    |   |
| Clamp                   |   |
| Compare                 |   |
| Concatenate             |   |
| Constant                |   |
| ConstantLiteral         |   |
| ConvertElementType      |   |
| Cos                     |   |
| Customcall              |   |
| Div                     |   |
| Dot                     |   |
| DotGeneral              |   |
| DynamicUpdateSlice      | Supports only for constant index  |
| Eq                      |   |
| Exp                     |   |
| Floor                   |   |
| Gather                  | Supports only disjoint start_index_map and remapped_offset_dims                         |
| Ge                      |   |
| GetTupleElement         |   |
| Gt                      |   |
| Iota                    |   |
| Le                      |   |
| Log                     |   |
| LogicalAnd              |   |
| LogicalNot              |   |
| Lt                      |   |
| Max                     |   |
| Min                     |   |
| Mul                     |   |
| Ne                      |   |
| Neg                     |   |
| Pad                     |   |
| Pow                     | Exponent argument must be a compile-time integer constant                               |
| Reduce                  | Min, Max, Add and Mul are the only supported computations. Init_values must be constant |
| Reshape                 |   |
| Rev (reverse)           |   |
| RngBitGenerator         | Ignores user seed   |
| RngUniform              |   |
| Rsqrt                   |   |
| Scatter                 |   |
| Select                  |   |
| SelectAndScatter        |   |

continues on next page

Table 7.1 – continued from previous page

| Supported XLA Operators | Notes |
|-------------------------|-------|
| ShiftRightLogical       |       |
| Sign                    |       |
| Sin                     |       |
| Slice                   |       |
| Sqrt                    |       |
| Sub                     |       |
| Tanh                    |       |
| Transpose               |       |
| Tuple                   |       |

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 7.1.2 Neuron Compiler for Inf1

*This document is relevant for:* Inf1

### API Reference Guide

*This document is relevant for:* Inf1

### Neuron compiler CLI Reference Guide (neuron-cc)

This document describes the command line interface of the Neuron compiler. This reference is not relevant for applications that run neuron-cc from within a machine learning framework (TensorFlow-Neuron for example) since these options are passed from the framework directly to neuron-cc.

Using neuron-cc on the command line may be desirable for applications that do not use a framework, or customize existing frameworks. It is also possible to supply CLI commands to the framework as options to be passed through to the compiler.

### Usage

Optional parameters are shown in square brackets. See the individual framework guides for the correct syntax.

## Neuron Compiler CLI

**neuron-cc** [options] <command> [parameters]

Common options for the Neuron CLI:

- **--verbose** (string) default="WARN":

Valid values:

- DEBUG
- INFO
- WARN
- ERROR

Use **neuron-cc** <command> **--help** for information on a specific command.

### Available Commands:

- compile
- list-operators

**neuron-cc** compile [parameters]

Compile a model for use on the AWS Inferentia Machine Learning Accelerator.

```
neuron-cc compile <file names> --framework <value> --io-config <value> [--
↪neuroncore-pipeline-cores <value>] [--enable-saturate-infinity] [--enable-fast-
↪loading-neuron-binaries] [--enable-fast-context-switch] [--fp32-cast cast-method]
↪[--fast-math cast-method] [--output <value>]
```

### Compile Parameters:

- **<file names>**: Input containing model specification. The number of arguments required varies between frameworks:
  - **TENSORFLOW**: A local filename or URI of a TensorFlow Frozen GraphDef (.pb); or the name of a local directory containing a TensorFlow SavedModel.

See <https://github.com/tensorflow/tensorflow/blob/master/tensorflow/core/framework/graph.proto> for the associated .proto schema for TensorFlow Frozen GraphDefs. See [https://www.tensorflow.org/guide/saved\\_model](https://www.tensorflow.org/guide/saved_model) for more information on the SavedModel format.

- **MXNET**: List of local filenames or URIs where input architecture .json file and parameter .param file are stored. These contains information related to the architecture of your graph and associated parameters, respectively.
- **--framework** (string): Framework in which the model was trained.

Valid values:

- TENSORFLOW
- MXNET
- XLA

- **--neuroncore-pipeline-cores** (int) (default=1): Number of neuron cores to be used in “NeuronCore Pipeline” mode. This is different from data parallel deployment (same model on multiple neuron cores). Refer to Runtime/Framework documentation for data parallel deployment options.

Compile for the given number of neuron cores so as to leverage NeuronCore Pipeline mode.

**Note:** This is not used to define the number of Neuron Cores to be used in a data parallel deployment (ie the same model on multiple Neuron Cores). That is a runtime/framework configuration choice.

- `--output` (string) (default="out.neff"): Filename where compilation output (NEFF archive) will be recorded.
- `--io-config` (string): Configuration containing the names and shapes of input and output tensors.

The io-config can be specified as a local filename, a URI, or a string containing the io-config itself.

The io-config must be formatted as a JSON object with two members "inputs" and "outputs". "inputs" is an object mapping input tensor names to an array of shape and data type. "outputs" is an array of output tensor names. Consider the following example:

```
{
  "inputs": {
    "input0:0": [[1,100,100,3], "float16"],
    "input1:0": [[1,100,100,3], "float16"]
  },
  "outputs": ["output:0"]
}
```

- `--enable-saturate-infinity`: Convert +/- infinity values to MAX/MIN\_FLOAT for certain computations that have a high risk of generating Not-a-Number (NaN) values. There is a potential performance impact during model execution when this conversion is enabled.
- `--enable-fast-loading-neuron-binaries`: Write the compilation output (NEFF archive) in uncompressed format which results in faster loading of the archive during inference.
- `--enable-fast-context-switch`: Optimize for faster model switching rather than inference latency. This results in overall faster system performance when your application switches between models frequently on the same neuron core (or set of cores). The optimization triggered by this option for example defers loading some weight constants until the start of inference.
- `--fast-math`: Controls tradeoff between performance and accuracy for fp32 operators. See more suggestions on how to use this option with the below arguments in *Mixed precision and performance-accuracy tuning (neuron-cc)*.
  - `all` (Default): enables all optimizations that improve performance. This option can potentially lower precision/accuracy.
  - `none`: Disables all optimizations that improve performance. This option will provide best precision/accuracy.
  - Tensor transpose options
    - \* `fast-relayout`: Only enables fast relayout optimization to improve performance by using the matrix multiplier for tensor transpose. The data type used for the transpose is either FP16 or BF16, which is controlled by the `fp32-cast-xxx` keyword.
    - \* `no-fast-relayout`: Disables fast relayout optimization which ensures that tensor transpose is bit-accurate (lossless) but slightly slower.
  - Casting options
    - \* `fp32-cast-all` (Default): Cast all FP32 operators to BF16 to achieve highest performance and preserve dynamic range. Same as setting `--fp32-cast all`.
    - \* `fp32-cast-all-fp16`: Cast all FP32 operators to FP16 to achieve speed up and increase precision versus BF16. Same setting as `--fp32-cast all-fp16`.
    - \* `fp32-cast-matmult`: Only cast FP32 operators that use Neuron Matmult engine to BF16 while using FP16 for matmult-based transpose to get better accuracy. Same as setting `--fp32-cast matmult`.
    - \* `fp32-cast-matmult-bf16`: Cast only FP32 operators that use Neuron Matmult engine (including matmult-based transpose) to BF16 to preserve dynamic range. Same as setting `--fp32-cast matmult-bf16`.

- \* `fp32-cast-matmult-fp16`: Cast only FP32 operators that use Neuron Matmult engine (including matmult-based transpose) to fp16 to better preserve precision. Same as setting `--fp32-cast matmult-fp16`.

**Important:**

- `all` and `none` are mutually exclusive
- `all` is equivalent to using `fp32-cast-all fast-relayout` (best performance)
- `none` is equivalent to using `fp32-cast-matmult-bf16 no-fast-relayout` (best accuracy)
- `fp32-cast-*` options are mutually exclusive
- `fast-relayout` and `no-fast-relayout` are mutually exclusive
- The `fp32-cast-*` and `*-fast-relayout` options will overwrite the default behavior in `all` and `none`.
- For backward compatibility, the `--fp32-cast` option has higher priority over `--fast-math`. It will overwrite the FP32 casting options in any of the `--fast-math` options if `--fp32-cast` option is present explicitly.

- `--fp32-cast` : Refine the automatic casting of fp32 tensors. This is being replaced by a newer `--fast-math`.

**Important:**

- `--fp32-cast` option is being deprecated and `--fast-math` will replace it in future releases.
- `--fast-math` is introducing the `no-fast-relayout` option to enable lossless transpose operation.

The `--fp32-cast` is an interface for controlling the performance and accuracy tradeoffs. Many of the `--fast-math` values invoke (override) it.

- `all` (default): Cast all FP32 operators to BF16 to achieve speed up and preserve dynamic range.
- `matmult`: Cast only FP32 operators that use Neuron Matmult engine to BF16 while using fp16 for matmult-based transpose to get better accuracy.
- `matmult-fp16`: Cast only FP32 operators that use Neuron Matmult engine (including matmult-based transpose) to fp16 to better preserve precision.
- `matmult-bf16`: Cast only FP32 operators that use Neuron Matmult engine (including matmult-based transpose) to BF16 to preserve dynamic range.
- `all-fp16`: Cast all FP32 operators to FP16 to achieve speed up and better preserve precision.

**Log Levels:**

Logs at levels “trace”, “debug”, and “info” will be written to STDOUT.

Logs at levels “warn”, “error”, and “fatal” will be written to STDERR.

**Exit Status**

**0** - Compilation succeeded

**>0** - An error occurred during compilation.

**Examples**

Compiling a saved TensorFlow model:

```
neuron-cc compile test_graph_tformatmul.pb --framework TENSORFLOW --io-
  ↪ config test_graph_tformatmul.config
```

Compiling a MXNet model:

```
neuron-cc compile lenet-symbol.json lenet-0001.params --framework MXNET --
↳neuroncore-pipeline-cores 2 --output file.neff
```

Compiling an XLA HLO:

```
neuron-cc compile bert-model.hlo --framework XLA --output file.neff
```

### **neuron-cc list-operators [parameters]**

Returns a newline ('n') separated list of operators supported by the NeuronCore.

- **TENSORFLOW**: Operators will be formatted according to the value passed to the associated REGISTER\_OP("OperatorName") macro.

See [https://www.tensorflow.org/guide/create\\_op#define\\_the\\_op\\_interface](https://www.tensorflow.org/guide/create_op#define_the_op_interface) for more information regarding operator registration in TensorFlow.

- **MXNET**: Operator names will be formatted according to the value passed to the associated NNVM\_REGISTER\_OP(operator\_name) macro.
- **XLA**: Operator names will be formatted according to the value used by XLA compiler in XlaBuilder.

See [https://www.tensorflow.org/xla/operation\\_semantics](https://www.tensorflow.org/xla/operation_semantics) for more information regarding XLA operator semantics in XLA interface.

```
neuron-cc list-operators --framework <value>
```

- **--framework** (string): Framework in which the operators were registered.

Valid values:

- TENSORFLOW
- MXNET
- XLA

### **Exit Status**

**0** - Call succeeded

**>0** - An error occurred

### **Example**

```
$ neuron-cc list-operators --framework TENSORFLOW
AddN
AdjustContrastv2
CheckNumbers
...
```

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Developer Guide

*This document is relevant for: Inf1*

### Mixed precision and performance-accuracy tuning (neuron-cc)

#### Table of contents

- *Neuron Hardware*
- *Performance-accuracy tradeoffs for models trained in FP32*
- *Compiler casting options*
  - *--fast-math option*

The Neuron Compiler supports machine learning models with FP32, FP16 and BF16 (Bfloat16) tensors and operators. The Neuron hardware supports a mix of 32 and 16 bit datatypes. The available auto-cast methods and their performance / accuracy trade-offs are explained in this document.

#### Neuron Hardware

The Neuron hardware supports matrix multiplication using FP16 or BF16 on its Matmult Engine, and accumulations using FP32. Similarly, operators such as activations or vector operations are supported using FP16, BF16 and FP32. Neuron supports tensor transpose in two ways - by fast matrix multiplication in FP16/BF16 or by slower byte-by-byte data movements.

#### Performance-accuracy tradeoffs for models trained in FP32

Models that are trained using FP32 data types can be deployed on Neuron through ahead of time compilation using the *Neuron Compiler*.

**By default**, the Neuron Compiler will **cast all FP32 tensors, weights and operations to BF16**. Only partial sums are left in FP32. The default, casting will generate the highest performance for a FP32 trained model.

Using the `--fast-math` CLI option, you can choose the right tradeoff between performance and accuracy. The tradeoff usually is between achieving high performance or optimal accuracy, and decision what settings to use will be application specific.

It is recommended that the you start with compiling the model to achieve the high performance (default), you can then test the accuracy of the application and, if needed, try the next higher precision casting option until the desired accuracy and performance are achieved. A typical flow can be:

1. You can compile without options (default) or with `--fast-math all` which will optimize for performance.
2. If accuracy is not sufficient you can try `--fast-math fp32-cast-matmult`
3. If accuracy is not sufficient you can try `--fast-math fp32-cast-matmult no-fast-relayout`
4. If accuracy is not sufficient you can try `--fast-math none` which will optimize for accuracy .

Between step 2 and step 3, and between step 3 and step 4 you have additional options that can provide different level of accuracy and which are explained in the below section.



Note that compiler has to preserve the input/output (i/o) tensor types requested by Framework, therefore no casting is done on the i/o tensors. Additional speedup can be obtained by casting them in the Framework prior compilation.

To learn how to use compiler command line interface (CLI) options with your application's framework, please see [PyTorch-Neuron trace python API](#), [tensorflow-ref-neuron-compile-api](#) and [TensorFlow 2.x \(tensorflow-neuron\) Tracing API](#).

## Compiler casting options

### --fast-math option

The `--fast-math` option is intended to replace the `--fp32-cast` option. It is recommended to start using or migrating to `--fast-math` option. The `--fast-math` option provides the same level of functionality as the `--fp32-cast` option in addition to the following:

- The `--fast-math` option introduces the `no-fast-relayout` option to enable lossless transpose operation. This was not possible with the `--fp32-cast` option.
- The `--fast-math` option provides finer control than the `--fp32-cast` option. The transpose operation and the cast operation are controlled independently:
  - `no-fast-relayout` and `fast-relayout` provide control for the transpose operation.
  - `fp32-cast-*` provide control for casting.

See the detailed list of the options in [Neuron compiler CLI Reference Guide \(neuron-cc\)](#).

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Misc (neuron-cc)

*This document is relevant for:* Inf1

## Neuron Compiler FAQ (neuron-cc)

### Table of contents

- [Where can I compile to Neuron?](#)
- [My current Neural Network is based on FP32, how can I use it with Neuron?](#)
- [What are some of the important compiler defaults I should be aware of?](#)
- [Which operators does Neuron support?](#)
- [Any operators that Neuron doesn't support?](#)
- [Will I need to recompile again if I updated runtime/driver version?](#)
- [I have a NEFF binary, how can I tell which compiler version](#)
- [How long does it take to compile?](#)

## Where can I compile to Neuron?

The one-time compilation step from the standard framework-level model to NEFF binary may be performed on any EC2 instance or even on-premises.

We recommend using a high-performance compute server of choice (C5 or z1d instance types), for the fastest compile times and ease of use with a prebuilt [DLAMI](#). Developers can also install Neuron in their own environments; this approach may work well for example when building a large fleet for inference, allowing the model creation, training and compilation to be done in the training fleet, with the NEFF files being distributed by a configuration management application to the inference fleet.

## My current Neural Network is based on FP32, how can I use it with Neuron?

Developers who want to train their models in FP32 for best accuracy can compile and deploy them with Neuron. The Neuron compiler automatically converts FP32 to internally supported datatypes, such as FP16 or BF16. You can find more details about FP32 data type support and performance and accuracy tuning in [Mixed precision and performance-accuracy tuning \(neuron-cc\)](#). The Neuron compiler preserves the application interface - FP32 inputs and outputs. Transferring such large tensors may become a bottleneck for your application. Therefore, you can improve execution time by casting the inputs and outputs to FP16 or BF16 in the ML framework prior to compilation for Inferentia.

## What are some of the important compiler defaults I should be aware of?

The compiler compiles the input graph for a single NeuronCore by default. Using the `--neuroncore-pipeline-cores` option directs the compiler to partition so as to run on a specified number of NeuronCores. This number can be less than the total available NeuronCores on an instance. See [inferentia-arch](#) for more information on NeuronCores.

## Which operators does Neuron support?

see [Neuron Supported operators](#).

You can also use the “`neuron-cc list-operators`” command on the cli to list the operators. See [neuron-cc-list-operators](#)

If your model contains operators missing from the above list, and you can't reach your performance goals, please post a message on the Neuron developer forum or open a github issue to let us know.

## Any operators that Neuron doesn't support?

Models with control-flow and dynamic shapes are not supported. You will need to partition the model using the framework prior to compilation. See the [neuron-cc](#).

## Will I need to recompile again if I updated runtime/driver version?

The compiler and runtime are committed to maintaining compatibility for major version releases with each other. The versioning is defined as major.minor, with compatibility for all versions with the same major number. If the versions mismatch, an error notification is logged and the load will fail. This will then require the model to be recompiled.

## I have a NEFF binary, how can I tell which compiler version

generated it? We will bring a utility out to help with this soon.

## How long does it take to compile?

It depends on the model and its size and complexity, but this generally takes a few minutes.

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Neuron Compiler (neuron-cc) for Inf1 Release Notes

### Table of contents

- *Introduction*
- *Known issues and limitations - updated 11/23/2022*
- *Neuron Compiler release [1.21.0.0]*
- *Neuron Compiler release [1.20.3.0]*
- *Neuron Compiler release [1.19.0.0]*
- *Neuron Compiler release [1.17.0.0]*
- *Neuron Compiler release [1.16.2.0]*
- *Neuron Compiler release [1.15.0.0]*
- *Neuron Compiler release [1.14.3.0]*
- *Neuron Compiler release [1.13.3.0]*
- *Neuron Compiler release [1.11.7.0]*
- *Neuron Compiler release [1.11.4.0]*
- *Neuron Compiler release [1.10.3.0]*
- *Neuron Compiler release [1.9.1.0]*
- *Neuron Compiler release [1.8.5.0]*
- *Neuron Compiler release [1.8.2.0]*
- *Neuron Compiler release [1.7.3.0]*
- *[1.6.13.0]*
- *[1.5.5.0]*

- [1.4.0.0]
- [1.3.0.0]
- [1.2.7.0]
- [1.2.2.0]
- [1.1.7.0]
- [1.0.24045.0]
- [1.0.20600.0]
- [1.0.18001.0]
- [1.0.17937.0]
- [1.0.16861.0]
- [1.0.15275.0]
- [1.0.12696.0]
- [1.0.9410.0]
- [1.0.7878.0]
- [1.0.6801.0]
- [1.0.5939.0]
- [1.0.5301.0]
- [1.0.4680.0]

## Introduction

This document lists the release notes for AWS Neuron compiler. The Neuron Compiler is an ahead-of-time compiler that ensures Neuron will optimally utilize the Inferentia chips.

Operator-support for each input format is provided directly from the compiler.

```
neuron-cc list-operators --framework {TENSORFLOW | MXNET | XLA}
```

The supported operators are also listed here:

Tensorflow: *TensorFlow Neuron (tensorflow-neuron (TF1.x)) Supported operators*

Pytorch: *PyTorch Neuron (torch-neuron) Supported operators*

XLA: *neuron-cc-ops-xla*

Apache MXNet: *Neuron Apache MXNet Supported operators*

## Known issues and limitations - updated 11/23/2022

- There is a known issue of increased latency and lower throughput when MLM head is compiled along with BERT model. The workaround is to compile them separately and feed the raw Bert into the head.
- *TensorFlow 2.x* - In this release supported operators are limited to BERT-like models, specifically no conv2d or reduce-window operators are available.
- *Control flow* Neuron only supports control flow operators which are static at compile time. For example static length RNN, top-k, sort.
- *Data layout* The Neuron compiler supports multiple data layout format (NCHW, NHWC, ...). Non-CNHW input/output data-layouts will require Neuron to insert additional transpose operations, causing a degradation in performance.
- *Primary inputs in NeuronCore Pipeline mode* When a neural network is executed in NeuronCore Pipeline mode, only the first operator in a neural network can receive primary inputs from the host.
- *Reduce data type* INT8 data type is not currently supported by the Neuron compiler.
- *NeuronCore Pipeline*: NeuronCorePipeline mode provides low-latency and high-throughput for small batch sizes. We recommend to start testing with batch=1 and gradually increase batch size to fine tune your model throughput and latency performance.
- *Large input tensors* support varies by model. On some models the large input tensors (eg 1024x1024) may result in lower performance or exceeding hardware or compile-time limits, especially on models where the large input tensor is used by many downstream operators. Workarounds may include use of smaller batch, see [Neuron Batching](#)
- *Conv2d operator* is mapped to Inferentia except for specific cases of extremely large tensors and specific parameters.
- *Conv3d operator* performance is limited when the operator has small number of input channels (< 64).
- FP64 and INT64 input and output tensors are not supported. Please cast to FP32/INT32 in the machine learning framework, prior compiling for Neuron.

## Neuron Compiler release [1.21.0.0]

Date: 12/21/2023

- Minor bug fixes.

## Neuron Compiler release [1.20.3.0]

Date: 10/26/2023

- Minor bug fixes.

### Neuron Compiler release [1.19.0.0]

Date: 09/15/2023

- Minor bug fixes.

### Neuron Compiler release [1.17.0.0]

Date: 7/19/2023

#### New in this release

- This release introduces a new `--enable-saturate-infinity` compiler option. A computation that can generate +/- infinity is at a high risk of generating Not-a-Number (NaN) values when the infinity value is used in subsequent computations. This option helps avoid this by converting +Inf/-Inf values to MAX/MIN\_FLOAT before operations that could produce NaN values for +Inf/-Inf inputs on the target architecture. While this option helps to avoid NaN values, there is a potential performance degradation that occurs during model execution when this conversion is enabled.
- Minor bug fixes.

### Neuron Compiler release [1.16.2.0]

Date: 6/14/2023

- Minor bug fixes.

### Neuron Compiler release [1.15.0.0]

Date: 05/01/2023

- Minor bug fixes.

### Neuron Compiler release [1.14.3.0]

Date: 04/19/2023

- Minor bug fixes.

### Neuron Compiler release [1.13.3.0]

Date: 11/23/2022

- Resolved long compile-times when compiling the YOLOv5 and YOLOv6 models. [GitHub · aws-neuron-sdk · #434]
- Improved the layout algorithm to resolve an issue compiling a transformer-based text recognition model. [GitHub · aws-neuron-sdk · #410]
- Support was added for additional XLA operators

**Neuron Compiler release [1.11.7.0]**

Date: 08/02/2022

- Fixed a bug for correct handling of mxnet dropout instruction when mode is set as 'training' while performing inference.

**Neuron Compiler release [1.11.4.0]**

Date: 04/29/2022

- Solved an issue that caused a "false positive" reporting of a data race that may occur due to address overlap.
- Minor bug fixes.

**Neuron Compiler release [1.10.3.0]**

Date: 03/25/2022

- Minor bug fixes.

**Neuron Compiler release [1.9.1.0]**

Date: 01/20/2022

- Fixed an issue with frontend compiler for fused operators that was reported in [github #362](#).

**Neuron Compiler release [1.8.5.0]**

Date: 01/05/2022

**New in this release**

- Minor bug fixes.

**Neuron Compiler release [1.8.2.0]**

Date: 12/15/2021

**New in this release**

- Performance enhancements as a result of improved layout and DMA optimizations.
- Minor bug fixes.

### Neuron Compiler release [1.7.3.0]

Date: 10/27/2021

#### New in this release

- The compiler's `list-operators` command can now display the supported TensorFlow 2.x operators.
- Support added for new operators in TensorFlow 1.x - ArgMax and ArgMin.
- Introducing the `--fast-math` option for better fine-tuning of accuracy/performance. See [Mixed precision and performance-accuracy tuning \(neuron-cc\)](#)

### [1.6.13.0]

Date 08/12/2021

#### New in this release

- TensorFlow 2.x - First support of TensorFlow 2.x. The support is limited to operators in BERT-like models and was tested with Huggingface BERT small, base, large and DistillBert.

#### Resolved issues

- Fixed compiler backend issue in `Tensor_tensor` argument distance, [github #269](#)

### [1.5.5.0]

Date 07/02/2021

#### Summary

- Robustness and performance improvements.

#### New in this release

- Added `--enable-fast-context-switch` option to optimize for faster model switching rather than inference latency.
- Deprecated support for ONNX
- Improved robustness of Conv3d
- Corrected compilation error “too many instructions” in DLRM model



**[1.4.0.0]**

Date 5/28/2021

**Summary**

- Performance improvements, and usability improvements.

**New in this release**

- Added uncompressed NEFF format for faster loading models prior inference. Enable it by `--enable-fast-loading-neuron-binaries`. Some cases of large models may be detrimentally impacted as it will not be compressed but many cases will benefit.
- Corrected compilation error in specific arguments of `ResizeBilinear` operator

**[1.3.0.0]**

Date 4/30/2021

**Summary**

- Performance improvements, new operators, and usability improvements.

**New in this release**

- Improved performance of batched CNN models like resnet50 with the default compiler options by 10%.
- Improved performance of bert base sequence 128 batch 6 by upto 16%
- Added support for group and depth wise convolution (with limited performance when the number of input channels is small).
- Added more detailed debug names to support for tensorboard.

**Resolved Issues**

- Corrected potential race condition in overwriting tiles of output tensors.
- Fixed various issues in pipelined inference by enabling fine grain partitioning by default.

**[1.2.7.0]**

Date 2/24/2021

**Summary**

Fix for CVE-2021-3177.

**[1.2.2.0]**

Date 1/30/2021

**Summary**

Added support for multiple new operators (see operators list) for TensorFlow and MXNET. Improved inference performance of language, object recognition models on single as well as multiple pipelined cores using neuroncore-pipeline.

**New in this release**

- The following models are now supported: Resnext 224x224, specific BERT variations applied to natural language processing and translation.
- A number of new operators is now supported on Inferentia, see the full lists *TensorFlow Neuron (tensorflow-neuron (TF1.x)) Supported operators* and *Neuron Apache MXNet Supported operators*
- Improved inference performance on yolov4 BERT base sequence 64 (on 16 pipelined cores) and openpose 184.

**Resolved Issues**

- Corrected a random failure to compile Resnet50 batch 5
- Corrected numerical inaccuracy in RSQRT and related operators for tensors with very large values ( $> 1e20$ )

**[1.1.7.0]**

Date 12/23/2020

**Summary**

Added support for PyTorch Yolo V4, a new Framework-visible progress bar and improved inference performance. We continue to streamline the compiler usability by removing the need for options passed to control behavior. We are aiming to remove the need for such options entirely. Some tutorials have been updated to reflect this, but Resnet50 remains in need of these options to achieve maximum performance. Other usability improvements have been added, such as the compiler progress bar. As always, please let us know if there are other areas that we can improve.

## New in this release

- Pytorch Yolo V4 is now supported.
- Added a compiler progress bar when compilation is invoked from the Framework. This allows the user to see that progress continues as compilation proceeds, which is useful when compilation takes several minutes. A dot is printed every 20 seconds.
- Improved inference performance of Tensorflow BERT base seq 256 batch 3 by 10% .

## Resolved Issues

- Resolved issue with depthwise convolution that manifests as a type check error

**[1.0.24045.0]**

Date 11/17/2020

## Summary

Improved performance for pipelined execution (NeuronCore Pipeline).

## New in this release

- NeuronCore Pipeline: improved partitioning to enable better static weights loading to cache.

## Resolved Issues

- `--static-weights` : No longer needed. As this is shown in some examples, please remove the option since the compiler now performs this auto-detection by default.
- `--num-neuroncores` renamed to `--neuroncore-pipeline-cores`. The prior option form is still functional (backwards compatible) and will be removed in future releases.
- `--batching_en`: Resolved compilation failure of ResNet50 FP32 batch 1 on Ubuntu16 when “`--batching_en`” was used.

**[1.0.20600.0]**

Date 9/22/2020

### Summary

Various performance improvements - both compilation time and inference speed of object recognition models.

- Compiler optimization ‘-O2’ option is now enabled by default.

### New in this release

- Improved inference performance of YOLO v3, YOLO v4, VGG16, SSD300. BERT models were improved by an additional 10%.
- Modified such that -O2 is now the default behavior and does not need to be specified. Note: some tutorials still explicitly specify “-O2”. These will be modified in forthcoming updates.

### Resolved Issues

- Sped up compilation of large models that were taking hours to sub-40 minute.

### [1.0.18001.0]

Date 8/08/2020

### Summary

Various performance improvements.

### New in this release

Improved performance of BERT base with -O2

### Resolved Issues

- n/a

### [1.0.17937.0]

Date 8/05/2020

## Summary

Various improvements.

**[1.0.16861.0]**

Date 7/16/2020

## Summary

This release has some bug fixes and some functional and performance improvements to support compilation of several neural networks.

## New in this release

This release

- Supports compilation of PoseNet, tested for images of specific resolutions upto 736.
- Update the -O2 with a new memory allocator to reduce spilling to DRAM
- Improved performance of the '-O2' on BERT base, and openpose pose network.

## Resolved Issues

- Resolved compilation error in Vgg16 batch 1

## Other Notes

- Some versions of Inception network may fail to compile in Tensorflow on Ubuntu 16 in conda environment. The symptom is neuron-cc backend data race error. As a workaround use Ubuntu 18, Amazon Linux 2, or virtual env, or use neuron-cc with flag -O2.

|   |
|---|
| <b>Warning:</b> Starting with Neuron 1.14.0, Ubuntu 16 is no longer supported |
|---|

**[1.0.15275.0]**

Date 6/11/2020

## Summary

This release has some bug fixes and some functional and performance improvements to support compilation of several neural networks.

## New in this release

This release

- Supports compilation of PoseNet for images of specific resolutions upto 400x400.
- Improves performance of resnet152.
- Supports a new command line option ‘-O2’ that can help with handling of large tensor inputs for certain models.
- increase NEFF versions to 1.0. This means new NEFFs compiled from this release forward are not compatible with older versions of Neuron Runtime prior to May, 2020 (1.0.6905.0) release. Please update the Neuron Runtime when using NEFF version 1.0.

## Resolved Issues

- Compilation issues on prosotron encoder, decoder neural networks.

## Other Notes

## Dependencies

- This version creates NEFF 1.0 thus may require update of neuron-rtd if older than May 2020 release.

dmlc\_nnm==1.0.2574.0      dmlc\_topi==1.0.2574.0      dmlc\_tvm==1.0.2574.0      inferentia\_hwm==1.0.1362.0  
islpy==2018.2

## [1.0.12696.0]

Date 5/11/2020

## Summary

Bug fixes and some functional and performance improvements to several neural networks.

## New in this release

- This version supports compilation of unmodified Tensorflow BERT with batch size 1, 4, 6 for input sequence 128.
- Improved Tensorflow BERT batch 4 sequence 128 performance to 45% of the accelerator peak (from 34%).
- Support for MXNET BERT base batch 8 compilation
- Support for TF Resnet152 batch 2 compilation
- Most compiler messages are migrated from cout to logging mechanisms with verbosity control

## Resolved Issues

- Fixed failure to compile unmodified Tensorflow BERT model for small batches
- Fixed run-to-run-variability in OneHot operator implementation
- Robustness improvements for ParallelWavenet and transformer decoder networks

## Other Notes

## Dependencies

```
dmlc_nnmvm==1.0.2356.0  
dmlc_topi==1.0.2356.0  
dmlc_tvm==1.0.2356.0  
inferentia_hwm==1.0.1294.0  
islpy==2018.2
```

## [1.0.9410.0]

Date 3/26/2020

## Summary

Bug fixes and some functional and performance improvements to several neural networks.

## New in this release

- Support compilation of modified SSD-300 (tensorflow-ssd300)
- Improved inference performance in natural language processing networks (such as prosotron encoder) by 45%

## Resolved Issues

- Eliminated redundant fp32 to bfloat16 cast on input and output tensors

## Known issues and limitations

- See previous releases.

### Other Notes

- Added support for faster iteration on recurrent networks (aka auto-loop)

### Dependencies

```
dmlc_nnvvm==1.0.2049.0
dmlc_topi==1.0.2049.0
pip install --upgrade dmlc_tvm==1.0.2049.0
inferentia_hwm==1.0.897.0
islpy==2018.2
```

[1.0.7878.0]

Date 2/27/2020

### Summary

Bug fixes and minor performance improvements.

### New in this release

None

### Resolved Issues

- Corrected image resize operator functionality
- Compiler internal enhancements made that will benefit models such as BERT

### Known issues and limitations

- See previous releases.

### Other Notes

### Dependencies

```
dmlc_nnvvm-1.0.1826.0
dmlc_topi-1.0.1826.0
dmlc_tvm-1.0.1826.0
inferentia_hwm-1.0.897.0
islpy-2018.2
```



**[1.0.6801.0]**

Date 1/27/2020

**Summary**

Bug fixes and some performance enhancement related to data movement for BERT-type neural networks.

**New in this release**

None

**Resolved Issues**

- Improved throughput for operators processed in the Neuron Runtime CPU. As an example: execution of 4 single NeuronCore NEFF models of ResNet50 v2 float16 batch = 5 in parallel on an inf1.1xlarge sped up by 30%.
- Corrected shape handling in Gather(TensorFlow)/Take(MXNet) operators that are processed by the Neuron Runtime in the Neuron Runtime vCPU, which resolves a possible crash in Neuron Compiler when compiling models with these operators with some shapes.
- Added support for TensorFlow *OneHot* operator (as a Neuron Runtime CPU operator).
- Added more internal checking for compiler correctness with newly defined error messages for this case.

“Internal ERROR: Data race between Op1 'Name1(...) [...]' and Op2 'Name2(...) [...]'”

- Fixed out-of-memory issue introduced in 1.0.5939.0 such that some large models (BERT) compiled on instances with insufficient host memory would cause the runtime to crash with an invalid NEFF. This is actually a compiler error, but due to additional script layers wrapping this in the [\[Broken\] Running TensorFlow BERT-Large with AWS Neuron](#), this would have likely been seen as a runtime error like this:

```
2020-01-09 13:40:26.002594: E tensorflow/core/framework/op_segment.cc:54] Create kernel_
↳failed: Invalid argument: neff is invalid
2020-01-09 13:40:26.002637: E tensorflow/core/common_runtime/executor.cc:642] Executor_
↳failed to create kernel. Invalid argument: neff is invalid
[[{{node bert/NeuronOp}}]]
```

**Known issues and limitations**

See previous release notes. Some tutorials show use of specific compiler options and flags, these are needed to help provide guidance to the compiler to achieve best performance in specific cases. Please do not use in cases other than as shown in the specific tutorial as results may not be defined. These options should be considered beta and will be removed over time.

### Other Notes

### Dependencies

```
dmlc_nnmvm-1.0.1619.0  
dmlc_topi-1.0.1619.0  
dmlc_tvm-1.0.1619.0  
inferentia_hwm-1.0.839.0  
islpy-2018.2
```

[1.0.5939.0]

Date 12/20/2019

### Summary

Bug fixes and some performance enhancement for NeuronCore Pipeline.

### New in this release

### Resolved Issues

- Fixed pipeline execution on more than 10 NeuronCores
- Improved NeuronCores Pipeline execution by improving data exchange efficiency between NeuronCores
- Added warning for unaligned memory access
- Fixed handling of cast on input FP32 tensor
- Improved handling of data layouts and transpose
- Improved dead-code elimination
- Improved efficiency of compute engine synchronization
- Improved efficiency of data transfers within the Neuron code

### Known issues and limitations

See previous release notes. Some tutorials show use of specific compiler options and flags, these are needed to help provide guidance to the compiler to achieve best performance in specific cases. Please do not use in cases other than as shown in the specific tutorial as results may not be defined. These options should be considered beta and will be removed over time.

## Other Notes

### Dependencies

- dmlc\_nnm-1.0.1416.0
- dmlc\_topi-1.0.1416.0
- dmlc\_tvm-1.0.1416.0
- inferentia\_hwm-1.0.720.0
- islpy-2018.2

### [1.0.5301.0]

Date 12/1/2019

## Summary

### New in this release

### Resolved Issues

- Added warning for unsupported operators and convolution sizes
- Added warning for unsupported layout / upsampling
- Added support for Relu6, AddV2, BatchMatmulV2 operators
- Added support for default MXNet outputs in `-io-config`
- Improved performance of batched inference for convolutional networks
- Fixed MatMult column size 1
- Fixed bf16 constant loading
- Fixed Conv2D tile accumulation

### Known Issues and Limitations

See previous release notes. Resolved issues are shown in Resolved Issues.

## Other Notes

Please install g++ on AMIs without g++ pre-installed (i.e. server AMIs):

```
# Ubuntu  
sudo apt-get install -y g++
```

```
# Amazon Linux  
sudo yum install -y gcc-c++
```

Supported Python versions:

- 3.5, 3.6, 3.7

Supported Linux distributions:

- Ubuntu 16, Ubuntu 18, Amazon Linux 2

### Dependencies

- dmlc\_nnvm-1.0.1328.0
- dmlc\_topi-1.0.1328.0
- dmlc\_tvm-1.0.1328.0
- inferentia\_hwm-1.0.674.0
- islpy-2018.2

### [1.0.4680.0]

Date: 11/25/2019

### New in this release

N/A, this is the first release.

### Resolved issues

N/A, this is the first release.

### Known issues and limitations

1. **Control flow** Inferentia has a limited support for control flow. In general, Neuron can only support control flow operators which are static at compile time, i.e. static length RNN, top-k, sort, ...
2. **Size of neural network** The size of neural network is influenced by a) type of neural network (CNN, LSTM, MLP) , b) number of layers, c) sizes of input (dimension of the tensors, batch size, ...). The current Neuron compiler release has a limitation in terms of the size of neural network it could effectively optimize. As a result, we limit CNN models (e.g. ResNet) to have an input size of up to 480x480 FP16, batch size of 4; LSTM models (e.g. GNMT) are limited to a time step limit of up to 900; MLP models (like BERT) are limited up to sequence-length equal 128, batch=8.
3. **Data layout** The Neuron compiler supports multiple data layout formats (NCHW, NHWC, ...). Non-CNHW input/output data-layouts will require Neuron to insert additional *transpose* operations, causing a degradation in performance.
4. **Object detection models** Computer-vision object detection and segmentation models are not supported by the current release.
5. **Reduce data type** INT8 data type is not currently supported by the Neuron compiler.

6. **Tensor residency** When a sub-graph that is executed on the host is communicating with a sub-graph that is executing on Neuron cores, tensors are copied via the communication queues between the host and Inferentia memory for each inference, which may result in end-to-end performance degradation.
7. **Primary inputs in NeuronCore Pipeline mode** When a neural network is executed in NeuronCore Pipeline mode, only the first operator in a neural network can receive primary inputs from the host.

## Other Notes

## Dependencies

- nnvm: dmlc\_nnvm-1.0.1219.0
- topi: dmlc\_topi-1.0.1219.0
- tvn: dmlc\_tvm-1.0.1219.0
- hwm: inferentia\_hwm-1.0.602.0
- islpy: islpy-2018.2+aws2018.x.73.0

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Neuron Supported operators

*This document is relevant for:* Inf1

## TensorFlow Neuron (tensorflow-neuron (TF1.x)) Supported operators

To see a list of supported operators for TensorFlow 1.x, run the following command:

```
neuron-cc list-operators --framework TENSORFLOW
```

## Neuron Compiler Release [1.9.1.0]

Date: 01/20/2022

Added

```
isNan
FusedBatchNormV3
```

### Neuron Compiler Release [1.7.3.0]

Added

ArgMax  
ArgMin

### Neuron Compiler Release [1.6.13.0]

No changes

### Neuron Compiler Release [1.5.5.0]

No changes

### Neuron Compiler Release [1.4.0.0]

No changes

### Neuron Compiler Release [1.3.0.0]

Added

Abs  
Cos  
DepthwiseConv2dNative  
Erf  
Rank  
Sin  
Size

### Neuron Compiler Release [1.2.7.0]

No changes

### Neuron Compiler Release [1.2.2.0]

Added

AdjustContrastv2  
AdjustSaturation  
BroadcastTo  
Cholesky  
Conv2DBackpropInput  
Conv3D  
CropAndResize

(continues on next page)

(continued from previous page)

```

FloorDiv
HSVToRGB
InvertPermutation
L2Loss
Log1p
MatrixBandPart
MatrixDiag
MatrixSetDiag
MatrixTriangularSolve
MaxPool3D
MirrorPad
RGBToHSV
Range
SoftmaxCrossEntropyWithLogits
SquaredDifference
StopGradient
Unpack
UnsortedSegmentSum

```

### Neuron Compiler Release [1.0.24045.0]

Added FloorDiv, Softplus, Unstack

### Neuron Compiler Release [1.0.18001]

No changes

### Neuron Compiler Release [1.0.16764]

Added:

```

LogSoftmax
Neg
ResizeBilinear
ResizeNearestNeighbor

```

### Neuron Compiler Release [1.0.15275]

Added

```

Neg

```

Removed

```

Log

```

(was inadvertently advertised as supported)

### Neuron Compiler Release [1.0.12696]

No changes

### Neuron Compiler Release [1.0.9410]

No changes

### Neuron Compiler Release [1.0.7878]

No changes

### Neuron Compiler Release [1.0.6801]

No changes

### Neuron Compiler Release [1.0.5939]

No changes

### Neuron Compiler Release [1.0.5301]

No changes

### Neuron Compiler Release [1.0.4680.0]

Add  
AddV2  
All  
AvgPool  
BatchMatMul  
BatchMatMulV2  
BatchToSpaceND  
BiasAdd  
Cast  
Ceil  
Concat  
ConcatV2  
Const  
Conv2D  
Equal  
Exp  
ExpandDims  
Fill  
Floor  
FusedBatchNorm  
Greater

(continues on next page)



(continued from previous page)

GreaterEqual  
Identity  
LRN  
LeakyRelu  
Less  
LessEqual  
Log  
LogicalAnd  
LogicalNot  
LogicalOr  
MatMul  
Max  
MaxPool  
Maximum  
Mean  
Min  
Minimum  
Mul  
NoOp  
NotEqual  
Pack  
Pad  
PadV2  
Placeholder  
Pow  
Prod  
RandomUniform  
RealDiv  
Reciprocal  
Relu  
Relu6  
Reshape  
ReverseV2  
Round  
Rsqrt  
Select  
Shape  
Sigmoid  
Sign  
Slice  
Softmax  
SpaceToBatchND  
Split  
SplitV  
Sqrt  
Square  
Squeeze  
StridedSlice  
Sub  
Sum  
Tanh  
Tile

(continues on next page)

(continued from previous page)

Transpose  
ZerosLike

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Neuron Apache MXNet Supported operators

To see a list of supported operators for MXNet, run the following command:

```
neuron-cc list-operators --framework MXNET
```

## Neuron Compiler Release [1.6.13.0]

Added

amp\_cast  
amp\_multicast

## Neuron Compiler Release [1.4.1.0]

No changes

## Neuron Compiler Release [1.4.0.0]

No changes

## Neuron Compiler Release [1.3.0.0]

No changes

## Neuron Compiler Release [1.2.7.0]

No changes

## Neuron Compiler Release [1.2.2.0]

No changes

**Neuron Compiler Release [1.2.0.0]**

Added

```
Deconvolution
LayerNorm
Pad
SwapAxis
_contrib_arange_like
_contrib_interleaved_matmul_encdec_qk
_contrib_interleaved_matmul_encdec_valatt
_contrib_interleaved_matmul_selfatt_qk
_contrib_interleaved_matmul_selfatt_valatt
arctan
broadcast_like
cos
erf
pad
sin
slice_axis
```

**Neuron Compiler Release [1.0.24045.0]**

Added `_contrib_div_sqrt_dim`, `broadcast_axis`

**Neuron Compiler Release [1.0.18001.0]**

No changes

**Neuron Compiler Release [1.0.17937.0]**

No changes

**Neuron Compiler Release [1.0.16861.0]**

Removed `log` (Was erroneously reported as added in previous release. )

**Neuron Compiler Release [1.0.15275]**

Added `log`

### Neuron Compiler Release [1.0.12696]

No changes

### Neuron Compiler Release [1.0.9410]

No changes

### Neuron Compiler Release [1.0.7878]

No changes

### Neuron Compiler Release [1.0.6801]

No changes

### Neuron Compiler Release [1.0.5939]

no changes

### Neuron Compiler Release [1.0.5301]

no changes

### Neuron Compiler Release [1.0.4680.0]

```
Activation
BatchNorm
Cast
Concat
Convolution
Convolution_v1
Dropout
Flatten
FullyConnected
LeakyReLU
Pooling
Pooling_v1
RNN
Reshape
SequenceMask
SliceChannel
Softmax
UpSampling
__add_scalar__
__div_scalar__
__mul_scalar__
```

(continues on next page)

(continued from previous page)

```
__pow_scalar__
__rdiv_scalar__
__rpow_scalar__
__rsub_scalar__
__sub_scalar__
_arange
_copy
_div_scalar
_equal_scalar
_full
_greater_equal_scalar
_greater_scalar
_lesser_equal_scalar
_lesser_scalar
_maximum
_maximum_scalar
_minimum
_minimum_scalar
_minus_scalar
_mul_scalar
_not_equal_scalar
_ones
_plus_scalar
_power_scalar
_rdiv_scalar
_rminus_scalar
_rnn_param_concat
_zeros
batch_dot
broadcast_add
broadcast_div
broadcast_equal
broadcast_greater
broadcast_greater_equal
broadcast_lesser
broadcast_lesser_equal
broadcast_maximum
broadcast_minimum
broadcast_mod
broadcast_mul
broadcast_not_equal
broadcast_sub
ceil
clip
concat
elemwise_add
elemwise_div
elemwise_mul
elemwise_sub
exp
expand_dims
flatten
```

(continues on next page)

(continued from previous page)

```
floor
gather_nd
log
log_softmax
max
mean
min
negative
ones_like
relu
repeat
reshape
reshape_like
reverse
rsqrt
sigmoid
slice
slice_like
softmax
split
sqrt
square
squeeze
stack
sum
tanh
tile
transpose
where
zeros_like
```

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## Neuron Compiler for Trn1 & Inf2

### API Reference Guide

- [Neuron Compiler CLI Reference Guide](#)

### Developer Guide

- [Mixed Precision and Performance-accuracy Tuning \(neuronx-cc\)](#)

### Misc

- [FAQ](#)
- [What's New](#)

## Neuron Compiler for Inf1

### API Reference Guide

- [Neuron compiler CLI Reference Guide \(neuron-cc\)](#)

### Developer Guide

- [Mixed precision and performance-accuracy tuning \(neuron-cc\)](#)

### Misc

- [FAQ](#)
- [What's New](#)
- [Neuron Supported operators](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 7.2 Neuron Kernel Interface (NKI) - Beta

Neuron Kernel Interface (NKI) is a bare-metal language and compiler for directly programming NeuronDevices available on AWS Trn/Inf instances. You can use NKI to develop, optimize and run new operators directly on NeuronCores while making full use of available compute and memory resources. NKI empowers ML developers to self-serve and invent new ways to use the NeuronCore hardware, starting NeuronCores v2 (Trainium1) and beyond.

NKI provides developers with direct access to the NeuronCore ISA (Instruction Set Architecture), accessible from a Python-based programming environment, which has syntax and tile-level semantics that are similar to [Triton](#) and [NumPy](#). This enables developers to get started quickly and optimize performance in a familiar environment, while at the same time get full control of the underlying hardware. At the hardware level, NeuronCore's tensorized memory access capability enables efficient reading and writing of multi-dimensional arrays on a per instruction basis, which makes NKI's tile-based programming highly suitable for the NeuronCore instruction set.

For comparison, before NKI was introduced, the only way to program NeuronDevices was through defining high-level ML models in frameworks such as [PyTorch](#) and [JAX](#). Neuron Compiler takes such high-level model definitions as input, performs multiple rounds of optimization, and eventually generates a NEFF (Neuron Executable File Format) that is executable on NeuronDevices. At a high level, Neuron Compiler runs the following optimization stages in order:

1. **Hardware-agnostic graph-level optimizations.** These transformations are done in the compiler front-end, using [XLA](#), including optimizations like constant propagation, re-materialization and operator fusion.
2. **Loop-level optimization.** Compiler turns the optimized graph from Step 1 into a series of loop nests and performs layout, tiling and loop fusion optimizations.
3. **Hardware intrinsics mapping.** Compiler maps the architecture-agnostic loop nests from Step 2 into architecture-specific instructions.
4. **Hardware-specific optimizations.** These optimizations are mainly done at the instruction level<sup>1</sup> in compiler back-end, with a key goal of reducing memory pressure and improving instruction-level parallelism. For example, memory allocation and instruction scheduling are done in this stage.

NKI kernels bypass the first 3 steps, and are compiled into IRs (intermediate representations) that the compiler's back-end (Step 4 above) can directly consume. Advanced features in NKI, such as direct allocation, also allow programmers to bypass certain compiler passes in Step 4. As a result, NKI developers can now have great control over NeuronDevices down to the instruction level. We highly recommend developers to study the underlying hardware architecture before optimizing performance of their NKI kernels. See the NKI guide below to learn more!

## 7.2.1 Guide

NKI guide is organized in four parts:

1. [API Reference Guide](#) has the NKI API reference manual.
2. [Writing Functional NKI Kernels](#) includes guides that are designed for NKI beginners to learn NKI key concepts and implement kernels to meet functionality requirements.
3. [Writing Performant NKI Kernels](#) includes a deep dive of NeuronDevice architecture and programmer's guides to optimize performance of NKI kernels.
4. [General Resources](#) include any miscellaneous guides.

### API Reference Guide

NKI API Reference Manual

### Writing Functional NKI Kernels

Getting Started with NKI  
NKI Tutorials

NKI Programming Model  
NKI Kernels

NKI Kernel as Framework Custom-Operator

---

<sup>1</sup> A small number of loop-level optimizations are performed after hardware intrinsic mappings in the current Beta release. Subject to future changes.



## Writing Performant NKI Kernels

NeuronDevice Architecture Guide

Profiling NKI kernels with Neuron Profile  
Direct Allocation Developer Guide

NKI Performance

## General Resources

NKI FAQ

NKI What's New

NKI Known Issues

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI API Reference Manual

Summary of different NKI API sets:

- **nki** top-level module contains APIs to decorate and simulate NKI kernels as well as NKI object types.
- **nki.language** consists of high-level compute and data movement APIs designed for ease-of-use. `nki.language` allows NKI programmers to transition from NumPy/Triton implementation to NKI quickly without the need to *fully* understand underlying NeuronDevice architecture. Most language APIs invoke one or more `nki.isa` APIs (that is, NeuronDevice hardware instructions) under the hood.
- **nki.isa** consists of low-level APIs that highly resemble hardware instructions in NeuronDevice ISA (instruction set architecture) designed to provide fine control over the hardware. These APIs expose all the programmable input parameters of the corresponding hardware instructions and also enforce the same tile-size and layout requirements as specified in NeuronDevice ISA.
- **nki.compiler** consists of features that control the compilation process of a NKI kernel.
- Other documents:
  - **NKI API Common Fields** documents common NKI API input parameters such as data types and masks, as well as common API behavior such as type promotion.
  - **NKI API Errors** captures common error types that are thrown by the NKI kernel compilation *frontend*, including syntax, tile-size and layout violation errors.

*This document is relevant for: Inf2, Trn1, Trn2*

## nki

### Decorators

|                              |   |
|------------------------------|---|
| <code>jit</code>             | This decorator compiles a function to run on NeuronDevices.   |
| <code>benchmark</code>       | Benchmark a NKI kernel on a NeuronDevice by using <code>nki.benchmark</code> as a decorator.          |
| <code>profile</code>         | Profile a NKI kernel on a NeuronDevice by using <code>nki.profile</code> as a decorator.              |
| <code>baremetal</code>       | Compile and run a NKI kernel on NeuronDevice without involving ML frameworks such as PyTorch and JAX. |
| <code>simulate_kernel</code> | Simulate a nki kernel on CPU using a built-in simulator in Neuron Compiler.                           |

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.jit

`nki.jit(func=None, mode='auto', **kwargs)`

This decorator compiles a function to run on NeuronDevices.

This decorator tries to automatically detect the current framework and compile the function as a custom operator of the current framework. To bypass the framework detection logic, you may specify the mode parameter explicitly.

### Parameters

- **func** – The function that define the custom op
- **mode** – The compilation mode, possible values: “jax”, “torchxla”, “baremetal”, “benchmark”, “simulation” and “auto”

Listing 7.1: An Example

```

from neuronxcc import nki
import neuronxcc.nki.language as nl

@nki.jit
def nki_tensor_tensor_add(a_tensor, b_tensor):
    c_tensor = nl.ndarray(a_tensor.shape, dtype=a_tensor.dtype, buffer=nl.shared_hbm)

    a = nl.load(a_tensor)
    b = nl.load(b_tensor)

    c = a + b

    nl.store(c_tensor, c)

    return c_tensor

```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.benchmark

`nki.benchmark(kernel=None, **kwargs)`

Benchmark a NKI kernel on a NeuronDevice by using `nki.benchmark` as a decorator. You must run this API on a Trn/Inf instance with NeuronDevices (v2 or beyond) attached and also `aws-neuronx-tools` installed on the host using the following steps:

```

# on Ubuntu
sudo apt-get install aws-neuronx-tools=2.* -y

# on Amazon Linux
sudo yum install aws-neuronx-tools-2.* -y

```

You may specify a path to save your NEFF file through input parameter `save_neff_name` and a path to save your NTFF file through `save_trace_name`. See [Profiling NKI kernels with Neuron Profile](#) for more information on how to visualize the execution trace for profiling purposes.

**Note:** Similar to `nki.baremetal`, The decorated function using `nki.benchmark` expects `numpy.ndarray` as

input/output tensors instead of ML framework tensor objects.

In addition to generating NEFF/NTFF files, this decorator also invokes `neuron-bench` to collect execution latency statistics of the NEFF file and prints the statistics to the console.

`neuron-bench` is a tool that launches the NEFF file on a `NeuronDevice` in a loop to collect end-to-end latency statistics. You may specify the number of warm-up iterations to skip benchmarking in input parameter `warmup`, and the number of benchmarking iterations in `iters`. Currently, `nki.benchmark` only supports benchmarking on a single `NeuronCore`, since NKI not yet supports collective compute. Note, `neuron-bench` measures not only the device latency but also the time taken to transfer data between host and device. However, the tool does not rely on any ML framework to launch the NEFF and therefore reports NEFF latency without any framework overhead.

#### Parameters

- **warmup** – The number of iterations for warmup execution (10 by default).
- **iters** – The number of iterations for benchmarking (100 by default).
- **save\_neff\_name** – Save the compiled neff file if specify a name (unspecified by default).
- **save\_trace\_name** – Save the trace (profile) file if specified a name (unspecified by default); at the moment, it requires that the `save_neff_name` is unspecified or specified as `'file.neff'`.
- **additional\_compile\_opt** – Additional Neuron compiler flags to pass in when compiling the kernel.

#### Returns

A function object that wraps the decorating function. A property `benchmark_result.nc_latency` is available after invocation. `get_latency_percentile(int)` of the property returns the specified percentile latency in microsecond(us). Available percentiles: [0, 1, 10, 25, 50, 90, 99, 100]

Listing 7.2: An Example

```
from neuronxcc.nki import benchmark
import neuronxcc.nki.language as nl
import numpy as np

@benchmark(warmup=10, iters = 100, save_neff_name='file.neff', save_trace_name=
↳ 'profile.ntff')
def nki_tensor_tensor_add(a_tensor, b_tensor):
    c_tensor = nl.ndarray(a_tensor.shape, dtype=a_tensor.dtype, buffer=nl.shared_hbm)

    a = nl.load(a_tensor)
    b = nl.load(b_tensor)

    c = a + b

    nl.store(c_tensor, c)

    return c_tensor

a = np.zeros([128, 1024], dtype=np.float32)
b = np.random.random_sample([128, 1024]).astype(np.float32)
c = nki_tensor_tensor_add(a, b)

metrics = nki_tensor_tensor_add.benchmark_result.nc_latency
```

(continues on next page)

(continued from previous page)

```
print("latency.p50 = " + str(metrics.get_latency_percentile(50)))
print("latency.p99 = " + str(metrics.get_latency_percentile(99)))
```

**Note:** `nki.benchmark` does not use the actual inputs passed into the benchmarked function when running the neff file. For instance, in the above example, the output `c` tensor is undefined and should not be used for numerical accuracy checks.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.profile

`nki.profile(func=None, **kwargs)`

Profile a NKI kernel on a NeuronDevice by using `nki.profile` as a decorator.

**Note:** Similar to `nki.baremetal`, The decorated function using `nki.benchmark` expects `numpy.ndarray` as input/output tensors instead of ML framework tensor objects.

### Parameters

- **working\_directory** – A path to working directory where profile artifacts are saved. This must be specified and must also be an absolute path.
- **save\_neff\_name** – Name of the saved neff file if specified (file.neff by default).
- **save\_trace\_name** – Name of the saved trace (profile) file if specified (profile.ntff by default)
- **additional\_compile\_opt** – Additional Neuron compiler flags to pass in when compiling the kernel.
- **overwrite** – Overwrite existing profile artifacts if set to True. Default is False.
- **profile\_nth** – Profiles the *profile\_nth* execution. Default is 1.

### Returns

None

Listing 7.3: An Example

```
from neuronxcc import nki
import neuronxcc.nki.language as nl

@nki.profile(working_directory="/home/ubuntu/profiles", save_neff_name='file.neff',
             save_trace_name='profile.ntff')
def nki_tensor_tensor_add(a_tensor, b_tensor):
    c_tensor = nl.ndarray(a_tensor.shape, dtype=a_tensor.dtype, buffer=nl.shared_hbm)

    a = nl.load(a_tensor)
    b = nl.load(b_tensor)

    c = a + b

    nl.store(c_tensor, c)
```

(continues on next page)

(continued from previous page)

```
return c_tensor
```

`nki.profile` will save `file.neff`, `profile.ntff`, along with json files containing a profile summary inside of the `working_directory`.

See [Profiling NKI kernels with Neuron Profile](#) for more information on how to visualize the execution trace for profiling purposes.

In addition, more information about *neuron-profile* can be found in its [documentation](#).

---

**Note:** `nki.profile` does not use the actual inputs passed into the profiled function when running the `neff` file. For instance, in the above example, the output `c` tensor is undefined and should not be used for numerical accuracy checks. The input tensors are used mainly to specify the shape of inputs.

---

*This document is relevant for:* `Inf2`, `Trn1`, `Trn2`

*This document is relevant for:* `Inf2`, `Trn1`, `Trn2`

## `nki.baremetal`

`nki.baremetal(kernel=None, **kwargs)`

Compile and run a NKI kernel on `NeuronDevice` without involving ML frameworks such as PyTorch and JAX. If you decorate your NKI kernel function with decorator `@nki.baremetal(...)`, you may call the NKI kernel function directly just like any other Python function. You must run this API on a `Trn/Inf` instance with `NeuronDevices` (v2 or beyond) attached.

---

**Note:** The decorated function using `nki.baremetal` expects `numpy.ndarray` as input/output tensors instead of ML framework tensor objects.

---

This decorator compiles the NKI kernel into an executable on `NeuronDevices` (NEFF) and also collects an execution trace (NTFF) by running the NEFF on the local `NeuronDevice`. See [Profiling NKI kernels with Neuron Profile](#) for more information on how to visualize the execution trace for profiling purposes.

Since `nki.baremetal` runs the compiled NEFF without invoking any ML framework, it is the fastest way to compile and run any NKI kernel standalone on `NeuronDevice`. Therefore, this decorator is useful for quickly iterating an early implementation of a NKI kernel to reach functional correctness before porting it to the ML framework and injecting the kernel into the full ML model. To iterate over NKI kernel performance quickly, NKI also provides `nki.benchmark` decorator which uses the same underlying mechanism as `nki.baremetal` but additionally collects latency statistics in different percentiles.

### Parameters

- **save\_neff\_name** – A file path to save your NEFF file. By default, this is unspecified, and the NEFF file will be deleted automatically after execution.
- **save\_trace\_name** – A file path to save your NTFF file. By default, this is unspecified, and the NTFF file will be deleted automatically after execution. Known issue: if `save_trace_name` is specified, `save_neff_name` must be set to “file.neff”.
- **additional\_compile\_opt** – Additional Neuron compiler flags to pass in when compiling the kernel.
- **artifacts\_dir** – A directory path to save Neuron compiler artifacts. The directory must be empty before running the kernel. A non-empty directory would lead to a compilation error.

**Returns**

None

Listing 7.4: An Example

```

from neuronxcc.nki import baremetal
import neuronxcc.nki.language as nl
import numpy as np

@baremetal(save_neff_name='file.neff', save_trace_name='profile.ntff')
def nki_tensor_tensor_add(a_tensor, b_tensor):
    c_tensor = nl.ndarray(a_tensor.shape, dtype=a_tensor.dtype, buffer=nl.shared_hbm)

    a = nl.load(a_tensor)
    b = nl.load(b_tensor)

    c = a + b

    nl.store(c_tensor, c)

    return c_tensor

a = np.zeros([128, 1024], dtype=np.float32)
b = np.random.random_sample([128, 1024]).astype(np.float32)
c = nki_tensor_tensor_add(a, b)

assert np.allclose(c, a + b)

```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.simulate\_kernel**

**nki.simulate\_kernel**(kernel, \*args, \*\*kwargs)

Simulate a nki kernel on CPU using a built-in simulator in Neuron Compiler. This simulation mode is especially useful for inspecting intermediate tensor values using *nki.language.device\_print* (see code example below).

---

**Note:** All input and output tensors to the kernel must be `numpy.ndarray` when using this `simulate_kernel` API.

---

To run the kernel on a NeuronCore instead, please refer to *Getting Started with NKI*.

**Parameters**

- **kernel** – The kernel to be simulated
- **args** – The args of the kernel
- **kwargs** – The kwargs of the kernel

**Returns**

Examples:

```

import neuronxcc.nki as nki
import neuronxcc.nki.language as nl
import numpy as np

```

(continues on next page)

(continued from previous page)

```

@nki.jit
def print_kernel():
    a = nl.ndarray([4, 4], dtype=nl.float32, buffer=nl.shared_hbm)

    # Create (4, 4) tensor in sbuf
    y = nl.zeros([4, 4], dtype=np.float32)

    # Print tensor y
    nl.device_print("value of y:", y)

    # Directly store tensor y as a single tile
    nl.store(a, value=y)

    return a

np.random.seed(0)

a = nki.simulate_kernel(print_kernel)

assert np.allclose(a, np.zeros([4, 4]))

```

*This document is relevant for: Inf2, Trn1, Trn2*

## Types

---

*tensor*

A tensor object represents a multidimensional, homogeneous array of fixed-size items

---

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.tensor

**class** `nki.tensor`

A tensor object represents a multidimensional, homogeneous array of fixed-size items

## Methods

|                     |  |
|---------------------|--|
| <i>assert_shape</i> | Assert that the tensor has the given shape.  |
| <i>astype</i>       | Copy of the tensor, cast to a specified type.  |
| <i>broadcast_to</i> | Broadcast tensor to a new shape based on numpy broadcast rules.                          |
| <i>expand_dims</i>  | Gives a new shape to a tensor by adding a dimension of size 1 at the specified position. |
| <i>reshape</i>      | Gives a new shape to an array without changing its data.                                 |
| <i>view</i>         | Return a new view of the tensor, reinterpret to a specified type.                        |

## Attributes

|                 |  |
|-----------------|--|
| <i>dtype</i>    | Data type of the tensor.               |
| <i>itemsize</i> | Length of one tensor element in bytes. |
| <i>ndim</i>     | Number of dimensions of the tensor.    |
| <i>shape</i>    | Shape of the tensor.                   |

### **assert\_shape**(*shape*)

Assert that the tensor has the given shape.

#### **Parameters**

**shape** – The expected shape.

#### **Returns**

The tensor.

### **astype**(*dtype*)

Copy of the tensor, cast to a specified type.

#### **Parameters**

**dtype** – The target dtype

#### **Returns**

the tensor with new type. Copy ALWAYS occur

### **broadcast\_to**(*shape*)

Broadcast tensor to a new shape based on numpy broadcast rules. The tensor object must be a tile or can be implicitly converted to a tile. A tensor can be implicitly converted to a tile iff the partition dimension is the highest dimension.

#### **Parameters**

**shape** – The new shape

#### **Returns**

Return a new view of the tensor, no copy will occur

### **property dtype**

Data type of the tensor.

### **expand\_dims**(*axis*)

Gives a new shape to a tensor by adding a dimension of size 1 at the specified position.

#### **Parameters**

**axis** – the position of the new dimension.

#### **Returns**

Return a new tensor with expanded shape



**property itemsize**

Length of one tensor element in bytes.

**property ndim**

Number of dimensions of the tensor.

**reshape**(*shape*)

Gives a new shape to an array without changing its data.

**Parameters**

**shape** – The new shape

**Returns**

Return a new view of the tensor, no copy will occur

**property shape**

Shape of the tensor.

**view**(*dtype*)

Return a new view of the tensor, reinterpret to a specified type.

**Returns**

A new tensor object refer to the original tensor data, NO copy will occur

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.language****Memory operations**

|                         |   |
|-------------------------|---|
| <i>load</i>             | Load a tensor from device memory (HBM) into on-chip memory (SBUF).  |
| <i>store</i>            | Store into a tensor on device memory (HBM) from on-chip memory (SBUF).                                      |
| <i>load_transpose2d</i> | Load a tensor from device memory (HBM) and 2D-transpose the data before storing into on-chip memory (SBUF). |
| <i>atomic_rmw</i>       | Perform an atomic read-modify-write operation on HBM data <code>dst = op(dst, value)</code>                 |
| <i>copy</i>             | Create a copy of the src tile.  |
| <i>broadcast_to</i>     | Broadcast the <code>src</code> tile to a new shape based on numpy broadcast rules.                          |

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.load

`nki.language.load(src, *, mask=None, dtype=None, **kwargs)`

Load a tensor from device memory (HBM) into on-chip memory (SBUF).

See [Memory hierarchy](#) for detailed information.

### Parameters

- **src** – HBM tensor to load the data from.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

### Returns

a new tile on SBUF with values from src.

```
import neuronxcc.nki.language as nl

# load from in_tensor[P, F] that is on HBM
# copy into data_tile[P, F] that is on SBUF
data_tile = nl.load(in_tensor)
...
```

---

**Note:** Partition dimension size can't exceed the hardware limitation of `nki.language.tile_size.pmax`, see [Tile size considerations](#).

---

Partition dimension has to be the first dimension in the index tuple of a tile. Therefore, data may need to be split into multiple batches to load/store, for example:

```
import neuronxcc.nki.language as nl

for i_b in nl.affine_range(4):
    data_tile = nl.zeros((128, 512), dtype=in_tensor.dtype)
    # load from in_tensor[4, 128, 512] one batch at a time
    # copy into data_tile[128, 512]
    i_p, i_f = nl.mgrid[0:128, 0:512]
    data_tile[i_p, i_f] = nl.load(in_tensor[i_b, i_p, i_f])
...
```

Also supports indirect DMA access with dynamic index values:

```
import neuronxcc.nki.language as nl
...

#####
↪ #####
# Indirect DMA read example 1:
# - data_tensor on HBM has shape [128 x 512].
# - idx_tensor on HBM has shape [64] (with values [0, 2, 4, 6, ...]).
# - idx_tensor values read from HBM and stored in SBUF idx_tile of shape [64 x 1]
# - data_tensor values read from HBM indexed by values in idx_tile
```

(continues on next page)

(continued from previous page)

```

# and store into SBUF data_tile of shape [64 x 512].
#####
↳#####
i_p = nl.arange(64)[: , None]
i_f = nl.arange(512)[None, :]

idx_tile = nl.load(idx_tensor[i_p]) # indices have to be in SBUF
data_tile = nl.load(data_tensor[idx_tile[i_p, 0], i_f])
...

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
↳#####
# Indirect DMA read example 2:
# - data_tensor on HBM has shape [128 x 512].
# - idx_tile on SBUF has shape [64 x 1] (with values [[0], [2], [4], ...] generated
↳by iota)
# - data_tensor values read from HBM indexed by values in idx_tile
# and store into SBUF data_tile of shape [64 x 512].
#####
↳#####
i_f = nl.arange(512)[None, :]

idx_expr = 2*nl.arange(64)[: , None]
idx_tile = nisa.iota(idx_expr, dtype=np.int32)
data_tile = nl.load(data_tensor[idx_tile, i_f])
...

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.store

`nki.language.store(dst, value, *, mask=None, **kwargs)`

Store into a tensor on device memory (HBM) from on-chip memory (SBUF).

See [Memory hierarchy](#) for detailed information.

### Parameters

- **dst** – HBM tensor to store the data into.
- **value** – An SBUF tile that contains the values to store.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

```

import neuronxcc.nki.language as nl
...

```

(continues on next page)

(continued from previous page)

```
# store into out_tensor[P, F] that is on HBM
# from data_tile[P, F] that is on SBUF
nl.store(out_tensor, data_tile)
```

**Note:** Partition dimension size can't exceed the hardware limitation of `nki.language.tile_size.pmax`, see *Tile size considerations*.

Partition dimension has to be the first dimension in the index tuple of a tile. Therefore, data may need to be split into multiple batches to load/store, for example:

```
import neuronxcc.nki.language as nl

for i_b in nl.affine_range(4):
    data_tile = nl.zeros((128, 512), dtype=in_tensor.dtype)

    ...
    # store into out_tensor[4, 128, 512] one batch at a time
    # from data_tile[128, 512]
    i_p, i_f = nl.mgrid[0:128, 0:512]
    nl.store(out_tensor[i_b, i_p, i_f], value=data_tile[i_p, i_f])
```

Also supports indirect DMA access with dynamic index values:

```
import neuronxcc.nki.language as nl
...

#####
# Indirect DMA write example 1:
# - data_tensor has shape [128 x 512].
# - idx_tensor on HBM has shape [64] (with values [0, 2, 4, 6, ...]).
# - idx_tensor values read from HBM and stored in SBUF idx_tile.
# - data_tile of shape [64 x 512] values written into
#   HBM data_tensor indexed by values in idx_tile.
#####
i_p = nl.arange(64)[: , None]
i_f = nl.arange(512)[None, :]
idx_tile = nl.load(idx_tensor[i_p]) # indices have to be in SB

nl.store(data_tensor[idx_tile[i_p, 0], i_f], value=data_tile[0:64, 0:512])
```

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
↪#####
# Indirect DMA write example 2:
# - data_tensor has shape [128 x 512].
# - idx_tile on SBUF has shape [64 x 1] (with values [[0], [2], [4], ...]) ↪
```

(continues on next page)

(continued from previous page)

```

→generated by iota)
# - data_tile of shape [64 x 512] values written into
#   HBM data_tensor indexed by values in idx_tile.
#####
→#####
idx_expr = 2*nl.arange(64)[: , None]
idx_tile = nisa.iota(idx_expr, dtype=np.int32)

nl.store(data_tensor[idx_tile, i_f], value=data_tile[0:64, 0:512])

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.load\_transpose2d

`nki.language.load_transpose2d(src, *, mask=None, dtype=None, **kwargs)`

Load a tensor from device memory (HBM) and 2D-transpose the data before storing into on-chip memory (SBUF).

### Parameters

- **src** – HBM tensor to load the data from.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

### Returns

a new tile on SBUF with values from src 2D-transposed.

```

import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor
...

# load from in_tensor[F, P] that is on HBM
# transpose and copy into local_tile[P, F] that is on SBUF
N, M = in_tensor.shape
local_tile: tensor[M, N] = nl.load_transpose2d(in_tensor)
...

```

---

**Note:** Partition dimension size can't exceed the hardware limitation of `nki.language.tile_size.pmax`, see [Tile size considerations](#).

---

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.atomic\_rmw

`nki.language.atomic_rmw(dst, value, op, *, mask=None, **kwargs)`

Perform an atomic read-modify-write operation on HBM data `dst = op(dst, value)`

### Parameters

- **dst** – HBM tensor with subscripts, only supports indirect dynamic indexing currently.
- **value** – tile or scalar value that is the operand to `op`.
- **op** – atomic operation to perform, only supports `np.add` currently.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

```
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor
...

value: tensor[N, M] = nl.load(value_tensor)

# dynamic indices have to be in SBUF, with shape [N, 1]
indices_tile: tensor[N, 1] = nl.load(indices_tensor)

ix = nl.arange(M)[None, :]

#####
# Atomic read-modify-write example:
# - read: values of rmw_tensor is indexed by values from indices_tile
# - modify: incremented by value
# - write: saved back into rmw_tensor
# resulting in rmw_tensor = rmw_tensor + value
#####
nl.atomic_rmw(rmw_tensor[indices_tile, ix], value=value, op=np.add)
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.copy

`nki.language.copy(src, *, mask=None, dtype=None, **kwargs)`

Create a copy of the src tile.

### Parameters

- **src** – the source of copy, must be a tile in SBUF or PSUM.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

### Returns

a new tile with the same layout as `src`, this new tile will be in SBUF, but can be also assigned to a PSUM tensor.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.broadcast\_to**

`nki.language.broadcast_to(src, *, shape, **kwargs)`

Broadcast the `src` tile to a new shape based on numpy broadcast rules. The `src` may also be a tensor object which may be implicitly converted to a tile. A tensor can be implicitly converted to a tile if the partition dimension is the outermost dimension. If `src.shape` is already the same as `shape`, this operation will simply return `src`.

**Parameters**

- **src** – the source of broadcast, a tile in SBUF or PSUM. May also be a tensor object.
- **shape** – the target shape for broadcasting.

**Returns**

a new tile broadcast along the partition dimension of `src`, this new tile will be in SBUF, but can be also assigned to a PSUM tensor.

```
import neuronxcc.nki.language as nl

#####
# Example 1: Load from in_tensor[P, F] that is on HBM and
# copy into out_tile[P, F] that is on SBUF by broadcasting
#####
...

...
# broadcast into out_tile[P, F] that is on SBUF
# from data_tile[P, F] that is on SBUF
in_tile = nl.load(in_tensor, dtype=in_tensor.dtype)
out_tile = nl.broadcast_to(in_tile, shape=(128, in_tensor.shape[1]))

# store output
nl.store(out_tensor, out_tile)
```

*This document is relevant for: Inf2, Trn1, Trn2*

## Creation operations

|                                     |   |
|-------------------------------------|---|
| <code>ndarray</code>                | Create a new tensor of given shape and dtype on the specified buffer.   |
| <code>empty_like</code>             | Create a new tensor with the same shape and type as a given tensor.   |
| <code>zeros</code>                  | Create a new tensor of given shape and dtype on the specified buffer, filled with zeros.  |
| <code>zeros_like</code>             | Create a new tensor of zeros with the same shape and type as a given tensor.  |
| <code>ones</code>                   | Create a new tensor of given shape and dtype on the specified buffer, filled with ones.   |
| <code>full</code>                   | Create a new tensor of given shape and dtype on the specified buffer, filled with initial value.                                  |
| <code>rand</code>                   | Generate a tile of given shape and dtype, filled with random values that are sampled from a uniform distribution between 0 and 1. |
| <code>random_seed</code>            | Sets a seed, specified by user, to the random number generator on HW.   |
| <code>shared_constant</code>        | Create a new tensor filled with the data specified by data array.   |
| <code>shared_identity_matrix</code> | Create a new identity tensor with specified data type.  |

*This document is relevant for: Inf2, Trn1, Trn2*

## `nki.language.ndarray`

`nki.language.ndarray(shape, dtype, *, buffer=None, name="", **kwargs)`

Create a new tensor of given shape and dtype on the specified buffer.

((Similar to `numpy.ndarray`))

### Parameters

- **shape** – the shape of the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, `sbuf`, `psum`, `hbm`), defaults to `sbuf`.
- **name** – the name of the tensor.

### Returns

a new tensor allocated on the buffer.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*



## nki.language.empty\_like

`nki.language.empty_like(a, dtype=None, *, buffer=None, name="", **kwargs)`

Create a new tensor with the same shape and type as a given tensor.

((Similar to [numpy.empty\\_like](#)))

### Parameters

- **a** – the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, *sbuf*, *psum*, *hbm*), defaults to *sbuf*.
- **name** – the name of the tensor.

### Returns

a tensor with the same shape and type as a given tensor.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.zeros

`nki.language.zeros(shape, dtype, *, buffer=None, name="", **kwargs)`

Create a new tensor of given shape and dtype on the specified buffer, filled with zeros.

((Similar to [numpy.zeros](#)))

### Parameters

- **shape** – the shape of the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, *sbuf*, *psum*, *hbm*), defaults to *sbuf*.
- **name** – the name of the tensor.

### Returns

a new tensor allocated on the buffer.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.zeros\_like

`nki.language.zeros_like(a, dtype=None, *, buffer=None, name="", **kwargs)`

Create a new tensor of zeros with the same shape and type as a given tensor.

((Similar to [numpy.zeros\\_like](#)))

### Parameters

- **a** – the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, *sbuf*, *psum*, *hbm*), defaults to *sbuf*.
- **name** – the name of the tensor.

### Returns

a tensor of zeros with the same shape and type as a given tensor.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.ones

`nki.language.ones(shape, dtype, *, buffer=None, name="", **kwargs)`

Create a new tensor of given shape and dtype on the specified buffer, filled with ones.

((Similar to [numpy.ones](#)))

### Parameters

- **shape** – the shape of the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, *sbuf*, *psum*, *hbm*), defaults to *sbuf*.
- **name** – the name of the tensor.

### Returns

a new tensor allocated on the buffer.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.full

`nki.language.full(shape, fill_value, dtype, *, buffer=None, name="", **kwargs)`

Create a new tensor of given shape and dtype on the specified buffer, filled with initial value.

((Similar to [numpy.full](#)))

### Parameters

- **shape** – the shape of the tensor.
- **fill\_value** – the initial value of the tensor.
- **dtype** – the data type of the tensor (see [Supported Data Types](#) for more information).
- **buffer** – the specific buffer (ie, *sbuf*, *psum*, *hbm*), defaults to *sbuf*.
- **name** – the name of the tensor.

### Returns

a new tensor allocated on the buffer.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.rand

`nki.language.rand(shape, dtype=<class 'numpy.float32'>, **kwargs)`

Generate a tile of given shape and dtype, filled with random values that are sampled from a uniform distribution between 0 and 1.

### Parameters

- **shape** – the shape of the tile.
- **dtype** – the data type of the tile (see [Supported Data Types](#) for more information).

### Returns

a tile with random values.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.random\_seed

`nki.language.random_seed(seed, *, mask=None, **kwargs)`

Sets a seed, specified by user, to the random number generator on HW. Using the same seed will generate the same sequence of random numbers when using together with the `random()` API

### Parameters

- **seed** – a 32-bit scalar value to use as the seed.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

none

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.shared\_constant

`nki.language.shared_constant(constant, dtype=None, **kwargs)`

Create a new tensor filled with the data specified by data array.

### Parameters

**constant** – the constant data to be filled into a tensor

### Returns

a tensor which contains the constant data

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.shared\_identity\_matrix

`nki.language.shared_identity_matrix(n, dtype=<class 'numpy.uint8'>, **kwargs)`

Create a new identity tensor with specified data type.

This function has the same behavior to [nki.language.shared\\_constant](#) but is preferred if the constant matrix is an identity matrix. The compiler will reuse all the identity matrices of the same dtype in the graph to save space.

### Parameters

- **n** – the number of rows(and columns) of the returned identity matrix
- **dtype** – the data type of the tensor, default to be `np.uint8` (see [Supported Data Types](#) for more information).

### Returns

a tensor which contains the identity tensor

*This document is relevant for: Inf2, Trn1, Trn2*

## Math operations

|                        |   |
|------------------------|---|
| <i>add</i>             | Add the inputs, element-wise.   |
| <i>subtract</i>        | Subtract the inputs, element-wise.  |
| <i>multiply</i>        | Multiply the inputs, element-wise.  |
| <i>divide</i>          | Divide the inputs, element-wise.  |
| <i>power</i>           | Elements of x raised to powers of y, element-wise.  |
| <i>maximum</i>         | Maximum of the inputs, element-wise.  |
| <i>minimum</i>         | Minimum of the inputs, element-wise.  |
| <i>max</i>             | Maximum of elements along the specified axis (or axes) of the input.                                |
| <i>min</i>             | Minimum of elements along the specified axis (or axes) of the input.                                |
| <i>mean</i>            | Arithmetic mean along the specified axis (or axes) of the input.                                    |
| <i>var</i>             | Variance along the specified axis (or axes) of the input.   |
| <i>sum</i>             | Sum of elements along the specified axis (or axes) of the input.                                    |
| <i>prod</i>            | Product of elements along the specified axis (or axes) of the input.                                |
| <i>all</i>             | Whether all elements along the specified axis (or axes) evaluate to True.                           |
| <i>abs</i>             | Absolute value of the input, element-wise.  |
| <i>negative</i>        | Numerical negative of the input, element-wise.  |
| <i>sign</i>            | Sign of the numbers of the input, element-wise.   |
| <i>trunc</i>           | Truncated value of the input, element-wise.   |
| <i>floor</i>           | Floor of the input, element-wise.   |
| <i>ceil</i>            | Ceiling of the input, element-wise.   |
| <i>mod</i>             | Integer Mod of $x / y$ , element-wise   |
| <i>fmod</i>            | Floor-mod of $x / y$ , element-wise.  |
| <i>exp</i>             | Exponential of the input, element-wise.   |
| <i>log</i>             | Natural logarithm of the input, element-wise.   |
| <i>cos</i>             | Cosine of the input, element-wise.  |
| <i>sin</i>             | Sine of the input, element-wise.  |
| <i>tan</i>             | Tangent of the input, element-wise.   |
| <i>tanh</i>            | Hyperbolic tangent of the input, element-wise.  |
| <i>arctan</i>          | Inverse tangent of the input, element-wise.   |
| <i>sqrt</i>            | Non-negative square-root of the input, element-wise.  |
| <i>rsqrt</i>           | Reciprocal of the square-root of the input, element-wise.   |
| <i>sigmoid</i>         | Logistic sigmoid activation function on the input, element-wise.                                    |
| <i>relu</i>            | Rectified Linear Unit activation function on the input, element-wise.                               |
| <i>gelu</i>            | Gaussian Error Linear Unit activation function on the input, element-wise.                          |
| <i>gelu_dx</i>         | Derivative of Gaussian Error Linear Unit (gelu) on the input, element-wise.                         |
| <i>gelu_apprx_tanh</i> | Gaussian Error Linear Unit activation function on the input, element-wise, with tanh approximation. |
| <i>silu</i>            | Sigmoid Linear Unit activation function on the input, element-wise.                                 |

continues on next page

Table 7.2 – continued from previous page

|                         |   |
|-------------------------|---|
| <code>silu_dx</code>    | Derivative of Sigmoid Linear Unit activation function on the input, element-wise. |
| <code>erf</code>        | Error function of the input, element-wise.  |
| <code>erf_dx</code>     | Derivative of the Error function (erf) on the input, element-wise.                |
| <code>softplus</code>   | Softplus activation function on the input, element-wise.                          |
| <code>mish</code>       | Mish activation function on the input, element-wise.                              |
| <code>square</code>     | Square of the input, element-wise.  |
| <code>softmax</code>    | Softmax activation function on the input, element-wise.                           |
| <code>rms_norm</code>   | Apply Root Mean Square Layer Normalization.                                       |
| <code>dropout</code>    | Randomly zeroes some of the elements of the input tile given a probability rate.  |
| <code>matmul</code>     | $x @ y$ matrix multiplication of $x$ and $y$ .                                    |
| <code>transpose</code>  | Transposes a 2D tile between its partition and free dimension.                    |
| <code>reciprocal</code> | Reciprocal of the the input, element-wise.  |

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.add

`nki.language.add(x, y, *, dtype=None, mask=None, **kwargs)`

Add the inputs, element-wise.

((Similar to [numpy.add](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has  $x + y$ , element-wise.

Examples:

```
import neuronxcc.nki.language as nl

a = nl.load(a_tensor[0:128, 0:512])
b = nl.load(b_tensor[0:128, 0:512])
# add a and b element-wise and store in c[128, 512]
c = nl.add(a, b)
nl.store(c_tensor[0:128, 0:512], c)

a = nl.load(a_tensor[0:128, 0:512])
b = 2.2
# add constant b to each element in a
c = nl.add(a, b)
```

(continues on next page)

(continued from previous page)

```

nl.store(c_tensor[0:128, 0:512], c)

a = nl.load(a_tensor[0:128, 0:512])
b = nl.load(b_tensor[0:128, 0:1])
# broadcast on free dimension -- [128, 1] is broadcasted to [128, 512]
c = nl.add(a, b)
nl.store(c_tensor[0:128, 0:512], c)

a = nl.load(a_tensor[0:128, 0:512])
b = nl.load(b_tensor[0:1, 0:512])
# broadcast on partition dimension -- [1, 512] is broadcasted to [128, 512]
c = nl.add(a, b)
nl.store(c_tensor[0:128, 0:512], c)

a = nl.load(a_tensor[0:128, 0:512])
b = nl.load(b_tensor[0:1, 0:1])
# broadcast on both dimensions -- [1, 1] is broadcasted to [128, 512]
c = nl.add(a, b)
nl.store(c_tensor[0:128, 0:512], c)

a = nl.load(a_tensor[0:128, 0:1])
b = nl.load(b_tensor[0:1, 0:512])
# broadcast on each dimensions -- [128, 1] and [1, 512] are broadcasted to [128, 512]
c = nl.add(a, b)
nl.store(c_tensor[0:128, 0:512], c)

```

---

**Note:** Broadcasting in the partition dimension is generally more expensive than broadcasting in free dimension. It is recommended to align your data to perform free dimension broadcast whenever possible.

---

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.subtract

`nki.language.subtract(x, y, *, dtype=None, mask=None, **kwargs)`

Subtract the inputs, element-wise.

((Similar to [numpy.subtract](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. **x.shape** and **y.shape** must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has  $x - y$ , element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.multiply**

`nki.language.multiply(x, y, *, dtype=None, mask=None, **kwargs)`

Multiply the inputs, element-wise.

((Similar to [numpy.multiply](#)))

**Parameters**

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has  $x * y$ , element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.divide**

`nki.language.divide(x, y, *, dtype=None, mask=None, **kwargs)`

Divide the inputs, element-wise.

((Similar to [numpy.divide](#)))

**Parameters**

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has  $x / y$ , element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.power

`nki.language.power(x, y, *, dtype=None, mask=None, **kwargs)`

Elements of `x` raised to powers of `y`, element-wise.

((Similar to [numpy.power](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values `x` to the power of `y`.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.maximum

`nki.language.maximum(x, y, *, dtype=None, mask=None, **kwargs)`

Maximum of the inputs, element-wise.

((Similar to [numpy.maximum](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has the maximum of each elements from `x` and `y`.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



## nki.language.minimum

`nki.language.minimum(x, y, *, dtype=None, mask=None, **kwargs)`

Minimum of the inputs, element-wise.

((Similar to [numpy.minimum](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has the minimum of each elements from x and y.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.max

`nki.language.max(x, axis, *, dtype=None, mask=None, keepdims=False, **kwargs)`

Maximum of elements along the specified axis (or axes) of the input.

((Similar to [numpy.max](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: `[1]`, `[1,2]`, `[1,2,3]`, `[1,2,3,4]`
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **keepdims** – If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

### Returns

a tile with the maximum of elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.min

`nki.language.min(x, axis, *, dtype=None, mask=None, keepdims=False, **kwargs)`

Minimum of elements along the specified axis (or axes) of the input.

((Similar to [numpy.min](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: `[1]`, `[1,2]`, `[1,2,3]`, `[1,2,3,4]`
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **keepdims** – If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

### Returns

a tile with the minimum of elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.mean

`nki.language.mean(x, axis, *, dtype=None, mask=None, keepdims=False, **kwargs)`

Arithmetic mean along the specified axis (or axes) of the input.

((Similar to [numpy.mean](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: `[1]`, `[1,2]`, `[1,2,3]`, `[1,2,3,4]`
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with the average of elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed. float32 intermediate and return values are used for integer inputs.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.var

`nki.language.var(x, axis, *, dtype=None, mask=None, **kwargs)`

Variance along the specified axis (or axes) of the input.

((Similar to [numpy.var](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: [1], [1,2], [1,2,3], [1,2,3,4]
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with the variance of the elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.sum

`nki.language.sum(x, axis, *, dtype=None, mask=None, keepdims=False, **kwargs)`

Sum of elements along the specified axis (or axes) of the input.

((Similar to [numpy.sum](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: [1], [1,2], [1,2,3], [1,2,3,4]
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **keepdims** – If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

### Returns

a tile with the sum of elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.prod

`nki.language.prod(x, axis, *, dtype=None, mask=None, keepdims=False, **kwargs)`

Product of elements along the specified axis (or axes) of the input.

((Similar to [numpy.prod](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: `[1]`, `[1,2]`, `[1,2,3]`, `[1,2,3,4]`
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **keepdims** – If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

### Returns

a tile with the product of elements along the provided axis. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.all

`nki.language.all(x, axis, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Whether all elements along the specified axis (or axes) evaluate to True.

((Similar to [numpy.all](#)))

### Parameters

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: `[1]`, `[1,2]`, `[1,2,3]`, `[1,2,3,4]`
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a boolean tile with the result. This return tile will have a shape of the input tile's shape with the specified axes removed.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.abs

`nki.language.abs(x, *, dtype=None, mask=None, **kwargs)`

Absolute value of the input, element-wise.

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has absolute values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.negative

`nki.language.negative(x, *, dtype=None, mask=None, **kwargs)`

Numerical negative of the input, element-wise.

((Similar to [numpy.negative](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has numerical negative values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.sign

`nki.language.sign(x, *, dtype=None, mask=None, **kwargs)`

Sign of the numbers of the input, element-wise.

((Similar to [numpy.sign](#)))

The sign function returns -1 if  $x < 0$ , 0 if  $x == 0$ , 1 if  $x > 0$ .

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has sign values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### **nki.language.trunc**

`nki.language.trunc(x, *, dtype=None, mask=None, **kwargs)`

Truncated value of the input, element-wise.

((Similar to [numpy.trunc](#)))

The truncated value of the scalar  $x$  is the nearest integer  $i$  which is closer to zero than  $x$  is. In short, the fractional part of the signed number  $x$  is discarded.

#### **Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### **Returns**

a tile that has truncated values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### **nki.language.floor**

`nki.language.floor(x, *, dtype=None, mask=None, **kwargs)`

Floor of the input, element-wise.

((Similar to [numpy.floor](#)))

The floor of the scalar  $x$  is the largest integer  $i$ , such that  $i \leq x$ .

#### **Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### **Returns**

a tile that has floor values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.ceil

`nki.language.ceil(x, *, dtype=None, mask=None, **kwargs)`

Ceiling of the input, element-wise.

((Similar to [numpy.ceil](#)))

The ceil of the scalar  $x$  is the smallest integer  $i$ , such that  $i \geq x$ .

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has ceiling values of  $x$ .

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.mod

`nki.language.mod(x, y, dtype=None, mask=None, **kwargs)`

Integer Mod of  $x / y$ , element-wise

Computes the remainder complementary to the `floor_divide` function. It is equivalent to the Python modulus  $x \% y$  and has the same sign as the divisor  $y$ .

((Similar to [numpy.mod](#)))

### Parameters

- **x** – a tile. If  $x$  is a scalar value it will be broadcast to the shape of  $y$ .  $x.shape$  and  $y.shape$  must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **y** – a tile or a scalar value.  $x.shape$  and  $y.shape$  must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values  $x \bmod y$ .

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.fmod

`nki.language.fmod(x, y, dtype=None, mask=None, **kwargs)`

Floor-mod of  $x / y$ , element-wise.

The remainder has the same sign as the dividend  $x$ . It is equivalent to the Matlab(TM) `rem` function and should not be confused with the Python modulus operator  $x \% y$ .

((Similar to [numpy.fmod](#)))

### Parameters

- **x** – a tile. If  $x$  is a scalar value it will be broadcast to the shape of  $y$ .  $x.shape$  and  $y.shape$  must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **y** – a tile or a scalar value.  $x.shape$  and  $y.shape$  must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values  $x \text{ fmod } y$ .

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.exp

`nki.language.exp(x, *, dtype=None, mask=None, **kwargs)`

Exponential of the input, element-wise.

((Similar to [numpy.exp](#)))

The `exp(x)` is  $e^x$  where  $e$  is the Euler's number = 2.718281...

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has exponential values of  $x$ .

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



## nki.language.log

`nki.language.log(x, *, dtype=None, mask=None, **kwargs)`

Natural logarithm of the input, element-wise.

((Similar to [numpy.log](#)))

It is the inverse of the exponential function, such that:  $\log(\exp(x)) = x$ . The natural logarithm base is  $e$ .

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has natural logarithm values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.cos

`nki.language.cos(x, *, dtype=None, mask=None, **kwargs)`

Cosine of the input, element-wise.

((Similar to [numpy.cos](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has cosine values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.sin

`nki.language.sin(x, *, dtype=None, mask=None, **kwargs)`

Sine of the input, element-wise.

((Similar to [numpy.sin](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has sine values of **x**.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.tan**

`nki.language.tan(x, *, dtype=None, mask=None, **kwargs)`

Tangent of the input, element-wise.

((Similar to [numpy.tan](#)))

**Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has tangent values of **x**.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.tanh**

`nki.language.tanh(x, *, dtype=None, mask=None, **kwargs)`

Hyperbolic tangent of the input, element-wise.

((Similar to [numpy.tanh](#)))

**Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has hyperbolic tangent values of **x**.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.arctan

`nki.language.arctan(x, *, dtype=None, mask=None, **kwargs)`

Inverse tangent of the input, element-wise.

((Similar to [numpy.arctan](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has inverse tangent values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.sqrt

`nki.language.sqrt(x, *, dtype=None, mask=None, **kwargs)`

Non-negative square-root of the input, element-wise.

((Similar to [numpy.sqrt](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has square-root values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.rsqrt

`nki.language.rsqrt(x, *, dtype=None, mask=None, **kwargs)`

Reciprocal of the square-root of the input, element-wise.

((Similar to [torch.rsqrt](#)))

$\text{rsqrt}(x) = 1 / \text{sqrt}(x)$

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has reciprocal square-root values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.language.sigmoid**

`nki.language.sigmoid(x, *, dtype=None, mask=None, **kwargs)`

Logistic sigmoid activation function on the input, element-wise.

((Similar to `torch.nn.functional.sigmoid`))

$\text{sigmoid}(x) = 1/(1+\exp(-x))$

**Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has sigmoid of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.language.relu**

`nki.language.relu(x, *, dtype=None, mask=None, **kwargs)`

Rectified Linear Unit activation function on the input, element-wise.

$\text{relu}(x) = (x) + \max(0, x)$

((Similar to `torch.nn.functional.relu`))

**Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has relu of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.gelu

`nki.language.gelu(x, *, dtype=None, mask=None, **kwargs)`

Gaussian Error Linear Unit activation function on the input, element-wise.

((Similar to [torch.nn.functional.gelu](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has gelu of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.gelu\_dx

`nki.language.gelu_dx(x, *, dtype=None, mask=None, **kwargs)`

Derivative of Gaussian Error Linear Unit (gelu) on the input, element-wise.

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has gelu\_dx of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.gelu\_aprx\_tanh

`nki.language.gelu_aprx_tanh(x, *, dtype=None, mask=None, **kwargs)`

Gaussian Error Linear Unit activation function on the input, element-wise, with tanh approximation.

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has gelu of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.silu

`nki.language.silu(x, *, dtype=None, mask=None, **kwargs)`

Sigmoid Linear Unit activation function on the input, element-wise.

((Similar to [torch.nn.functional.silu](#)))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has silu of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.silu\_dx

`nki.language.silu_dx(x, *, dtype=None, mask=None, **kwargs)`

Derivative of Sigmoid Linear Unit activation function on the input, element-wise.

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has silu\_dx of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.erf

`nki.language.erf(x, *, dtype=None, mask=None, **kwargs)`

Error function of the input, element-wise.

((Similar to [torch.erf](#)))

$\text{erf}(x) = 2/\sqrt{\pi} \cdot \text{integral}(\exp(-t^2), t=0..x)$ .

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has erf of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.erf\_dx

`nki.language.erf_dx(x, *, dtype=None, mask=None, **kwargs)`

Derivative of the Error function (erf) on the input, element-wise.

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has erf\_dx of x.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.softplus

`nki.language.softplus(x, *, dtype=None, mask=None, **kwargs)`

Softplus activation function on the input, element-wise.

Softplus is a smooth approximation to the ReLU activation, defined as:

$$\text{softplus}(x) = \log(1 + \exp(x))$$

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has softplus of x.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.mish

`nki.language.mish(x, *, dtype=None, mask=None, **kwargs)`

Mish activation function on the input, element-wise.

Mish: A Self Regularized Non-Monotonic Neural Activation Function is defined as:

$$\text{mish}(x) = x * \tanh(\text{softplus}(x))$$

see: <https://arxiv.org/abs/1908.08681>

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has mish of **x**.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.square**

`nki.language.square(x, *, dtype=None, mask=None, **kwargs)`

Square of the input, element-wise.

((Similar to [numpy.square](#)))

**Parameters**

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has square of **x**.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.softmax**

`nki.language.softmax(x, axis, *, dtype=None, compute_dtype=None, mask=None, **kwargs)`

Softmax activation function on the input, element-wise.

((Similar to [torch.nn.functional.softmax](#)))

**Parameters**

- **x** – a tile.
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: [1], [1,2], [1,2,3], [1,2,3,4]
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **compute\_dtype** – (optional) dtype for the internal computation - currently ``dtype`` and ``compute_dtype`` behave the same, both sets internal compute and return dtype.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

a tile that has softmax of **x**.



*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.rms\_norm

```
nki.language.rms_norm(x, w, axis, n, epsilon=1e-06, *, dtype=None, compute_dtype=None, mask=None,
                      **kwargs)
```

Apply Root Mean Square Layer Normalization.

### Parameters

- **x** – input tile
- **w** – weight tile
- **axis** – axis along which to compute the root mean square (rms) value
- **n** – total number of values to calculate rms
- **epsilon** – epsilon value used by rms calculation to avoid divide-by-zero
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **compute\_dtype** – (optional) dtype for the internal computation - currently `dtype` and `compute\_dtype` behave the same, both sets internal compute and return dtype.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

`x / RMS(x) * w`

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.dropout

```
nki.language.dropout(x, rate, *, dtype=None, mask=None, **kwargs)
```

Randomly zeroes some of the elements of the input tile given a probability rate.

### Parameters

- **x** – a tile.
- **rate** – a scalar value or a tile with 1 element, with the probability rate.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with randomly zeroed elements of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.matmul

`nki.language.matmul(x, y, *, transpose_x=False, mask=None, **kwargs)`

$x @ y$  matrix multiplication of  $x$  and  $y$ .

((Similar to `numpy.matmul`))

---

**Note:** For optimal performance on hardware, use `nki.isa.nc_matmul()` or call `nki.language.matmul` with `transpose_x=True`. Use `nki.isa.nc_matmul` also to access low-level features of the Tensor Engine.

---



---

**Note:** Implementation details: `nki.language.matmul` calls `nki.isa.nc_matmul` under the hood. `nc_matmul` is neuron specific customized implementation of `matmul` that computes  $x.T @ y$ , as a result, `matmul(x, y)` lowers to `nc_matmul(transpose(x), y)`. To avoid this extra transpose instruction being inserted, use  $x.T$  and `transpose_x=True` inputs to this `matmul`.

---

### Parameters

- **x** – a tile on SBUF (partition dimension  $\leq 128$ , free dimension  $\leq 128$ ),  $x$ 's free dimension must match  $y$ 's partition dimension.
- **y** – a tile on SBUF (partition dimension  $\leq 128$ , free dimension  $\leq 512$ )
- **transpose\_x** – Defaults to `False`. If `True`,  $x$  is treated as already transposed. If `False`, an additional transpose will be inserted to make  $x$ 's partition dimension the contract dimension of the `matmul` to align with the Tensor Engine.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

$x @ y$  or  $x.T @ y$  if `transpose_x=True`

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.transpose

`nki.language.transpose(x, *, dtype=None, mask=None, **kwargs)`

Transposes a 2D tile between its partition and free dimension.

### Parameters

- **x** – 2D input tile
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has the values of the input tile with its partition and free dimensions swapped.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.reciprocal

`nki.language.reciprocal(x, *, dtype=None, mask=None, **kwargs)`

Reciprocal of the the input, element-wise.

((Similar to [numpy.reciprocal](#)))

`reciprocal(x) = 1 / x`

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has reciprocal values of **x**.

*This document is relevant for: Inf2, Trn1, Trn2*

## Bitwise operations

|  |  |
|--|--|
| <a href="#"><code>bitwise_and</code></a> | Bitwise AND of the two inputs, element-wise. |
| <a href="#"><code>bitwise_or</code></a>  | Bitwise OR of the two inputs, element-wise.  |
| <a href="#"><code>bitwise_xor</code></a> | Bitwise XOR of the two inputs, element-wise. |
| <a href="#"><code>invert</code></a>      | Bitwise NOT of the input, element-wise.      |
| <a href="#"><code>left_shift</code></a>  | Bitwise left-shift x by y, element-wise.     |
| <a href="#"><code>right_shift</code></a> | Bitwise right-shift x by y, element-wise.    |

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.bitwise\_and

`nki.language.bitwise_and(x, y, *, dtype=None, mask=None, **kwargs)`

Bitwise AND of the two inputs, element-wise.

((Similar to [numpy.bitwise\\_and](#)))

Computes the bit-wise AND of the underlying binary representation of the integers in the input tiles. This function implements the C/Python operator `&`

### Parameters

- **x** – a tile or a scalar value of integer type.
- **y** – a tile or a scalar value of integer type. **x.shape** and **y.shape** must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values **x** & **y**.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## **nki.language.bitwise\_or**

`nki.language.bitwise_or(x, y, *, dtype=None, mask=None, **kwargs)`

Bitwise OR of the two inputs, element-wise.

((Similar to [numpy.bitwise\\_or](#)))

Computes the bit-wise OR of the underlying binary representation of the integers in the input tiles. This function implements the C/Python operator `|`

### **Parameters**

- **x** – a tile or a scalar value of integer type.
- **y** – a tile or a scalar value of integer type. **x.shape** and **y.shape** must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### **Returns**

a tile that has values  $x \mid y$ .

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## **nki.language.bitwise\_xor**

`nki.language.bitwise_xor(x, y, *, dtype=None, mask=None, **kwargs)`

Bitwise XOR of the two inputs, element-wise.

((Similar to [numpy.bitwise\\_xor](#)))

Computes the bit-wise XOR of the underlying binary representation of the integers in the input tiles. This function implements the C/Python operator `^`

### **Parameters**

- **x** – a tile or a scalar value of integer type.
- **y** – a tile or a scalar value of integer type. **x.shape** and **y.shape** must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### **Returns**

a tile that has values  $x \wedge y$ .

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.invert

`nki.language.invert(x, *, dtype=None, mask=None, **kwargs)`

Bitwise NOT of the input, element-wise.

((Similar to [numpy.invert](#)))

Computes the bit-wise NOT of the underlying binary representation of the integers in the input tile. This ufunc implements the C/Python operator `~`

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with bitwise NOT `x` element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.left\_shift

`nki.language.left_shift(x, y, *, dtype=None, mask=None, **kwargs)`

Bitwise left-shift `x` by `y`, element-wise.

((Similar to [numpy.left\\_shift](#)))

Computes the bit-wise left shift of the underlying binary representation of the integers in the input tiles. This function implements the C/Python operator `<<`

### Parameters

- **x** – a tile or a scalar value of integer type.
- **y** – a tile or a scalar value of integer type. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values `x << y`.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.right\_shift

`nki.language.right_shift(x, y, *, dtype=None, mask=None, **kwargs)`

Bitwise right-shift `x` by `y`, element-wise.

((Similar to [numpy.right\\_shift](#)))

Computes the bit-wise right shift of the underlying binary representation of the integers in the input tiles. This function implements the C/Python operator `>>`

### Parameters

- **x** – a tile or a scalar value of integer type.
- **y** – a tile or a scalar value of integer type. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile that has values `x >> y`.

*This document is relevant for:* Inf2, Trn1, Trn2

## Logical operations

|                               |   |
|-------------------------------|---|
| <a href="#">equal</a>         | Element-wise boolean result of <code>x == y</code> .    |
| <a href="#">not_equal</a>     | Element-wise boolean result of <code>x != y</code> .    |
| <a href="#">greater</a>       | Element-wise boolean result of <code>x &gt; y</code> .  |
| <a href="#">greater_equal</a> | Element-wise boolean result of <code>x &gt;= y</code> . |
| <a href="#">less</a>          | Element-wise boolean result of <code>x &lt; y</code> .  |
| <a href="#">less_equal</a>    | Element-wise boolean result of <code>x &lt;= y</code> . |
| <a href="#">logical_and</a>   | Element-wise boolean result of <code>x AND y</code> .   |
| <a href="#">logical_or</a>    | Element-wise boolean result of <code>x OR y</code> .    |
| <a href="#">logical_xor</a>   | Element-wise boolean result of <code>x XOR y</code> .   |
| <a href="#">logical_not</a>   | Element-wise boolean result of <code>NOT x</code> .     |

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.equal

`nki.language.equal(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of `x == y`.

((Similar to [numpy.equal](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the

input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);

- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

a tile with boolean result of  $x == y$  element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### nki.language.not\_equal

`nki.language.not_equal(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of  $x != y$ .

((Similar to [numpy.not\\_equal](#)))

#### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

a tile with boolean result of  $x != y$  element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### nki.language.greater

`nki.language.greater(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of  $x > y$ .

((Similar to [numpy.greater](#)))

#### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

a tile with boolean result of  $x > y$  element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.greater\_equal

`nki.language.greater_equal(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of  $x \geq y$ .

((Similar to [numpy.greater\\_equal](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x \geq y$  element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.less

`nki.language.less(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of  $x < y$ .

((Similar to [numpy.less](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x < y$  element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



## nki.language.less\_equal

```
nki.language.less_equal(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)
```

Element-wise boolean result of  $x \leq y$ .

((Similar to [numpy.less\\_equal](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x \leq y$  element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.logical\_and

```
nki.language.logical_and(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)
```

Element-wise boolean result of  $x \text{ AND } y$ .

((Similar to [numpy.logical\\_and](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x \text{ AND } y$  element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.logical\_or

```
nki.language.logical_or(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)
```

Element-wise boolean result of  $x$  OR  $y$ .

((Similar to [numpy.logical\\_or](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x$  OR  $y$  element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.logical\_xor

```
nki.language.logical_xor(x, y, *, dtype=<class 'bool'>, mask=None, **kwargs)
```

Element-wise boolean result of  $x$  XOR  $y$ .

((Similar to [numpy.logical\\_xor](#)))

### Parameters

- **x** – a tile or a scalar value.
- **y** – a tile or a scalar value. `x.shape` and `y.shape` must be [broadcastable](#) to a common shape, that will become the shape of the output.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of  $x$  XOR  $y$  element-wise.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.logical\_not

`nki.language.logical_not(x, *, dtype=<class 'bool'>, mask=None, **kwargs)`

Element-wise boolean result of NOT `x`.

((Similar to `numpy.logical_not`))

### Parameters

- **x** – a tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with boolean result of NOT `x` element-wise.

*This document is relevant for: Inf2, Trn1, Trn2*

## Tensor manipulation operations

|                          |  |
|--------------------------|--|
| <code>ds</code>          | Construct a dynamic slice for simple tensor indexing.  |
| <code>arange</code>      | Return contiguous values within a given interval, used for indexing a tensor to define a tile.   |
| <code>mgrid</code>       | Same as NumPy <code>mgrid</code> : "An instance which returns a dense (or fleshed out) mesh-grid when indexed, so that each returned argument has the same shape." |
| <code>expand_dims</code> | Expand the shape of a tile.  |

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.ds

`nki.language.ds(start, size)`

Construct a dynamic slice for simple tensor indexing.

```
import neuronxcc.nki.language as nl
...

@nki.jit(mode="simulation")
def example_kernel(in_tensor):
    out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
                             buffer=nl.shared_hbm)
    for i in nl.affine_range(in_tensor.shape[1] // 512):
        tile = nl.load(in_tensor[:, (i * 512):((i + 1) * 512)])
        # Same as above but use ds (dynamic slice) instead of the native
        # slice syntax
        tile = nl.load(in_tensor[:, nl.ds(i * 512, 512)])
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.arange

`nki.language.arange(*args)`

Return contiguous values within a given interval, used for indexing a tensor to define a tile.

((Similar to [numpy.arange](#)))

**arange can be called as:**

- `arange(stop)`: Values are generated within the half-open interval `[0, stop)` (the interval including zero, excluding stop).
- `arange(start, stop)`: Values are generated within the half-open interval `[start, stop)` (the interval including start, excluding stop).

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.mgrid

`nki.language.mgrid = Ellipsis`

Same as NumPy `mgrid`: “An instance which returns a dense (or fleshed out) mesh-grid when indexed, so that each returned argument has the same shape. The dimensions and number of the output arrays are equal to the number of indexing dimensions.”

Complex numbers are not supported in the step length.

((Similar to [numpy.mgrid](#)))

```
import neuronxcc.nki.language as nl
...

i_p, i_f = nl.mgrid[0:128, 0:512]
tile = nl.load(in_tensor[i_p, i_f])
...
nl.store(out_tensor[i_p, i_f], tile)
```

```
import neuronxcc.nki.language as nl
...

grid = nl.mgrid[0:128, 0:512]
tile = nl.load(in_tensor[grid.p, grid.x])
...
nl.store(out_tensor[grid.p, grid.x], tile)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.expand\_dims

`nki.language.expand_dims(data, axis)`

Expand the shape of a tile. Insert a new axis that will appear at the `axis` position in the expanded tile shape. Currently only supports expanding dimensions after the last index of the tile.

((Similar to [numpy.expand\\_dims](#)))

### Parameters

- **data** – a tile input
- **axis** – int or tuple/list of ints. Position in the expanded axes where the new axis (or axes) is placed; must be free dimensions, not partition dimension (0); Currently only supports axis (or axes) after the last index.

### Returns

a tile with view of input `data` with the number of dimensions increased.

*This document is relevant for: Inf2, Trn1, Trn2*

## Indexing/Searching operations

|   |  |
|---|--|
| <a href="#"><code>where</code></a>            | Return elements chosen from <code>x</code> or <code>y</code> depending on condition. |
| <a href="#"><code>gather_flattened</code></a> | Gather elements from <code>data</code> according to the <code>indices</code> .       |

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.where

`nki.language.where(condition, x, y, *, dtype=None, mask=None, **kwargs)`

Return elements chosen from `x` or `y` depending on condition.

((Similar to [numpy.where](#)))

### Parameters

- **condition** – if True, yield `x`, otherwise yield `y`.
- **x** – a tile with values from which to choose if condition is True.
- **y** – a tile or a numerical value from which to choose if condition is False.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

a tile with elements from `x` where condition is True, and elements from `y` otherwise.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.language.gather\_flattened**

`nki.language.gather_flattened(data, indices, *, mask=None, dtype=None, **kwargs)`

Gather elements from `data` according to the `indices`.

This instruction gathers elements from the `data` tensor using integer indices provided in the `indices` tensor. For each element in the `indices` tensor, it retrieves the corresponding value from the `data` tensor using the index value to select from the free dimension of `data`. The gather instruction effectively performs up to 128 parallel gather operations, with each operation using the corresponding partition of `data` and `indices`.

The output tensor has the same shape as the `indices` tensor, with each output element containing the value from `data` at the position specified by the corresponding index. Out of bounds indices will return garbage values.

Both `data` and `indices` must be 2-, 3-, or 4-dimensional. The `indices` tensor must contain uint32 values.

For indexing purposes, all free dimensions are flattened and indexed as the same “row”. Consider this example:

```
data =
[[[1., 2.],
  [3., 4.]],
 [[5., 6.],
  [7., 8.]]]
indices =
[[[0, 1],
  [1, 3]],
 [[3, 3],
  [1, 0]]]
nl.gather_flattened(data, indices) produces this result:
[[[1., 2.],
  [2., 4.]],
 [[8., 8.],
  [6., 5.]]]
```

With the exception of handling out-of-bounds indices, this behavior is equivalent to:

```
indices_flattened = indices.reshape(indices.shape[0], -1)
data_flattened = data.reshape(data.shape[0], -1)
result = np.take_along_axis(data_flattened, indices_flattened, axis=-1)
result.reshape(indices.shape)
```

((Similar to [torch.gather\\_flattened](#)))

**Parameters**

- **data** – the source tensor to gather values from
- **indices** – tensor containing uint32 indices to gather across the flattened free dimension.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

**Returns**

a tensor with the same shape as `indices` containing gathered values from `data`

Example:

```

import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Gather values from a tensor using indices
#####
# Create source tensor
N = 32
M = 64
data = nl.rand((N, M), dtype=nl.float32)

# Create indices tensor - gather every 5th element
indices = nl.zeros((N, 10), dtype=nl.uint32)
for i in nl.static_range(N):
    for j in nl.static_range(10):
        indices[i, j] = j * 5

# Gather values from data according to indices
result = nl.gather_flattened(data=data, indices=indices)

```

*This document is relevant for: Inf2, Trn1, Trn2*

## Collective communication operations

---

*all\_reduce*

Apply reduce operation over multiple SPMD programs.

---

*This document is relevant for: Inf2, Trn1, Trn2*

### nki.language.all\_reduce

`nki.language.all_reduce(x, op, program_axes, *, dtype=None, mask=None, parallel_reduce=True, asynchronous=False, **kwargs)`

Apply reduce operation over multiple SPMD programs.

#### Parameters

- **x** – a tile.
- **op** – numpy ALU operator to use to reduce over the input tile.
- **program\_axes** – a single axis or a tuple of axes along which the reduction operation is performed.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **parallel\_reduce** – optional boolean parameter whether to turn on parallel reduction. Enable parallel reduction consumes additional memory.
- **asynchronous** – Defaults to False. If True, caller should synchronize before reading final result, e.g. using `nki.sync_thread`.

#### Returns

the reduced resulting tile

*This document is relevant for: Inf2, Trn1, Trn2*

## Iterators

|                               |  |
|-------------------------------|--|
| <code>static_range</code>     | Create a sequence of numbers for use as loop iterators in NKI, resulting in a fully unrolled loop. |
| <code>affine_range</code>     | Create a sequence of numbers for use as <b>parallel</b> loop iterators in NKI.                     |
| <code>sequential_range</code> | Create a sequence of numbers for use as <b>sequential</b> loop iterators in NKI.                   |

This document is relevant for: Inf2, Trn1, Trn2

### nki.language.static\_range

`nki.language.static_range(*args)`

Create a sequence of numbers for use as loop iterators in NKI, resulting in a fully unrolled loop. Unlike `affine_range` or `sequential_range`, Neuron compiler will fully unroll the loop during NKI kernel tracing.

Notes:

- Due to loop unrolling, compilation time may go up significantly compared to `affine_range` or `sequential_range`.
- On-chip memory (SBUF) usage may also go up significantly compared to `affine_range` or `sequential_range`.
- No loop-level optimizations will be performed in the compiler.
- `static_range` should only be used as a fall-back option for debugging purposes when `affine_range` or `sequential_range` is giving functionally incorrect results or undesirable performance characteristics.

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

### nki.language.affine\_range

`nki.language.affine_range(*args, **kwargs)`

Create a sequence of numbers for use as **parallel** loop iterators in NKI. `affine_range` should be the default loop iterator choice, when there is **no** loop carried dependency. Note, associative reductions are **not** considered loop carried dependencies in this context. A concrete example of associative reduction is multiple `nl.matmul` or `nisa.nc_matmul` calls accumulating into the same output buffer defined outside of this loop level (see code example #2 below).

When the above conditions are not met, we recommend using `sequential_range` instead.

Notes:

- Using `affine_range` prevents Neuron compiler from unrolling the loops until entering compiler backend, which typically results in better compilation time compared to the fully unrolled iterator `static_range`.
- Using `affine_range` also allows Neuron compiler to perform additional loop-level optimizations, such as loop vectorization in current release. The exact type of loop-level optimizations applied is subject to changes in future releases.
- Since each kernel instance only runs on a single NeuronCore, `affine_range` does **not** parallelize different loop iterations across multiple NeuronCores. However, different iterations could be parallelized/pipelined on different compute engines within a NeuronCore depending on the invoked instructions (engines) and data dependency in the loop body.



```

1 import neuronxcc.nki.language as nl
2
3 #####
4 # Example 1: No loop carried dependency
5 # Input/Output tensor shape: [128, 2048]
6 # Load one tile ([128, 512]) at a time, square the tensor element-wise,
7 # and store it into output tile
8 #####
9
10 # Every loop instance works on an independent input/output tile.
11 # No data dependency between loop instances.
12 for i_input in nl.affine_range(input.shape[1] // 512):
13     offset = i_input * 512
14     input_sb = nl.load(input[0:input.shape[0], offset:offset+512])
15     result = nl.multiply(input_sb, input_sb)
16     nl.store(output[0:input.shape[0], offset:offset+512], result)
17
18 #####
19 # Example 2: Matmul output buffer accumulation, a type of associative reduction
20 # Input tensor shapes for nl.matmul: xT[K=2048, M=128] and y[K=2048, N=128]
21 # Load one tile ([128, 128]) from both xT and y at a time, matmul and
22 # accumulate into the same output buffer
23 #####
24
25 result_psum = nl.zeros((128, 128), dtype=nl.float32, buffer=nl.psum)
26 for i_K in nl.affine_range(xT.shape[0] // 128):
27     offset = i_K * 128
28     xT_sbuf = nl.load(offset:offset+128, 0:xT.shape[1])
29     y_sbuf = nl.load(offset:offset+128, 0:y.shape[1])
30
31     result_psum += nl.matmul(xT_sbuf, y_sbuf, transpose_x=True)

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.sequential\_range

`nki.language.sequential_range(*args, **kwargs)`

Create a sequence of numbers for use as **sequential** loop iterators in NKI. `sequential_range` should be used when there is a loop carried dependency. Note, associative reductions are **not** considered loop carried dependencies in this context. See [affine\\_range](#) for an example of such associative reduction.

Notes:

- Inside a NKI kernel, any use of Python `range(...)` will be replaced with `sequential_range(...)` by Neuron compiler.
- Using `sequential_range` prevents Neuron compiler from unrolling the loops until entering compiler backend, which typically results in better compilation time compared to the fully unrolled iterator [static\\_range](#).
- Using `sequential_range` informs Neuron compiler to respect inter-loop dependency and perform much more conservative loop-level optimizations compared to `affine_range`.
- Using `affine_range` instead of `sequential_range` in case of loop carried dependency incorrectly is considered unsafe and could lead to numerical errors.

```

1  import neuronxcc.nki.language as nl
2
3  #####
4  # Example 1: Loop carried dependency from tiling tensor_tensor_scan
5  # Both sbuf tensor input0 and input1 shapes: [128, 2048]
6  # Perform a scan operation between the two inputs using a tile size of [128, 512]
7  # Store the scan output to another [128, 2048] tensor
8  #####
9
10 # Loop iterations communicate through this init tensor
11 init = nl.zeros((128, 1), dtype=input0.dtype)
12
13 # This loop will only produce correct results if the iterations are performed in_
  ↪ order
14 for i_input in nl.sequential_range(input0.shape[1] // 512):
15     offset = i_input * 512
16
17     # Depends on scan result from the previous loop iteration
18     result = nisa.tensor_tensor_scan(input0[:, offset:offset+512],
19                                     input1[:, offset:offset+512],
20                                     initial=init,
21                                     op0=nl.multiply, op1=nl.add)
22
23     nl.store(output[0:input0.shape[0], offset:offset+512], result)
24
25     # Prepare initial result for scan in the next loop iteration
26     init[:, :] = result[:, 511]

```

This document is relevant for: Inf2, Trn1, Trn2

## Memory Hierarchy

|                    |   |
|--------------------|---|
| <i>par_dim</i>     | Mark a dimension explicitly as a partition dimension.   |
| <i>psum</i>        | PSUM - Only visible to each individual kernel instance in the SPMD grid, alias of <code>nki.compiler.psum.auto_alloc()</code>         |
| <i>sbuf</i>        | State Buffer - Only visible to each individual kernel instance in the SPMD grid, alias of <code>nki.compiler.sbuf.auto_alloc()</code> |
| <i>hbm</i>         | HBM - Alias of <i>private_hbm</i>   |
| <i>private_hbm</i> | HBM - Only visible to each individual kernel instance in the SPMD grid  |
| <i>shared_hbm</i>  | Shared HBM - Visible to all kernel instances in the SPMD grid   |

This document is relevant for: Inf2, Trn1, Trn2

**nki.language.par\_dim****nki.language.par\_dim = Ellipsis**

Mark a dimension explicitly as a partition dimension.

*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf2, Trn1, Trn2***nki.language.psum****nki.language.psum = Ellipsis**PSUM - Only visible to each individual kernel instance in the SPMD grid, alias of `nki.compiler.psum.auto_alloc()`*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf2, Trn1, Trn2***nki.language.sbuf****nki.language.sbuf = Ellipsis**State Buffer - Only visible to each individual kernel instance in the SPMD grid, alias of `nki.compiler.sbuf.auto_alloc()`*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf2, Trn1, Trn2***nki.language.hbm****nki.language.hbm = Ellipsis**HBM - Alias of `private_hbm`*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf2, Trn1, Trn2***nki.language.private\_hbm****nki.language.private\_hbm = Ellipsis**

HBM - Only visible to each individual kernel instance in the SPMD grid

*This document is relevant for: Inf2, Trn1, Trn2**This document is relevant for: Inf2, Trn1, Trn2*

### nki.language.shared\_hbm

`nki.language.shared_hbm = Ellipsis`

Shared HBM - Visible to all kernel instances in the SPMD grid

*This document is relevant for:* Inf2, Trn1, Trn2

### Others

|                           |  |
|---------------------------|--|
| <code>program_id</code>   | Index of the current SPMD program along the given axis in the launch grid.                         |
| <code>num_programs</code> | Number of SPMD programs along the given axes in the launch grid.                                   |
| <code>program_ndim</code> | Number of dimensions in the SPMD launch grid.  |
| <code>spmd_dim</code>     | Create a dimension in the SPMD launch grid of a NKI kernel with sub-dimension tiling.              |
| <code>nc</code>           | Create a logical neuron core dimension in launch grid.   |
| <code>device_print</code> | Print a message with a String <code>prefix</code> followed by the value of a tile <code>x</code> . |
| <code>loop_reduce</code>  | Apply reduce operation over a loop.  |

*This document is relevant for:* Inf2, Trn1, Trn2

### nki.language.program\_id

`nki.language.program_id(axis)`

Index of the current SPMD program along the given axis in the launch grid.

#### Parameters

**axis** – The axis of the ND launch grid.

#### Returns

The program id along `axis` in the launch grid

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### nki.language.num\_programs

`nki.language.num_programs(axes=None)`

Number of SPMD programs along the given axes in the launch grid. If `axes` is not provided, returns the total number of programs.

#### Parameters

**axes** – The axes of the ND launch grid. If not provided, returns the total number of programs along the entire launch grid.

#### Returns

The number of SPMD(single process multiple data) programs along `axes` in the launch grid

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.language.program\_ndim

`nki.language.program_ndim()`

Number of dimensions in the SPMD launch grid.

### Returns

The number of dimensions in the launch grid, i.e. the number of axes

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.language.spmd\_dim

`nki.language.spmd_dim = Ellipsis`

Create a dimension in the SPMD launch grid of a NKI kernel with sub-dimension tiling.

A key use case for `spmd_dim` is to shard an existing NKI kernel over multiple NeuronCores without modifying the internal kernel implementation. Suppose we have a kernel, `nki_spmd_kernel`, which is launched with a 2D SPMD grid, (4, 2). We can shard the first dimension of the launch grid (size 4) over two physical NeuronCores by directly manipulating the launch grid as follows:

```
import neuronxcc.nki.language as nl

@nki.jit
def nki_spmd_kernel(a):
    b = nl.ndarray(a.shape, dtype=a.dtype, buffer=nl.shared_hbm)
    i = nl.program_id(0)
    j = nl.program_id(1)

    a_tile = nl.load(a[i, j])
    nl.store(b[i, j], a_tile)

    return b

#####
# Example 1: Let compiler decide how to distribute the instances of spmd kernel
#####
dst = nki_spmd_kernel[4, 2](src)

#####
# Example 2: Distribute SPMD kernel instances to physical NeuronCores with
# explicit annotations. Expected physical NeuronCore assignments:
#   Physical NC [0]: kernel[0, 0], kernel[0, 1], kernel[1, 0], kernel[1, 1]
#   Physical NC [1]: kernel[2, 0], kernel[2, 1], kernel[3, 0], kernel[3, 1]
#####
dst = nki_spmd_kernel[nl.spmd_dim(nl.nc(2), 2), 2](src)
dst = nki_spmd_kernel[nl.nc(2) * 2, 2](src) # syntactic sugar

#####
# Example 3: Distribute SPMD kernel instances to physical NeuronCores with
# explicit annotations. Expected physical NeuronCore assignments:
#   Physical NC [0]: kernel[0, 0], kernel[0, 1], kernel[2, 0], kernel[2, 1]
```

(continues on next page)

(continued from previous page)

```
# Physical NC [1]: kernel[1, 0], kernel[1, 1], kernel[3, 0], kernel[3, 1]
#####
dst = nki_spmc_kernel[nl.spmc_dim(2, nl.nc(2)), 2](src)
dst = nki_spmc_kernel[2 * nl.nc(2), 2](src) # syntactic sugar
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.nc

### nki.language.nc = Ellipsis

Create a logical neuron core dimension in launch grid.

The instances of spmd kernel will be distributed to different physical neuron cores on the annotated dimension.

```
# Let compiler decide how to distribute the instances of spmd kernel
c = kernel[2, 2](a, b)

import neuronxcc.nki.language as nl

# Distribute the kernel to physical neuron cores around the first dimension
# of the spmd grid.
c = kernel[nl.nc(2), 2](a, b)
# This means:
# Physical NC [0]: kernel[0, 0], kernel[0, 1]
# Physical NC [1]: kernel[1, 0], kernel[1, 1]
```

Sometimes the size of a spmd dimension is bigger than the number of available physical neuron cores. We can control the distribution with the following syntax:

```
import neuronxcc.nki.language as nl

@nki.jit
def nki_spmc_kernel(a):
    b = nl.ndarray(a.shape, dtype=a.dtype, buffer=nl.shared_hbm)
    i = nl.program_id(0)
    j = nl.program_id(1)

    a_tile = nl.load(a[i, j])
    nl.store(b[i, j], a_tile)

    return b

#####
# Example 1: Let compiler decide how to distribute the instances of spmd kernel
#####
dst = nki_spmc_kernel[4, 2](src)

#####
# Example 2: Distribute SPMD kernel instances to physical NeuronCores with
```

(continues on next page)

(continued from previous page)

```
# explicit annotations. Expected physical NeuronCore assignments:
#   Physical NC [0]: kernel[0, 0], kernel[0, 1], kernel[1, 0], kernel[1, 1]
#   Physical NC [1]: kernel[2, 0], kernel[2, 1], kernel[3, 0], kernel[3, 1]
#####
dst = nki_spmd_kernel[nl.spmd_dim(nl.nc(2), 2), 2](src)
dst = nki_spmd_kernel[nl.nc(2) * 2, 2](src) # syntactic sugar

#####
# Example 3: Distribute SPMD kernel instances to physical NeuronCores with
# explicit annotations. Expected physical NeuronCore assignments:
#   Physical NC [0]: kernel[0, 0], kernel[0, 1], kernel[2, 0], kernel[2, 1]
#   Physical NC [1]: kernel[1, 0], kernel[1, 1], kernel[3, 0], kernel[3, 1]
#####
dst = nki_spmd_kernel[nl.spmd_dim(2, nl.nc(2)), 2](src)
dst = nki_spmd_kernel[2 * nl.nc(2), 2](src) # syntactic sugar
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.device\_print

`nki.language.device_print(prefix, x, *, mask=None, **kwargs)`

Print a message with a String `prefix` followed by the value of a tile `x`. Printing is currently only supported in kernel simulation mode (see [nki.simulate\\_kernel](#) for a code example).

### Parameters

- **prefix** – prefix of the print message
- **x** – data to print out
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

None

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.language.loop\_reduce

`nki.language.loop_reduce(x, op, loop_indices, *, dtype=None, mask=None, **kwargs)`

Apply reduce operation over a loop. This is an ideal instruction to compute a high performance `reduce_max` or `reduce_min`.

Note: The destination tile is also the rhs input to `op`. For example,

```
b = nl.zeros((N_TILE_SIZE, M_TILE_SIZE), dtype=float32, buffer=nl.sbuf)
for k_i in affine_range(NUM_K_BLOCKS):

    # Skipping over multiple nested loops here.
    # a, is a psum tile from a matmul accumulation group.
    b = nl.loop_reduce(a, op=np.add, loop_indices=[k_i], dtype=nl.float32)
```

is the same as:

```

b = nl.zeros((N_TILE_SIZE, M_TILE_SIZE), dtype=nl.float32, buffer=nl.sbuf)
for k_i in affine_range(NUM_K_BLOCKS):

    # Skipping over multiple nested loops here.
    # a, is a psum tile from a matmul accumulation group.
    b = nisa.tensor_tensor(data1=b, data2=a, op=np.add, dtype=nl.float32)

```

If you are trying to use this instruction only for accumulating results on SBUF, consider simply using the += operator instead.

The `loop_indices` list enables the compiler to recognize which loops this reduction can be optimized across as part of any aggressive loop-level optimizations it may perform.

#### Parameters

- **x** – a tile.
- **op** – numpy ALU operator to use to reduce over the input tile.
- **loop\_indices** – a single loop index or a tuple of loop indices along which the reduction operation is performed. Can be numbers or `loop_index` objects coming from `nl.affine_range`.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

the reduced resulting tile

*This document is relevant for:* Inf2, Trn1, Trn2

## Data Types

|                          |  |
|--------------------------|--|
| <code>tfloat32</code>    | 32-bit floating-point number (1S,8E,10M) |
| <code>bfloat16</code>    | 16-bit floating-point number (1S,8E,7M)  |
| <code>float8_e4m3</code> | 8-bit floating-point number (1S,4E,3M)   |
| <code>float8_e5m2</code> | 8-bit floating-point number (1S,5E,2M)   |

*This document is relevant for:* Inf2, Trn1, Trn2

## `nki.language.tfloat32`

```
nki.language.tfloat32 = dtype('V4')
```

32-bit floating-point number (1S,8E,10M)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



**nki.language.bfloat16**

```
nki.language.bfloat16 = dtype(bfloat16)
```

16-bit floating-point number (1S,8E,7M)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.float8\_e4m3**

```
nki.language.float8_e4m3 = dtype(float8_e4m3)
```

8-bit floating-point number (1S,4E,3M)

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.float8\_e5m2**

```
nki.language.float8_e5m2 = dtype(float8_e5m2)
```

8-bit floating-point number (1S,5E,2M)

*This document is relevant for:* Inf2, Trn1, Trn2

**Constants**

|                  |                      |
|------------------|----------------------|
| <i>tile_size</i> | Tile size constants. |
| <i>fp32</i>      | FP32 Constants       |

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.language.tile\_size**

```
class nki.language.tile_size
```

Tile size constants.

**Attributes**

|                                   |   |
|-----------------------------------|---|
| <code>bn_stats_fmax</code>        | Maximum free dimension of BN_STATS  |
| <code>gemm_moving_fmax</code>     | Maximum free dimension of the moving operand of General Matrix Multiplication on Tensor Engine.     |
| <code>gemm_stationary_fmax</code> | Maximum free dimension of the stationary operand of General Matrix Multiplication on Tensor Engine. |
| <code>pmax</code>                 | Maximum partition dimension of a tile.  |
| <code>psum_fmax</code>            | Maximum free dimension of a tile on PSUM buffer.  |
| <code>psum_min_align</code>       | The minimum byte alignment requirement for PSUM free dimension address.                             |
| <code>sbuf_min_align</code>       | The minimum byte alignment requirement for SBUF free dimension address.                             |

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.language.fp32**

**class** `nki.language.fp32`

FP32 Constants

**Attributes**

|                  |   |
|------------------|---|
| <code>min</code> | FP32 Bit pattern (0xff7fffff) representing the minimum (or maximum negative) FP32 value |
|------------------|---|

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa**

## NKI ISA

|                                |   |
|--------------------------------|---|
| <i>nc_matmul</i>               | Compute stationary.T @ moving matrix multiplication using Tensor Engine.  |
| <i>nc_transpose</i>            | Perform a 2D transpose between the partition axis and the free axis of input data, i.e., a PF-transpose, using Tensor or Vector Engine.   |
| <i>activation</i>              | Apply an activation function on every element of the input tile using Scalar Engine.  |
| <i>activation_reduce</i>       | Perform the same computation as <code>nisa.activation</code> and also a reduction along the free dimension of the <code>nisa.activation</code> result using Scalar Engine.                              |
| <i>tensor_reduce</i>           | Apply a reduction operation to the free axes of an input data tile using Vector Engine.   |
| <i>tensor_partition_reduce</i> | Apply a reduction operation across partitions of an input data tile using GpSimd Engine.  |
| <i>tensor_tensor</i>           | Perform an element-wise operation of input two tiles using Vector Engine or GpSimd Engine.  |
| <i>tensor_tensor_scan</i>      | Perform a scan operation of two input tiles using Vector Engine.  |
| <i>scalar_tensor_tensor</i>    | Apply up to two math operators using Vector Engine: (data <op0> operand0) <op1> operand1.   |
| <i>tensor_scalar</i>           | Apply up to two math operators to the input data tile by broadcasting scalar/vector operands in the free dimension using Vector or Scalar or GpSimd Engine: (data <op0> operand0) <op1> operand1.       |
| <i>tensor_scalar_reduce</i>    | Perform the same computation as <code>nisa.tensor_scalar</code> with one math operator and also a reduction along the free dimension of the <code>nisa.tensor_scalar</code> result using Vector Engine. |
| <i>tensor_copy</i>             | Create a copy of <code>src</code> tile within NeuronCore on-chip SRAMs using Vector, Scalar or GpSimd Engine.   |
| <i>tensor_copy_dynamic_src</i> | Create a copy of <code>src</code> tile within NeuronCore on-chip SRAMs using Vector or Scalar or GpSimd Engine, with <code>src</code> located at a dynamic offset within each partition.                |
| <i>tensor_copy_dynamic_dst</i> | Create a copy of <code>src</code> tile within NeuronCore on-chip SRAMs using Vector or Scalar or GpSimd Engine, with <code>dst</code> located at a dynamic offset within each partition.                |
| <i>tensor_copy_predicated</i>  | Conditionally copy elements from the <code>src</code> tile to the destination tile on SBUF / PSUM based on a <code>predicate</code> using Vector Engine.  |
| <i>reciprocal</i>              | Compute reciprocal of each element in the input data tile using Vector Engine.  |
| <i>iota</i>                    | Build a constant literal in SBUF using GpSimd Engine, rather than transferring the constant literal values from the host to device.   |
| <i>dropout</i>                 | Randomly replace some elements of the input tile data with zeros based on input probabilities using Vector Engine.  |
| <i>affine_select</i>           | Select elements between an input tile <code>on_true_tile</code> and a scalar value <code>on_false_value</code> according to a boolean predicate tile using GpSimd Engine.                               |
| <i>range_select</i>            | Select elements from <code>on_true_tile</code> based on comparison with bounds using Vector Engine.   |

## 7.2. Neuron Kernel Interface (NKI) - Beta

Initialize a tile filled with a compile-time constant value using Vector or GpSimd Engine.

*bn\_stats*

Compute mean- and variance-related statistics for each partition of an input tile data in parallel using Vector Engine.

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.nc\_matmul

```
nki.isa.nc_matmul(stationary, moving, *, is_stationary_onezero=False, is_moving_onezero=False,
                  is_transpose=False, tile_position=(), tile_size=(), mask=None, **kwargs)
```

Compute  $\text{stationary.T} @ \text{moving}$  matrix multiplication using Tensor Engine.

The `nc_matmul` instruction *must* read inputs from SBUF and write outputs to PSUM. Therefore, the `stationary` and `moving` must be SBUF tiles, and the result tile is a PSUM tile.

The `nc_matmul` instruction currently supports `float8_e4m3`/`float8_e5m2`/`bfloat16`/`float16`/`tfloat32`/`float32` input data types as listed in [Supported Data Types](#). The `matmul` accumulation and results are always in `float32`.

The Tensor Engine imposes special layout constraints on the input tiles. First, the partition axis sizes of the `stationary` and `moving` tiles must be identical and  $\leq 128$ , which corresponds to the contraction dimension of the matrix multiplication. Second, the free axis sizes of `stationary` and `moving` tiles must be  $\leq 128$  and  $\leq 512$ , respectively. For example, `stationary.shape = (128, 126)`; `moving.shape = (128, 512)` and `nc_matmul(stationary, moving)` returns a tile of `shape = (126, 512)`. For more information about the `matmul` layout, see [Tensor Engine](#).

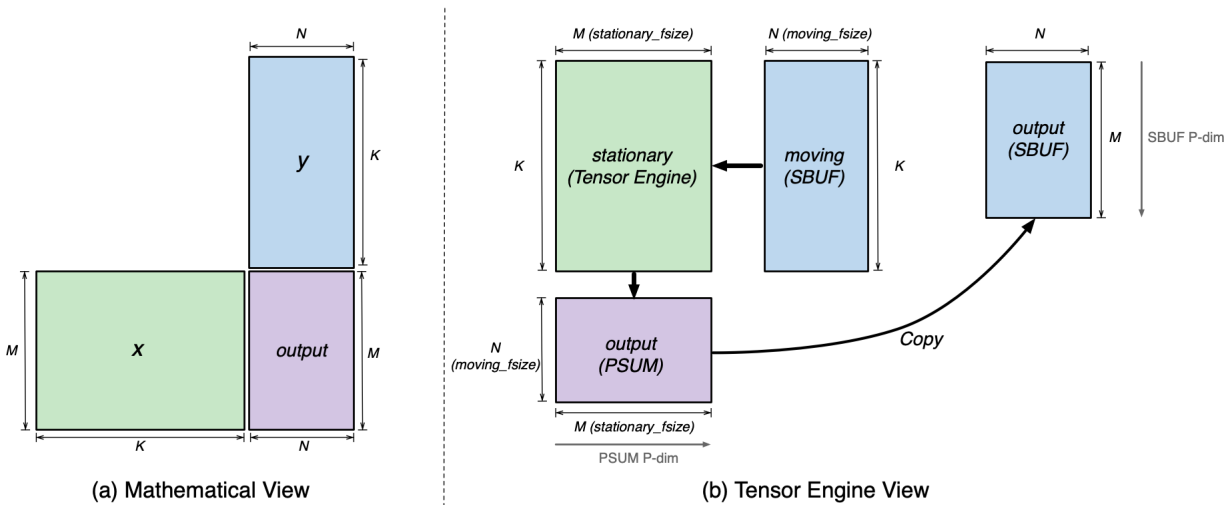


Fig. 7.1:  $M \times K \times N$  Matrix Multiplication Visualization.

If the contraction dimension of the matrix multiplication exceeds 128, you may accumulate multiple `nc_matmul` instruction output tiles into the same PSUM tile. See example code snippet below.

### Estimated instruction cost:

The Tensor Engine has complex performance characteristics given its data flow and pipeline design. The below formula is the *average* `nc_matmul` cost assuming many `nc_matmul` instructions of the same shapes running back-to-back on the engine:

| Cost ( <i>Tensor Engine Cycles</i> )                           | Condition   |
|--|---|
| $\max(\min(64, N_{\text{stationary}}), N_{\text{moving}})$     | input data type is one of float8_e4m3/float8_e5m2/bfloat16/float16/tfloat32 |
| $4 * \max(\min(64, N_{\text{stationary}}), N_{\text{moving}})$ | input data type is float32  |

where,

- **N\_stationary** is the number of elements per partition in stationary tile.
- **N\_moving** is the number of elements per partition in moving tile.

The Tensor Engine, as a systolic array with 128 rows and 128 columns of processing elements (PEs), could be underutilized for small `nc_matmul` instructions, i.e., the stationary tile has small free axis size or small partition axis size (e.g. 32, 64). In such a case, the Tensor Engine allows PE tiling, i.e., multiple small `nc_matmul` instructions to execute in parallel on the PE array, to improve compute throughput. PE tiling is enabled by setting `tile_position` and `tile_size`. `tile_position` indicates the PE tile starting position (row position, column position) for a `nc_matmul` instruction in the PE array. `tile_size` indicates the PE tile size (row size, column size) to hold by a `nc_matmul` instruction starting from the `tile_position`. For example, setting `tile_position` to (0, 0) and `tile_size` to (128, 128) means using full PE array.

Requirements on `tile_position` and `tile_size` are:

1. `tile_position` and `tile_size` must be both set to enable PE tiling.
2. The type of values in `tile_position` and `tile_size` must be integer or affine expression.
3. Values in `tile_position` and `tile_size` must be multiple of 32.
4. `tile_size` must be larger than or equal to accessed stationary tile size.
5. Both the row and column sizes in `tile_size` cannot be 32 for NeuronCore-v2.

#### Parameters

- **stationary** – the stationary operand on SBUF; layout: (partition axis <= 128, free axis <= 128)
- **moving** – the moving operand on SBUF; layout: (partition axis <= 128, free axis <= 512)
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **is\_stationary\_onezero** – hints to the compiler whether the stationary operand is a tile with ones/zeros only; setting this field explicitly could lead to 2x better performance if stationary tile is in float32; the field has no impact for non-float32 stationary.
- **is\_moving\_onezero** – hints to the compiler if the moving operand is a tile with ones/zeros only; setting this field explicitly could lead to 2x better performance if moving tile is in float32; the field has no impact for non-float32 moving.
- **is\_transpose** – hints to the compiler that this is a transpose operation with moving as an identity matrix.
- **tile\_position** – a 2D tuple (row, column) for the start PE tile position to run `nc_matmul`.
- **tile\_size** – a 2D tuple (row, column) for the PE tile size to hold by `nc_matmul` starting from `tile_position`.

#### Returns

a tile on PSUM that has the result of matrix multiplication of `stationary` and `moving` tiles; layout: partition axis comes from free axis of `stationary`, while free axis comes from free axis of `moving`.

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
```

(continues on next page)

(continued from previous page)

```
#####
# Example 1:
# multiply matrix a of shape (128, 128) and matrix b of shape (128, 512)
# to get matrix c in PSUM of shape (128, 512)
#####
a_mgrid = nl.mgrid[0:128, 0:128]
b_mgrid = nl.mgrid[0:128, 0:512]
c_mgrid = nl.mgrid[0:128, 0:512]

a = nl.load(a_tensor[a_mgrid.p, a_mgrid.x])
b = nl.load(b_tensor[b_mgrid.p, b_mgrid.x])

c_psum = nisa.nc_matmul(a[a_mgrid.p, a_mgrid.x], b[b_mgrid.p, b_mgrid.x])

nl.store(c_tensor[c_mgrid.p, c_mgrid.x], c_psum)

#####
# Example 2:
# multiply matrix d of shape (256, 128) and matrix e of shape (256, 512)
# to get matrix f in PSUM of shape (128, 512) using psum accumulation
#####
d_mgrid = nl.mgrid[0:128, 0:128]
e_mgrid = nl.mgrid[0:128, 0:512]
f_mgrid = nl.mgrid[0:128, 0:512]

f_psum = nl.zeros((128, 512), nl.float32, buffer=nl.psum)

for i_contract in nl.affine_range(2):
    d = nl.load(d_tensor[i_contract * 128 + d_mgrid.p, d_mgrid.x])
    e = nl.load(e_tensor[i_contract * 128 + e_mgrid.p, e_mgrid.x])
    f_psum += nisa.nc_matmul(d[d_mgrid.p, d_mgrid.x], e[e_mgrid.p, e_mgrid.x])

nl.store(f_tensor[f_mgrid.p, f_mgrid.x], f_psum)

#####
# Example 3:
# perform batched matrix multiplication on matrix g of shape (16, 64, 64)
# and matrix h of shape (16, 64, 512) to get matrix i of (16, 64, 512)
# using Tensor Engine PE tiling mode.
#####
g_mgrid = nl.mgrid[0:64, 0:64]
h_mgrid = nl.mgrid[0:64, 0:512]
i_mgrid = nl.mgrid[0:64, 0:512]

for i in nl.affine_range(4):
    for j in nl.affine_range(4):
        g = nl.load(g_tensor[i * 4 + j, g_mgrid.p, g_mgrid.x])
        h = nl.load(h_tensor[i * 4 + j, h_mgrid.p, h_mgrid.x])
        i_psum = nisa.nc_matmul(g, h, tile_position=((i % 2) * 64, (j % 2) * 64), tile_
↪size=(64, 64))
        nl.store(i_tensor[i * 4 + j, i_mgrid.p, i_mgrid.x], i_psum)
```

(continues on next page)

(continued from previous page)

```
return c_tensor, f_tensor, i_tensor
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.nc\_transpose

`nki.isa.nc_transpose(data, *, mask=None, dtype=None, engine=engine.unknown, **kwargs)`

Perform a 2D transpose between the partition axis and the free axis of input `data`, i.e., a PF-transpose, using Tensor or Vector Engine. If the `data` tile has more than one free axes, this API implicitly collapses all free axes into one axis and then performs a 2D PF-transpose.

In NeuronCore, both Tensor and Vector Engine can perform a PF-transpose, but they support different input shapes. Tensor Engine `nc_transpose` can handle an input tile of shape (128, 128) or smaller, while Vector Engine can handle shape (32, 32) or smaller. Therefore, when the input tile shape is (32, 32) or smaller, we have an option to run it on either engine, which is controlled by the `engine` field. If no `engine` is specified, Neuron Compiler will automatically select an engine based on the input shape. Note, similar to other Tensor Engine instructions, the Tensor Engine `nc_transpose` must read the input tile from SBUF and write the transposed result to PSUM. On the other hand, Vector Engine `nc_transpose` can read/write from/to either SBUF or PSUM.

Note, PF-transpose on Tensor Engine is done by performing a matrix multiplication between `data` as the stationary tensor and an identity matrix as the moving tensor. See [architecture guide](#) for more information. On NeuronCore-v2, such matmul-style transpose is not bit-accurate if the input `data` contains NaN/Inf. You may consider replacing NaN/Inf with regular floats (`float_max/float_min/zeros`) in the input matrix before calling `nc_transpose(engine=nki.isa.constants.engine.tensor)`.

### Estimated instruction cost:

| Cost ( <i>Engine Cycles</i> )   | Condition  |
|---------------------------------|--|
| <code>max(MIN_II, N)</code>     | engine set to <code>nki.isa.constants.engine.vector</code>   |
| <code>max(P, min(64, F))</code> | engine set to <code>nki.isa.constants.engine.tensor</code> and assuming many back-to-back <code>nc_transpose</code> of the same shape on Tensor Engine |

where,

- `N` is the number of elements per partition in `data`.
- `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.
- `P` is partition axis size of `data`.
- `F` is the number of elements per partition in `data`.

### Parameters

- **`data`** – the input tile to be transposed
- **`mask`** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **`dtype`** – if specified and it's different from the data type of input tile `data`, an additional `nki.isa.cast` instruction will be inserted to cast the transposed data into the target `dtype` (see [Supported Data Types](#) for more information)
- **`engine`** – specify which engine to use for transpose: `nki.isa.tensor_engine` or `nki.isa.vector_engine`; by default, the best engine will be selected for the given input tile shape

**Returns**

a tile with transposed result of input data tile

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
# Example 1: transpose tile a of shape (128, 64)
#####
i_p_a = nl.arange(128)[: , None]
i_f_a = nl.arange(64)[None, :]
aT = nisa.nc_transpose(a[i_p_a, i_f_a])

#####
# Example 2: transpose tile b of shape (32, 2) using Vector Engine
#####
i_p_b = nl.arange(32)[: , None]
i_f_b = nl.arange(2)[None, :]
bT = nisa.nc_transpose(b[i_p_b, i_f_b], engine=nisa.vector_engine)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

**nki.isa.activation**

`nki.isa.activation`(*op*, *data*, \*, *bias*=None, *scale*=1.0, *reduce\_op*=None, *reduce\_res*=None, *reduce\_cmd*=*reduce\_cmd.idle*, *mask*=None, *dtype*=None, \*\*kwargs)

Apply an activation function on every element of the input tile using Scalar Engine. The activation function is specified in the *op* input field (see [Supported Activation Functions for NKI ISA](#) for a list of supported activation functions and their valid input ranges).

The activation instruction can optionally multiply the input data by a scalar or vector *scale* and then add another vector *bias* before the activation function is applied, at no additional performance cost:

$$output = f_{act}(data * scale + bias)$$

When the *scale* is a scalar, it must be a compile-time constant. In this case, the *scale* is broadcasted to all the elements in the input data tile. When the *scale/bias* is a vector, it must have the same partition axis size as the input data tile and only one element per partition. In this case, the element of *scale/bias* within each partition is broadcasted to elements of the input data tile in the same partition.

There are 128 registers on the scalar engine for storing reduction results, corresponding to the 128 partitions of the input. The scalar engine can reduce along free dimensions without extra performance penalty, and store the result of reduction into these registers. The reduction is done after the activation function is applied.

$$output = f_{act}(data * scale + bias) \text{ } accu\_registers = reduce\_op(accu\_registers, reduce\_op(output, axis = \langle FreeAxis \rangle)$$

These registers are shared between `activation` and `activation_accu` calls, and the state of them can be controlled via the `reduce_cmd` parameter.

- `nisa.reduce_cmd.reset`: Reset the accumulators to zero



- `nisa.reduce_cmd.idle`: Do not use the accumulators
- `nisa.reduce_cmd.reduce`: keeps accumulating over the current value of the accumulator
- `nisa.reduce_cmd.reset_reduce`: Resets the accumulators then immediately accumulate the results of the current instruction into the accumulators

We can choose to read out the current values stored in the register by passing in a tensor in the `reduce_res` arguments. Reading out the accumulator will incur a small overhead.

Note that `activation_accu` can also change the state of the registers. It's user's responsibility to ensure correct ordering. It's recommended to not mixing the use of `activation_accu` and `activation`, when `reduce_cmd` is not set to `idle`.

Note, the Scalar Engine always performs the math operations in float32 precision. Therefore, the engine automatically casts the input data tile to float32 before performing multiply/add/activate specified in the activation instruction. The engine is also capable of casting the float32 math results into another output data type specified by the `dtype` field at no additional performance cost. If `dtype` field is not specified, Neuron Compiler will set output data type of the instruction to be the same as input data type of data. On the other hand, the `scale` parameter must have a float32 data type, while the `bias` parameter can be float32/float16/bfloat16.

The input data tile can be an SBUF or PSUM tile. Similarly, the instruction can write the output tile into either SBUF or PSUM, which is specified using the `buffer` field. If not specified, `nki.language.sbuf` is selected by default.

#### Estimated instruction cost:

$\max(\text{MIN\_II}, N)$  Scalar Engine cycles, where

- $N$  is the number of elements per partition in data.
- $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

#### Parameters

- **op** – an activation function (see [Supported Activation Functions for NKI ISA](#) for supported functions)
- **data** – the input tile; layout: (partition axis  $\leq 128$ , free axis)
- **bias** – a vector with the same partition axis size as `data` for broadcast add (after broadcast multiply with `scale`)
- **scale** – a scalar or a vector with the same partition axis size as `data` for broadcast multiply
- **reduce\_op** – the reduce operation to perform on the free dimension of the activation result
- **reduce\_res** – a tile of shape  $(\text{data.shape}[0], 1)$ , where  $\text{data.shape}[0]$  is the partition axis size of the input data tile. The result of `sum(ReductionResult)` is written in-place into the tensor.
- **reduce\_cmd** – an enum member from `nisa.reduce_cmd` to control the state of reduction registers
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

output tile of the activation instruction; layout: same as input data tile

Example:

```
import neuronxcc.nki.language as nl
import neuronxcc.nki.isa as nisa
```

(continues on next page)

(continued from previous page)

```
#####
# Example 1: perform exponential function on matrix a of shape (128, 1024)
#####
a = nl.load(a_tensor)
activated_a = nisa.activation(op=nl.exp, data=a)
nl.store(a_act_tensor, activated_a)

#####
# Example 2: perform the following operations to matrix b of shape (128, 512)
# using a single activation instruction: np.square(b * 2.0) + c
# 1) compute `np.square(b * 2.0 + c)`
# 2) cast 1) results into bfloat16
#####
b = nl.load(b_tensor)
c = nl.load(c_tensor)
activated_b = nisa.activation(op=np.square, data=b, bias=c, scale=2.0,
                             dtype=nl.bfloat16)
nl.store(b_act_tensor, activated_b)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

### nki.isa.activation\_reduce

`nki.isa.activation_reduce(op, data, *, reduce_op, reduce_res, bias=None, scale=1.0, mask=None, dtype=None, **kwargs)`

Perform the same computation as `nisa.activation` and also a reduction along the free dimension of the `nisa.activation` result using Scalar Engine. The results for the reduction is stored in the `reduce_res`.

This API is equivalent to calling `nisa.activation` with `reduce_cmd=nisa.reduce_cmd.reset_reduce` and passing in `reduce_res`. This API is kept for backward compatibility, we recommend using `nisa.activation` moving forward.

Refer to [nisa.activation](#) for semantics of `op/data/bias/scale`.

In addition to [nisa.activation](#) computation, this API also performs a reduction along the free dimension(s) of the [nisa.activation](#) result, at a small additional performance cost. The reduction result is returned in `reduce_res` in-place, which must be a SBUF/PSUM tile with the same partition axis size as the input tile data and one element per partition. On NeuronCore-v2, the `reduce_op` can only be an addition, `np.add` or `nl.add`.

There are 128 registers on the scalar engine for storing reduction results, corresponding to the 128 partitions of the input. These registers are shared between `activation` and `activation_accu` calls. This instruction first resets those registers to zero, performs the reduction on the value after activation function is applied, stores the results into the registers, then reads out the reduction results from the register, eventually store them into `reduce_res`.

Note that `nisa.activation` can also change the state of the register. It's user's responsibility to ensure correct ordering. It's the best practice to not mixing the use of `activation_reduce` and `activation`.

Reduction axis is not configurable in this API. If the input tile has multiple free axis, the API will reduce across all of them.

Mathematically, this API performs the following computation:

$$\begin{aligned} output &= f_{act}(data * scale + bias) \\ reduce\_res &= reduce\_op(output, axis = < FreeAxis >) \end{aligned}$$

#### Estimated instruction cost:

$\max(\text{MIN\_II}, N) + \text{MIN\_II}$  Scalar Engine cycles, where

- $N$  is the number of elements per partition in `data`, and
- $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

#### Parameters

- **op** – an activation function (see [Supported Activation Functions for NKI ISA](#) for supported functions)
- **data** – the input tile; layout: (partition axis  $\leq 128$ , free axis)
- **reduce\_op** – the reduce operation to perform on the free dimension of the activation result
- **reduce\_res** – a tile of shape  $(\text{data.shape}[0], 1)$ , where  $\text{data.shape}[0]$  is the partition axis size of the input data tile. The result of `sum(ReductionResult)` is written in-place into the tensor.
- **bias** – a vector with the same partition axis size as `data` for broadcast add (after broadcast multiply with `scale`)
- **scale** – a scalar or a vector with the same partition axis size as `data` for broadcast multiply
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

output tile of the activation instruction; layout: same as input data tile

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### `nki.isa.tensor_reduce`

`nki.isa.tensor_reduce(op, data, axis, *, mask=None, dtype=None, negate=False, keepdims=False, **kwargs)`

Apply a reduction operation to the free axes of an input data tile using Vector Engine.

The reduction operator is specified in the `op` input field (see [Supported Math Operators for NKI ISA](#) for a list of supported reduction operators). There are two types of reduction operators: 1) bitvec operators (e.g., `bitwise_and`, `bitwise_or`) and 2) arithmetic operators (e.g., `add`, `subtract`, `multiply`). For bitvec operators, the input/output data types must be integer types and Vector Engine treats all input elements as bit patterns without any data type casting. For arithmetic operators, there is no restriction on the input/output data types, but the engine automatically casts input data types to float32 and performs the reduction operation in float32 math. The float32 reduction results are cast to the target data type specified in the `dtype` field before written into the output tile. If the `dtype` field is not specified, it is default to be the same as input tile data type.

When the reduction `op` is an arithmetic operator, the instruction can also multiply the output reduction results by `-1.0` before writing into the output tile, at no additional performance cost. This behavior is controlled by the `negate` input field.

The reduction axes are specified in the `axis` field using a list of integer(s) to indicate axis indices. The reduction axes can contain up to four free axes and must start at the most minor free axis. Since axis 0 is the partition axis in

a tile, the reduction axes must contain axis 1 (most-minor). In addition, the reduction axes must be consecutive: e.g., [1, 2, 3, 4] is a legal axis field, but [1, 3, 4] is not.

Since this instruction only supports free axes reduction, the output tile must have the same partition axis size as the input data tile. To perform a partition axis reduction, we can either:

1. invoke a `nki.isa.nc_transpose` instruction on the input tile and then this `reduce` instruction to the transposed tile, or
2. invoke `nki.isa.nc_matmul` instructions to multiply a `nki.language.ones([128, 1], dtype=data.dtype)` vector with the input tile.

**Estimated instruction cost:**

| Cost (Vector Engine Cycles) | Condition   |
|-----------------------------|---|
| N/2                         | both input and output data types are <code>bfloat16</code> and the reduction operator is add or maximum |
| N                           | otherwise   |

where,

- N is the number of elements per partition in data.
- MIN\_II is the minimum instruction initiation interval for small input tiles. MIN\_II is roughly 64 engine cycles.

#### Parameters

- **op** – the reduction operator (see [Supported Math Operators for NKI ISA](#) for supported reduction operators)
- **data** – the input tile to be reduced
- **axis** – int or tuple/list of ints. The axis (or axes) along which to operate; must be free dimensions, not partition dimension (0); can only be the last contiguous dim(s) of the tile: [1], [1,2], [1,2,3], [1,2,3,4]
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **negate** – if True, reduction result is multiplied by `-1.0`; only applicable when op is an arithmetic operator
- **keepdims** – If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

#### Returns

output tile of the reduction result

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import numpy as np
...

#####
# Example 1: reduce add tile a of shape (128, 512)
# in the free dimension and return
# reduction result in tile b of shape (128, 1)
#####
```

(continues on next page)

(continued from previous page)

```
i_p_a = nl.arange(128)[: , None]
i_f_a = nl.arange(512)[None, :]

b = nisa.tensor_reduce(np.add, a[i_p_a, i_f_a], axis=[1])
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.isa.tensor\_partition\_reduce

`nki.isa.tensor_partition_reduce(op, data, *, mask=None, dtype=None, **kwargs)`

Apply a reduction operation across partitions of an input data tile using GpSimd Engine.

### Parameters

- **op** – the reduction operator (add, max, bitwise\_or, bitwise\_and)
- **data** – the input tile to be reduced
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

### Returns

output tile with reduced result

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import numpy as np
...

#####
# Example 1: reduce add tile a of shape (128, 32, 4)
# in the partition dimension and return
# reduction result in tile b of shape (1, 32, 4)
#####
a = nl.load(a_tensor[0:128, 0:32, 0:4])
b = nisa.tensor_partition_reduce(np.add, a)
nl.store(b_tensor[0:1, 0:32, 0:4], b)

#####
# Example 2: reduce add tile a of shape (b, p, f1, ...)
# in the partition dimension p and return
# reduction result in tile b of shape (b, 1, f1, ...)
#####
for i in nl.affine_range(a_tensor.shape[0]):
    a = nl.load(a_tensor[i])
    b = nisa.tensor_partition_reduce(np.add, a)
    nl.store(b_tensor[i], b)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.tensor\_tensor**

`nki.isa.tensor_tensor(data1, data2, op, *, dtype=None, mask=None, engine=engine.unknown, **kwargs)`

Perform an element-wise operation of input two tiles using Vector Engine or GpSimd Engine. The two tiles must have the same partition axis size and the same number of elements per partition.

The element-wise operator is specified using the `op` field and can be any *binary* operator supported by NKI (see [Supported Math Operators for NKI ISA](#) for details) that runs on the Vector Engine, or can be `np.power/nl.power` that runs on the GpSimd Engine. For bitvec operators, the input/output data types must be integer types and Vector Engine treats all input elements as bit patterns without any data type casting. For arithmetic operators, there is no restriction on the input/output data types, but the engine automatically casts input data types to float32 and performs the element-wise operation in float32 math. The float32 results are cast to the target data type specified in the `dtype` field before written into the output tile. If the `dtype` field is not specified, it is default to be the same as the data type of `data1` or `data2`, whichever has the higher precision.

Since GpSimd Engine cannot access PSUM, the input or output tiles cannot be in PSUM if `op` is `np.power/nl.power` (see [NeuronCore-v2 Compute Engines](#) for details). Otherwise, the output tile can be in either SBUF or PSUM. However, the two input tiles, `data1` and `data2` cannot both reside in PSUM. The three legal cases are:

1. Both `data1` and `data2` are in SBUF.
2. `data1` is in SBUF, while `data2` is in PSUM.
3. `data1` is in PSUM, while `data2` is in SBUF.

Note, if you need broadcasting capability in the free dimension for either input tile, you should consider using [nki.isa.tensor\\_scalar](#) API instead, which has better performance than `nki.isa.tensor_tensor` in general.

**Estimated instruction cost:**

See below table for `tensor_tensor` performance when it runs on Vector Engine.

| Cost ( <i>Vector Engine Cycles</i> ) | Condition  |
|--------------------------------------|--|
| $\max(\text{MIN\_II}, N)$            | one input tile is in PSUM and the other is in SBUF   |
| $\max(\text{MIN\_II}, N)$            | all of the below: <ul style="list-style-type: none"> <li>• both input tiles are in SBUF,</li> <li>• input/output data types are all <code>bfloat16</code>,</li> <li>• the operator is add, multiply or subtract,</li> <li>• Input tensor data is contiguous along the free dimension (that is, stride in each partition is 1 element)</li> </ul> |
| $\max(\text{MIN\_II}, 2N)$           | otherwise  |

where,

- `N` is the number of elements per partition in `data1/data2`.
- `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

**Parameters**

- **`data1`** – lhs input operand of the element-wise operation
- **`data2`** – rhs input operand of the element-wise operation
- **`op`** – a binary math operator (see [Supported Math Operators for NKI ISA](#) for supported operators)
- **`mask`** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **`dtype`** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);

- **engine** – (optional) the engine to use for the operation: `nki.isa.vector_engine`, `nki.isa.gpsimd_engine` or `nki.isa.unknown_engine` (default, let compiler select best engine based on the input tile shape).

#### Returns

an output tile of the element-wise operation

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor
...

#####
# Example 1: add two tiles, a and b, of the same
# shape (128, 512) element-wise and get
# the addition result in tile c
#####
a: tensor[128, 512] = nl.load(a_tensor)
b: tensor[128, 512] = nl.load(b_tensor)

c: tensor[128, 512] = nisa.tensor_tensor(a, b, op=nl.add)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

### nki.isa.tensor\_tensor\_scan

`nki.isa.tensor_tensor_scan(data0, data1, initial, op0, op1, reverse0=False, reverse1=False, *, dtype=None, mask=None, **kwargs)`

Perform a scan operation of two input tiles using Vector Engine.

Mathematically, the `tensor_tensor_scan` instruction on Vector Engine performs the following computation per partition:

```
# Let's assume we work with numpy, and data0 and data1 are 2D (with shape[0] being
↳ the partition axis)
import numpy as np

result = np.ndarray(data0.shape, dtype=data0.dtype)
result[:, 0] = op1(op0(data0[:, 0], initial), data1[:, 0])

for i in range(1, data0.shape[1]):
    result[:, i] = op1(op0(data0[:, i], result[:, i-1]), data1[:, i])
```

The two input tiles (`data0` and `data1`) must have the same partition axis size and the same number of elements per partition. The third input `initial` can either be a float32 compile-time scalar constant that will be broadcasted in the partition axis of `data0/data1`, or a tile with the same partition axis size as `data0/data1` and one element per partition.

The two input tiles, `data0` and `data1` cannot both reside in PSUM. The three legal cases are:

1. Both `data1` and `data2` are in SBUF.
2. `data1` is in SBUF, while `data2` is in PSUM.

3. `data1` is in PSUM, while `data2` is in SBUF.

The scan operation supported by this API has two programmable math operators in `op0` and `op1` fields. Both `op0` and `op1` can be any binary arithmetic operator supported by NKI (see [Supported Math Operators for NKI ISA](#) for details). We can optionally reverse the input operands of `op0` by setting `reverse0` to True (or `op1` by setting `reverse1`). Reversing operands is useful for non-commutative operators, such as subtract.

Input/output data types can be any supported NKI data type (see [Supported Data Types](#)), but the engine automatically casts input data types to float32 and performs the computation in float32 math. The float32 results are cast to the target data type specified in the `dtype` field before written into the output tile. If the `dtype` field is not specified, it is default to be the same as the data type of `data0` or `data1`, whichever has the highest precision.

#### Estimated instruction cost:

`max(MIN_II, 2N)` Vector Engine cycles, where

- `N` is the number of elements per partition in `data0/data1`.
- `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

#### Parameters

- **`data0`** – lhs input operand of the scan operation
- **`data1`** – rhs input operand of the scan operation
- **`initial`** – starting state of the scan; can be a SBUF/PSUM tile with 1 element/partition or a scalar compile-time constant
- **`op0`** – a binary arithmetic math operator (see [Supported Math Operators for NKI ISA](#) for supported operators)
- **`op1`** – a binary arithmetic math operator (see [Supported Math Operators for NKI ISA](#) for supported operators)
- **`reverse0`** – reverse ordering of inputs to `op0`; if false, `data0` is the lhs of `op0`; if true, `data0` is the rhs of `op0`
- **`reverse1`** – reverse ordering of inputs to `op1`; if false, `data1` is the rhs of `op1`; if true, `data1` is the lhs of `op1`
- **`mask`** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **`dtype`** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);

#### Returns

an output tile of the scan operation

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl

#####
# Example 1: scan two tiles, a and b, of the same
# shape (128, 1024) using multiply/add and get
# the scan result in tile c
#####
c = nl.ndarray(shape=(128, 1024), dtype=nl.float32)

c[:, 0:512] = nisa.tensor_tensor_scan(a[:, 0:512], b[:, 0:512],
                                     initial=0, op0=np.multiply, op1=np.add)

c[:, 512:1024] = nisa.tensor_tensor_scan(a[:, 512:1024], b[:, 512:1024],
```

(continues on next page)



(continued from previous page)

```
initial=c[:, 511],
op0=np.multiply, op1=np.add)
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.isa.scalar\_tensor\_tensor

```
nki.isa.scalar_tensor_tensor(*, data, op0, operand0, op1, operand1, reverse0=False, reverse1=False,
                             dtype=None, mask=None, **kwargs)
```

Apply up to two math operators using Vector Engine: (data <op0> operand0) <op1> operand1.

data input can be an SBUF or PSUM tile of 2D shape. operand0 can be SBUF or PSUM tile of shape (data.shape[0], 1), i.e., vector, or a compile-time constant scalar. operand1 can be SBUF or PSUM tile of shape (data.shape[0], data.shape[1]) (i.e., has to match data shape), note that operand1 and data can't both be on PSUM.

### Estimated instruction cost:

| Cost (Vector Engine Cycles) | Condition   |
|-----------------------------|---|
| N                           | data and operand1 are both bfloat16, op0=nl.subtract and op1=nl.multiply, and N is even |
| 2*N                         | otherwise   |

where,

- N is the number of elements per partition in data.

### Parameters

- **data** – the input tile
- **op0** – the first math operator used with operand0 (see [Supported Math Operators for NKI ISA](#) for supported operators)
- **operand0** – a scalar constant or a tile of shape (data.shape[0], 1), where data.shape[0] is the partition axis size of the input data tile.
- **reverse0** – reverse ordering of inputs to op0; if false, operand0 is the rhs of op0; if true, operand0 is the lhs of op0.
- **op1** – the second math operator used with operand1 (see [Supported Math Operators for NKI ISA](#) for supported operators).
- **operand1** – a tile of shape with the same partition and free dimension as data input.
- **reverse1** – reverse ordering of inputs to op1; if false, operand1 is the rhs of op1; if true, operand1 is the lhs of op1.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

an output tile of (data <op0> operand0) <op1> operand1 computation

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

**nki.isa.tensor\_scalar**

```
nki.isa.tensor_scalar(data, op0, operand0, reverse0=False, op1=None, operand1=None, reverse1=False, *,
                      dtype=None, mask=None, engine=engine.unknown, **kwargs)
```

Apply up to two math operators to the input `data` tile by broadcasting scalar/vector operands in the free dimension using Vector or Scalar or GpSimd Engine: (`data` `<op0>` `operand0`) `<op1>` `operand1`.

The input `data` tile can be an SBUF or PSUM tile. Both `operand0` and `operand1` can be SBUF or PSUM tiles of shape (`data.shape[0]`, 1), i.e., vectors, or compile-time constant scalars.

`op1` and `operand1` are optional, but must be `None` (default values) when unused. Note, performing one operator has the same performance cost as performing two operators in the instruction.

When the operators are non-commutative (e.g., subtract), we can reverse ordering of the inputs for each operator through:

- `reverse0 = True`: `tmp_res = operand0 <op0> data`
- `reverse1 = True`: `operand1 <op1> tmp_res`

The `tensor_scalar` instruction supports two types of operators: 1) bitvec operators (e.g., `bitwise_and`) and 2) arithmetic operators (e.g., `add`). See [Supported Math Operators for NKI ISA](#) for the full list of supported operators. The two operators, `op0` and `op1`, in a `tensor_scalar` instruction must be of the same type (both bitvec or both arithmetic). If bitvec operators are used, the `tensor_scalar` instruction must run on Vector Engine. Also, the input/output data types must be integer types, and input elements are treated as bit patterns without any data type casting.

If arithmetic operators are used, the `tensor_scalar` instruction can run on Vector or Scalar or GpSimd Engine. However, each engine supports limited arithmetic operators (see :ref:tbl-aluop). The Scalar Engine on `trn2` only supports a subset of the operator combination:

- `op0=np.multiply` and `op1=np.add`
- `op0=np.multiply` and `op1=None`
- `op0=add` and `op1=None`

Also, arithmetic operators impose no restriction on the input/output data types, but the engine automatically casts input data types to float32 and performs the operators in float32 math. The float32 computation results are cast to the target data type specified in the `dtype` field before written into the output tile, at no additional performance cost. If the `dtype` field is not specified, it is default to be the same as input tile data type.

**Estimated instruction cost:**

`max(MIN_II, N)` Vector or Scalar Engine cycles, where

- `N` is the number of elements per partition in `data`.
- `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

**Parameters**

- **`data`** – the input tile
- **`op0`** – the first math operator used with `operand0` (see [Supported Math Operators for NKI ISA](#) for supported operators)
- **`operand0`** – a scalar constant or a tile of shape (`data.shape[0]`, 1), where `data.shape[0]` is the partition axis size of the input data tile
- **`reverse0`** – reverse ordering of inputs to `op0`; if false, `operand0` is the rhs of `op0`; if true, `operand0` is the lhs of `op0`
- **`op1`** – the second math operator used with `operand1` (see [Supported Math Operators for NKI ISA](#) for supported operators); this operator is optional
- **`operand1`** – a scalar constant or a tile of shape (`data.shape[0]`, 1), where `data.shape[0]` is the partition axis size of the input data tile
- **`reverse1`** – reverse ordering of inputs to `op1`; if false, `operand1` is the rhs of `op1`; if true, `operand1` is the lhs of `op1`

- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **engine** – (optional) the engine to use for the operation: `nki.isa.vector_engine`, `nki.isa.scalar_engine`, `nki.isa.gpsimd_engine` (only allowed for `rsqrt`) or `nki.isa.unknown_engine` (default, let compiler select best engine based on the input tile shape).

#### Returns

an output tile of (data <op0> operand0) <op1> operand1 computation

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import numpy as np
...

#####
# Example 1: subtract 1.0 from all elements of tile a of
# shape (128, 512) and get the output tile in b
#####
i_p = nl.arange(128)[: , None]
i_f = nl.arange(512)[None, :]

b = nisa.tensor_scalar(a[i_p, i_f], np.subtract, 1.0)

#####
# Example 2: broadcast 1.0 into a shape of (128, 512) and subtract
# it with tile c to get output tile d
#####
i_p = nl.arange(128)[: , None]
i_f = nl.arange(512)[None, :]

d = nisa.tensor_scalar(c[i_p, i_f], np.subtract, 1.0, reverse0=True)

#####
# Example 3: broadcast multiply tile e with vector f and
# then broadcast add with scalar 2.5;
# tile e has a shape of (64, 1024) and vector f has a shape of (64, 1)
#####
i_p_ef = nl.arange(64)[: , None]
i_f_e = nl.arange(1024)[None, :]
i_f_f = nl.arange(1)[None, :]

g = nisa.tensor_scalar(e[i_p_ef, i_f_e], op0=np.multiply, operand0=f[i_p_ef, i_f_f],
↪ op1=np.add, operand1=2.5)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

**nki.isa.tensor\_scalar\_reduce**

```
nki.isa.tensor_scalar_reduce(*, data, op0, operand0, reduce_op, reduce_res, reverse0=False, dtype=None,
                             mask=None, **kwargs)
```

Perform the same computation as `nisa.tensor_scalar` with one math operator and also a reduction along the free dimension of the `nisa.tensor_scalar` result using Vector Engine.

Refer to [nisa.tensor\\_scalar](#) for semantics of `data/op0/operand0`. Unlike regular `nisa.tensor_scalar` where two operators are supported, only one operator is supported in this API. Also, `op0` can only be arithmetic operation in [Supported Math Operators for NKI ISA](#). Bitvec operators are not supported in this API.

In addition to [nisa.tensor\\_scalar](#) computation, this API also performs a reduction along the free dimension(s) of the [nisa.tensor\\_scalar](#) result, at a small additional performance cost. The reduction result is returned in `reduce_res` in-place, which must be a SBUF/PSUM tile with the same partition axis size as the input tile `data` and one element per partition. The `reduce_op` can be any of `nl.add`, `nl.subtract`, `nl.multiply`, `nl.max` or `nl.min`.

Reduction axis is not configurable in this API. If the input tile has multiple free axis, the API will reduce across all of them.

$$\begin{aligned} result &= data < op0 > operand0 \\ reduce\_res &= reduce\_op(dst, axis = < FreeAxis >) \end{aligned}$$

**Estimated instruction cost:**

$\max(\text{MIN\_II}, N) + \text{MIN\_II}$  Vector Engine cycles, where

- $N$  is the number of elements per partition in `data`, and
- $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

**Parameters**

- **data** – the input tile
- **op0** – the math operator used with `operand0` (any arithmetic operator in [Supported Math Operators for NKI ISA](#) is allowed)
- **operand0** – a scalar constant or a tile of shape `(data.shape[0], 1)`, where `data.shape[0]` is the partition axis size of the input data tile
- **reverse0** – (*not supported yet*) reverse ordering of inputs to `op0`; if false, `operand0` is the rhs of `op0`; if true, `operand0` is the lhs of `op0`. *<– currently not supported yet.*
- **reduce\_op** – the reduce operation to perform on the free dimension of `data <op0> operand0`
- **reduce\_res** – a tile of shape `(data.shape[0], 1)`, where `data.shape[0]` is the partition axis size of the input data tile. The result of `reduce_op(data <op0> operand0)` is written in-place into the tile.
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

an output tile of `(data <op0> operand0)` computation

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.isa.tensor\_copy

`nki.isa.tensor_copy(src, *, mask=None, dtype=None, engine=engine.unknown, **kwargs)`

Create a copy of `src` tile within NeuronCore on-chip SRAMs using Vector, Scalar or GpSimd Engine.

The output tile has the same partition axis size and also the same number of elements per partition as the input tile `src`.

All three compute engines, Vector, Scalar and GpSimd Engine can perform tensor copy. However, their copy behavior is slightly different across engines:

- Scalar Engine on NeuronCore-v2 performs copy by first casting the input tile to FP32 internally and then casting from FP32 to the output dtype (`dtype`, or `src.dtype` if `dtype` is not specified). Therefore, users should be cautious with assigning this instruction to Scalar Engine when the input data type cannot be precisely cast to FP32 (e.g., INT32).
- Both GpSimd and Vector Engine can operate in two modes: (1) bit-accurate copy when input and output data types are the same or (2) intermediate FP32 cast when input and output data types differ, similar to Scalar Engine.

In addition, since GpSimd Engine cannot access PSUM in NeuronCore, Scalar or Vector Engine must be chosen when the input or output tile is in PSUM (see [NeuronCore-v2 Compute Engines](#) for details). By default, this API returns a tile in SBUF, unless the returned value is assigned to a pre-declared PSUM tile.

### Estimated instruction cost:

$\max(\text{MIN\_II}, N)$  engine cycles, where  $N$  is the number of elements per partition in the input tile, and  $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

#### Parameters

- **src** – the source of copy, must be a tile in SBUF or PSUM.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **engine** – (optional) the engine to use for the operation: `nki.isa.vector_engine`, `nki.isa.scalar_engine`, `nki.isa.gpsimd_engine` or `nki.isa.unknown_engine` (default, compiler selects best engine based on engine workload).

#### Returns

a tile with the same content and partition axis size as the `src` tile.

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
# Example 1: Copy over the tensor to another tensor using the Vector engine.
#####
x = nl.load(in_tensor)
x_copy = nisa.tensor_copy(x, engine=nisa.vector_engine)
nl.store(out_tensor, value=x_copy)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

**nki.isa.tensor\_copy\_dynamic\_src**

`nki.isa.tensor_copy_dynamic_src(src, *, mask=None, dtype=None, engine=engine.unknown, **kwargs)`

Create a copy of `src` tile within NeuronCore on-chip SRAMs using Vector or Scalar or GpSimd Engine, with `src` located at a dynamic offset within each partition.

Both source and destination tiles can be in either SBUF or PSUM. By default, this API returns a tile in SBUF, unless the returned value is assigned to a pre-declared PSUM tile.

The source and destination tiles must also have the same number of partitions and the same number of elements per partition.

The dynamic offset must be a scalar value resided in SBUF. If you have a list of dynamic offsets for gathering tiles in SBUF/PSUM, you may loop over each offset and call `tensor_copy_dynamic_src` once per offset.

**Estimated instruction cost:**

`max(MIN_II_DYNAMIC, N)` engine cycles, where:

- `N` is the number of elements per partition in the `src` tile,
- `MIN_II_DYNAMIC` is the minimum instruction initiation interval for instructions with dynamic source location. `MIN_II_DYNAMIC` is roughly 600 engine cycles.

**Parameters**

- **src** – the source of copy, must be a tile in SBUF or PSUM that is dynamically indexed within each partition.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **engine** – (optional) the engine to use for the operation: `nki.isa.vector_engine`, `nki.isa.gpsimd_engine`, `nki.isa.scalar_engine` or `nki.isa.unknown_engine` (default, let compiler select best engine).
- **return** – the modified destination of copy.

Example:

```
import neuronxcc.nki.typing as nt
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
↪#####
# TensorCopyDynamicSrc example 0:
# - src_tensor in HBM of shape [128, 512]
# - offsets in HBM of shape [1, 64] (with values [4, 5, 6, 7, ...])
# - Gather tiles of shape [128, 1] from src_tensor into out_tensor using offsets
#####
↪#####

# Load src_tensor and offsets into SBUF
src_tensor_sbuf: nt.tensor[128, 512] = nl.load(src_tensor)
offsets_sbuf: nt.tensor[1, 64] = nl.load(offsets)
```

(continues on next page)

(continued from previous page)

```

# Copy into output tensor in SBUF
out_sbuf: nt.tensor[128, 64] = nl.ndarray([128, 64], dtype=src_tensor.dtype,
                                          buffer=nl.sbuf)

# Static indices to access a tile of shape [128, 1];
# Add dynamic offsets to iy for tensor_copy_dynamic_src
ix, iy = nl.mgrid[0:128, 0:1]

for idx in nl.affine_range(offsets_sbuf.shape[1]):
    out_sbuf[ix, idx] = nisa.tensor_copy_dynamic_src(
        src_tensor_sbuf[ix, offsets_sbuf[0, idx] + iy])

nl.store(out_tensor, value=out_sbuf)
...

import neuronxcc.nki.typing as nt
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

#####
→#####
# TensorCopyDynamicSrc example 1:
# - src_tensor in HBM of shape [128, 512, 4]
# - offsets in HBM of shape [1 x 8] (with values [4, 5, 6, 7, ...]) to index into
#   second axis of src_tensor
# - Gather tiles of shape [128, 4] from src_tensor into out_tensor using offsets
#####
→#####

# Load src_tensor and offsets into SBUF
src_tensor_sbuf: nt.tensor[128, 512, 4] = nl.load(src_tensor)
offsets_sbuf: nt.tensor[1, 8] = nl.load(offsets)

# Copy into output tensor in SBUF
out_sbuf: nt.tensor[128, 8, 4] = nl.ndarray([128, 8, 4], dtype=src_tensor.dtype,
                                          buffer=nl.sbuf)

# Static indices to access a tile of shape [128, 1, 4];
# Use dynamic offsets directly to index the second axis for tensor_copy_dynamic_src
ix, _, iz = nl.mgrid[0:128, 0:1, 0:4]

for idx in nl.affine_range(offsets.shape[1]):
    out_sbuf[ix, idx, iz] = nisa.tensor_copy_dynamic_src(
        src_tensor_sbuf[ix, offsets_sbuf[0, idx], iz])

nl.store(out_tensor, value=out_sbuf)
...

```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### nki.isa.tensor\_copy\_dynamic\_dst

`nki.isa.tensor_copy_dynamic_dst(*, dst, src, mask=None, dtype=None, engine=engine.unknown, **kwargs)`

Create a copy of `src` tile within NeuronCore on-chip SRAMs using Vector or Scalar or GpSimd Engine, with `dst` located at a dynamic offset within each partition.

Both source and destination tiles can be in either SBUF or PSUM.

The source and destination tiles must also have the same number of partitions and the same number of elements per partition.

The dynamic offset must be a scalar value resided in SBUF. If you have a list of dynamic offsets for scattering tiles in SBUF/PSUM, you may loop over each offset and call `tensor_copy_dynamic_dst` once per offset.

#### Estimated instruction cost:

`max(MIN_II_DYNAMIC, N)` engine cycles, where:

- `N` is the number of elements per partition in the `src` tile,
- `MIN_II_DYNAMIC` is the minimum instruction initiation interval for instructions with dynamic destination location. `MIN_II_DYNAMIC` is roughly 600 engine cycles.

#### Parameters

- **dst** – the destination of copy, must be a tile in SBUF or PSUM that is dynamically indexed within each dimension.
- **src** – the source of copy, must be a tile in SBUF or PSUM.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **engine** – (optional) the engine to use for the operation: `nki.isa.vector_engine`, `nki.isa.gpsimd_engine`, `nki.isa.scalar_engine` or `nki.isa.unknown_engine` (default, let compiler select best engine).

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### nki.isa.tensor\_copy\_predicated

`nki.isa.tensor_copy_predicated(*, src, dst, predicate, mask=None, dtype=None, reverse_pred=False, **kwargs)`

Conditionally copy elements from the `src` tile to the destination tile on SBUF / PSUM based on a predicate using Vector Engine.

This instruction provides low-level control over conditional data movement on NeuronCores, optimized for scenarios where only selective copying of elements is needed. Either `src` or `predicate` may be in PSUM, but not both simultaneously. Both `src` and `predicate` are permitted to be in SBUF.

Shape and data type constraints:

1. `src` (if it is a tensor), `dst`, and `predicate` must occupy the same number of partitions and same number of elements per partition.
2. `predicate` must be of type `uint8`, `uint16`, or `uint32`.
3. `src` and `dst` must share the same data type.

#### Behavior:

- Where `predicate` is `True`: The corresponding elements from `src` are copied to `dst` tile. If `src` is a scalar, the scalar is copied to the `dst` tile.



- Where predicate is False: The corresponding values in *dst* tile are unmodified

#### Estimated instruction cost:

| Cost (Vector Engine Cycles) | Condition   |
|-----------------------------|---|
| $\max(\text{MIN\_II}, N)$   | If <i>src</i> is from SBUF and predicate is from PSUM or the other way around |
| $\max(\text{MIN\_II}, 2N)$  | If both <i>src</i> and <i>dst</i> are in SBUF                                 |

- *N* is the number of elements per partition in *src* tile
- *MIN\_II* is the minimum instruction initiation interval for small input tiles. *MIN\_II* is roughly 64 engine cycles.

#### Parameters

- **src** – The source tile or number to copy elements from when predicate is True
- **dst** – The destination tile to copy elements to
- **predicate** – A tile that determines which elements to copy
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **reverse\_pred** – A boolean that reverses the effect of predicate.

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Conditionally copies elements from the `on_true` tile to
# SBUF/PSUM destination tile using Vector Engine, where copying occurs
# only at positions where the predicate evaluates to True.
#####

...
pre_tile: tensor[128, 512] = nl.load(predicate)
src_tile: tensor[128, 512] = nl.load(on_true_tensor)

ix, iy = nl.mgrid[0:128, 0:512]
dst_tile: tensor[128, 512] = nl.zeros(shape=src_tile.shape, dtype=src_tile.dtype)
dst_tile[ix, iy] = nl.load(on_false_tensor)

nisa.tensor_copy_predicated(src=src_tile, dst=dst_tile, predicate=pre_tile)
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.isa.reciprocal

`nki.isa.reciprocal(data, *, dtype=None, mask=None, **kwargs)`

Compute reciprocal of each element in the input data tile using Vector Engine.

### Estimated instruction cost:

$\max(\text{MIN\_II}, 8 \cdot N)$  Vector Engine cycles, where  $N$  is the number of elements per partition in `data`, and `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

### Parameters

- **data** – the input tile
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

### Returns

an output tile of reciprocal computation

Example:

```
import neuronxcc.nki as nki
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
...

x = nl.load(in_tensor[nl.mgrid[0:128, 0:512]])

y = nisa.reciprocal(x)
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.isa.iota

`nki.isa.iota(expr, dtype, *, mask=None, **kwargs)`

Build a constant literal in SBUF using GpSimd Engine, rather than transferring the constant literal values from the host to device.

The `iota` instruction takes an affine expression of `nki.language.arange()` indices as the input pattern to generate constant index values (see examples below for more explanation). The index values are computed in 32-bit integer math. The GpSimd Engine is capable of casting the integer results into any desirable data type (specified by `dtype`) before writing them back to SBUF, at no additional performance cost.

### Estimated instruction cost:

$150 + N$  GpSimd Engine cycles, where  $N$  is the number of elements per partition in the output tile.

### Parameters

- **expr** – an input affine expression of `nki.language.arange()`
- **dtype** – output data type of the generated constant literal (see [Supported Data Types](#) for more information)
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

**Returns**

an output tile in SBUF

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Generate tile a of 512 constant values in SBUF partition 0
# that start at 0 and increment by 1:
#####
# a = [0, 1, ..., 511]
expr_a = nl.arange(0, 512)[None, :]
a: tensor[1, 512] = nisa.iota(expr_a, dtype=nl.int32)

#####
# Example 2: Generate tile b of 128 constant values across SBUF partitions
# that start at 0 and increment by 1, with one value per partition:
# b = [[0],
#      [1],
#      ...,
#      [127]]
#####
expr_b = nl.arange(0, 128)[:, None]
b: tensor[128, 1] = nisa.iota(expr_b, dtype=nl.int32)

#####
# Example 3: Generate tile c of 512 constant values in SBUF partition 0
# that start at 0 and decrement by 1:
# c = [0, -1, ..., -511]
#####
expr_c = expr_a * -1
c: tensor[1, 512] = nisa.iota(expr_c, dtype=nl.int32)

#####
# Example 4: Generate tile d of 128 constant values across SBUF
# partitions that start at 5 and increment by 2
#####
# d = [[5],
#      [7],
#      ...,
#      [259]]
expr_d = 5 + expr_b * 2
d: tensor[128, 1] = nisa.iota(expr_d, dtype=nl.int32)

#####
# Example 5: Generate tile e of shape [128, 512] by
# broadcast-add expr_a and expr_b
# e = [[0, 1, ..., 511],
#      [1, 2, ..., 512],
#      ...
#      [127, 2, ..., 638]]
```

(continues on next page)

(continued from previous page)

```
#####
e: tensor[128, 512] = nisa.iota(expr_a + expr_b, dtype=nl.int32)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.dropout

`nki.isa.dropout(data, prob, *, mask=None, dtype=None, **kwargs)`

Randomly replace some elements of the input tile `data` with zeros based on input probabilities using Vector Engine. The probability of replacing input elements with zeros (i.e., drop probability) is specified using the `prob` field: - If the probability is 1.0, all elements are replaced with zeros. - If the probability is 0.0, all elements are kept with their original values.

The `prob` field can be a scalar constant or a tile of shape `(data.shape[0], 1)`, where each partition contains one drop probability value. The drop probability value in each partition is applicable to the input data elements from the same partition only.

Data type of the input data tile can be any valid NKI data types (see [Supported Data Types](#) for more information). However, data type of `prob` has restrictions based on the data type of `data`:

- If data type of `data` is any of the integer types (e.g., `int32`, `int16`), `prob` data type must be `float32`
- If data type of `data` is any of the float types (e.g., `float32`, `bfloat16`), `prob` data can be any valid float type

The output data type of this instruction is specified by the `dtype` field. The output data type must match the input data type of `data` if input data type is any of the integer types. Otherwise, output data type can be any valid NKI data types. If output data type is not specified, it is default to be the same as input data type.

### Estimated instruction cost:

`max(MIN_II, N)` Vector Engine cycles, where `N` is the number of elements per partition in `data`, and `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

#### Parameters

- **data** – the input tile
- **prob** – a scalar or a tile of shape `(data.shape[0], 1)` to indicate the probability of replacing elements with zeros
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

#### Returns

an output tile of the dropout result

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: From an input tile a of shape [128, 512], dropout its values
# with probabilities in tile b of shape [128, 1] and store the result in c.
#####
a: tensor[128, 512] = nl.load(a_tensor)
```

(continues on next page)

(continued from previous page)

```

b: tensor[128, 1] = nl.load(b_tensor)

c: tensor[128, 512] = nisa.dropout(a, prob=b)

nl.store(c_tensor, c)

#####
# Example 2: From an input tile a, dropout its values
# with probability of 0.2 and store the result in b.
#####
a = nl.load(in_tensor)

b = nisa.dropout(a, prob=0.2)

nl.store(out_tensor, b)

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.affine\_select

`nki.isa.affine_select(pred, on_true_tile, on_false_value, *, mask=None, dtype=None, **kwargs)`

Select elements between an input tile `on_true_tile` and a scalar value `on_false_value` according to a boolean predicate tile using GpSimd Engine. The predicate tile is calculated on-the-fly in the engine by evaluating an affine expression element-by-element as indicated in `pred`.

`pred` must meet the following requirements:

- It must not depend on any runtime variables that can't be resolved at compile-time.
- It can't be multiple masks combined using logical operators such as `&` and `|`.

For a complex predicate that doesn't meet the above requirements, consider using [nl.where](#).

The input tile `on_true_tile`, the calculated boolean predicate tile expressed by `pred`, and the returned output tile of this instruction must have the same shape. If the predicate value of a given position is `True`, the corresponding output element will take the element from `on_true_tile` in the same position. If the predicate value of a given position is `False`, the corresponding output element will take the value of `on_false_value`.

A common use case for `affine_select` is to apply a causal mask on the attention scores for transformer decoder models.

This instruction allows any float or 8-bit/16-bit integer data types for both the input data tile and output tile (see [Supported Data Types](#) for more information). The output tile data type is specified using the `dtype` field. If `dtype` is not specified, the output data type will be the same as the input data type of `data`. However, the data type of `on_false_value` must be float32, regardless of the input/output tile data types.

### Estimated instruction cost:

`GPSIMD_START` + `N` GpSimd Engine cycles, where `N` is the number of elements per partition in `on_true_tile` and `GPSIMD_START` is the instruction startup overhead on GpSimdE, roughly 150 engine cycles.

#### Parameters

- **pred** – an affine expression that defines the boolean predicate
- **on\_true\_tile** – an input tile for selection with a `True` predicate value
- **on\_false\_value** – a scalar value for selection with a `False` predicate value
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tiles, or whichever input type has the highest precision (see [NKI Type Promotion](#) for more information);

#### Returns

an output tile with values selected from either `on_true_tile` or `on_false_value` according to the following equation:  $\text{output}[x] = (\text{pred}[x] > 0) ? \text{on\_true\_tile}[x] : \text{on\_false\_value}$

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl

#####
# Example 1: Take tile a of shape [128, 128] and replace its
# upper triangle with -9984.0;
#####
ix, iy = nl.mgrid[0:128, 0:128]
a = nl.load(a_tensor[ix, iy])

b = nisa.affine_select(pred=(iy < ix), on_true_tile=a[ix, iy], on_false_value=-9984.
↪0)

nl.store(b_tensor[ix, iy], b)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

### nki.isa.range\_select

```
nki.isa.range_select(*, on_true_tile, comp_op0, comp_op1, bound0, bound1, reduce_cmd=reduce_cmd.idle,
                    reduce_res=None, reduce_op=<function amax>, range_start=0,
                    on_false_value=<property object>, mask=None, dtype=None, **kwargs)
```

Select elements from `on_true_tile` based on comparison with bounds using Vector Engine.

For each element in `on_true_tile`, compares its free dimension index + `range_start` against `bound0` and `bound1` using the specified comparison operators (`comp_op0` and `comp_op1`). If both comparisons evaluate to True, copies the element to the output; otherwise uses `on_false_value`.

Additionally performs a reduction operation specified by `reduce_op` on the results, storing the reduction result in `reduce_res`.

#### Note on numerical stability:

In self-attention, we often have this instruction sequence: `range_select` (VectorE) -> `reduce_res` -> `activation` (ScalarE). When `range_select` outputs a full row of `fill_value`, caution is needed to avoid NaN in the activation instruction that subtracts the output of `range_select` by `reduce_res` (max value):

- If `dtype` and `reduce_res` are both FP32, we should not hit any NaN issue since `FP32_MIN - FP32_MIN = 0`. Exponentiation on 0 is stable (1.0 exactly).
- If `dtype` is FP16/BF16/FP8, the `fill_value` in the output tile will become `-INF` since HW performs a downcast from `FP32_MIN` to a smaller dtype. In this case, you must make sure `reduce_res` uses FP32 dtype to avoid NaN in activation. NaN can be avoided because activation always upcasts input tiles to FP32 to perform math operations: `-INF - FP32_MIN = -INF`. Exponentiation on `-INF` is stable (0.0 exactly).

#### Constraints:

The comparison operators must be one of:

- `np.equal`
- `np.less`
- `np.less_equal`
- `np.greater`
- `np.greater_equal`

Partition dim sizes must match across `on_true_tile`, `bound0`, and `bound1`:

- `bound0` and `bound1` must have one element per partition
- `on_true_tile` must be one of the FP dtypes, and `bound0/bound1` must be FP32 types.

The comparison with `bound0`, `bound1`, and free dimension index is done in FP32. Make sure `range_start` + free dimension index is within  $2^{24}$  range.

#### Estimated instruction cost:

`max(MIN_II, N)` Vector Engine cycles, where:

- `N` is the number of elements per partition in `on_true_tile`, and
- `MIN_II` is the minimum instruction initiation interval for small input tiles.
- `MIN_II` is roughly 64 engine cycles.

#### Numpy equivalent:

```
indices = np.zeros(on_true_tile.shape)
indices[:] = range_start + np.arange(on_true_tile[0].size)

mask = comp_op0(indices, bound0) & comp_op1(indices, bound1)
select_out_tile = np.where(mask, on_true_tile, on_false_value)
reduce_tile = reduce_op(select_out_tile, axis=1, keepdims=True)
```

#### Parameters

- **`on_true_tile`** – input tile containing elements to select from
- **`on_false_value`** – constant value to use when selection condition is False. Due to HW constraints, this must be FP32\_MIN FP32 bit pattern
- **`comp_op0`** – first comparison operator
- **`comp_op1`** – second comparison operator
- **`bound0`** – tile with one element per partition for first comparison
- **`bound1`** – tile with one element per partition for second comparison
- **`reduce_op`** – reduction operator to apply on across the selected output. Currently only `np.max` is supported.
- **`reduce_res`** – optional tile to store reduction results.
- **`range_start`** – starting base offset for index array for the free dimension of `on_true_tile` Defaults to 0, and must be a compiler time integer.
- **`mask`** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **`dtype`** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

#### Returns

output tile with selected elements

Example:

```
import neuronxcc.nki as nki
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import numpy as np
...
```

(continues on next page)

(continued from previous page)

```
#####
# Example 1: # Select elements where
# bound0 <= range_start + index < bound1 and compute max reduction
#
# on_false_value must be nl.fp32.min
#####
on_true_tile = nl.load(on_true[...])
bound0_tile = nl.load(bound0[...])
bound1_tile = nl.load(bound1[...])

reduce_res_tile = nl.ndarray((on_true.shape[0], 1), dtype=nl.float32, buffer=nl.
↪sbuf)
result = nl.ndarray(on_true.shape, dtype=nl.float32, buffer=nl.sbuf)

result[...] = nisa.range_select(
    on_true_tile=on_true_tile,
    comp_op0=compare_op0,
    comp_op1=compare_op1,
    bound0=bound0_tile,
    bound1=bound1_tile,
    reduce_cmd=nisa.reduce_cmd.reset_reduce,
    reduce_res=reduce_res_tile,
    reduce_op=np.max,
    range_start=range_start,
    on_false_value=nl.fp32.min
)

nl.store(select_res[...], value=result[...])
nl.store(reduce_result[...], value=reduce_res_tile[...])
```

Alternatively, `reduce_cmd` can be used to chain multiple calls to the same accumulation register to accumulate across multiple `range_select` calls. For example:

```
import neuronxcc.nki as nki
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import numpy as np
...

#####
# Example 2.a: Initialize reduction with first range_select
# Notice we don't pass reduce_res since the accumulation
# register keeps track of the accumulation until we're ready to
# read it. Also we use reset_reduce in order to "clobber" or zero
# out the accumulation register before we start accumulating.
#
# Note: Since the type of these tensors are fp32, we use nl.fp32.min
# for on_false_value due to HW constraints.
#####
on_true_tile = nl.load(on_true[...])
bound0_tile = nl.load(bound0[...])
```

(continues on next page)



(continued from previous page)

```

bound1_tile = nl.load(bound1[...])

reduce_res_sbuf = nl.ndarray((on_true.shape[0], 1), dtype=np.float32, buffer=nl.
    ↪sbuf)
result_sbuf = nl.ndarray(on_true.shape, dtype=np.float32, buffer=nl.sbuf)

result_sbuf[...] = nisa.range_select(
    on_true_tile=on_true_tile,
    comp_op0=compare_op0,
    comp_op1=compare_op1,
    bound0=bound0_tile,
    bound1=bound1_tile,
    reduce_cmd=nisa.reduce_cmd.reset_reduce,
    reduce_op=np.max,
    range_start=range_start,
    on_false_value=nl.fp32.min
)

#####
# Example 2.b: Chain multiple range_select operations
# with reduction in an affine loop. Adding ones just lets us ensure the reduction
# gets updated with new values.
#####
ones = nl.full(on_true.shape, fill_value=1, dtype=np.float32, buffer=nl.sbuf)
# we are going to loop as if we're tiling on the partition dimension
iteration_step_size = on_true_tile.shape[0]

# Perform chained operations using an affine loop index for range_start
for i in range(1, 2):
    # Update input values
    on_true_tile[...] = nl.add(on_true_tile, ones)

    # Continue reduction with updated values
    # notice, we still don't have reduce_res specified
    result_sbuf[...] = nisa.range_select(
        on_true_tile=on_true_tile,
        comp_op0=compare_op0,
        comp_op1=compare_op1,
        bound0=bound0_tile,
        bound1=bound1_tile,
        reduce_cmd=nisa.reduce_cmd.reduce,
        reduce_op=np.max,
        # we can also use index expressions for setting the start of the range
        range_start=range_start + (i * iteration_step_size),
        on_false_value=nl.fp32.min
    )

range_start = range_start + (2 * iteration_step_size)
#####
# Example 2.c: Final iteration, we actually want the results to
# return to the user so we pass reduce_res argument so the
# reduction will be written from the accumulation

```

(continues on next page)

(continued from previous page)

```
# register to reduce_res_tile
#####
on_true_tile[...] = nl.add(on_true_tile, ones)
result_sbuf[...] = nisa.range_select(
    on_true_tile=on_true_tile,
    comp_op0=compare_op0,
    comp_op1=compare_op1,
    bound0=bound0_tile,
    bound1=bound1_tile,
    reduce_cmd=nisa.reduce_cmd.reduce,
    reduce_res=reduce_res_sbuf [...],
    reduce_op=np.max,
    range_start=range_start,
    on_false_value=nl.fp32.min
)

nl.store(select_res [...], value=result_sbuf [...])
nl.store(reduce_result [...], value=reduce_res_sbuf [...])
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.memset

`nki.isa.memset(shape, value, dtype, *, mask=None, engine=engine.unknown, **kwargs)`

Initialize a tile filled with a compile-time constant value using Vector or GpSimd Engine. The shape of the tile is specified in the `shape` field and the initialized value in the `value` field. The `memset` instruction supports all valid NKI dtypes (see [Supported Data Types](#)).

### Parameters

- **shape** – the shape of the output tile; layout: (partition axis, free axis). Note that `memset` ignores `nl.par_dim()` and always treats the first dimension as the partition dimension.
- **value** – the constant value to initialize with
- **dtype** – data type of the output tile (see [Supported Data Types](#) for more information)
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **engine** – specify which engine to use for `memset`: `nki.isa.vector_engine` or `nki.isa.gpsimd_engine`; `nki.isa.unknown_engine` by default, lets compiler select the best engine for the given input tile shape

### Returns

a tile with shape `shape` whose elements are initialized to `value`.

### Estimated instruction cost:

Given `N` is the number of elements per partition in the output tile, and `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

- If the initialized value is zero and output data type is `bfloat16/float16`, `max(MIN_II, N/2)` Vector Engine cycles;
- Otherwise, `max(MIN_II, N)` Vector Engine cycles

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
```

(continues on next page)

(continued from previous page)

```

...

#####
# Example 1: Initialize a float32 tile a of shape (128, 128)
# with a value of 0.2
#####
a = nisa.memset(shape=(128, 128), value=0.2, dtype=nl.float32)

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.bn\_stats

`nki.isa.bn_stats(data, *, mask=None, dtype=None, **kwargs)`

Compute mean- and variance-related statistics for each partition of an input tile `data` in parallel using Vector Engine.

The output tile of the instruction has 6 elements per partition:

- the count of the even elements (of the input tile elements from the same partition)
- the mean of the even elements
- `variance * count` of the even elements
- the count of the odd elements
- the mean of the odd elements
- `variance * count` of the odd elements

To get the final mean and variance of the input tile, we need to pass the above `bn_stats` instruction output into the `bn_aggr` instruction, which will output two elements per partition:

- mean (of the original input tile elements from the same partition)
- variance

Due to hardware limitation, the number of elements per partition (i.e., free dimension size) of the input `data` must not exceed 512 (`nl.tile_size.bn_stats_fmax`). To calculate per-partition mean/variance of a tensor with more than 512 elements in free dimension, we can invoke `bn_stats` instructions on each 512-element tile and use a single `bn_aggr` instruction to aggregate `bn_stats` outputs from all the tiles. Refer to Example 2 for an example implementation.

Vector Engine performs the above statistics calculation in float32 precision. Therefore, the engine automatically casts the input `data` tile to float32 before performing float32 computation and is capable of casting the float32 computation results into another data type specified by the `dtype` field, at no additional performance cost. If `dtype` field is not specified, the instruction will cast the float32 results back to the same data type as the input `data` tile.

### Estimated instruction cost:

$\max(\text{MIN\_II}, N)$  Vector Engine cycles, where  $N$  is the number of elements per partition in `data` and `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

#### Parameters

- **data** – the input tile (up to 512 elements per partition)
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

#### Returns

an output tile with 6-element statistics per partition

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Calculate the mean and variance for each partition
# of tile a with shape (128, 128)
#####
a: tensor[128, 128] = nl.load(a_tensor)
stats_a: tensor[128, 6] = nisa.bn_stats(a)
mean_var_a: tensor[128, 2] = nisa.bn_aggr(stats_a)

# Extract mean and variance
mean_a = mean_var_a[:, 0]
var_a = mean_var_a[:, 1]
nl.store(mean_a_tensor, mean_a)
nl.store(var_a_tensor, var_a)

# #####
# # Example 2: Calculate the mean and variance for each partition of
# # tile b with shape [128, 1024]
# #####
b: tensor[128, 1024] = nl.load(b_tensor)

# Run bn_stats in two tiles because b has 1024 elements per partition,
# but bn_stats has a limitation of nl.tile_size.bn_stats_fmax
# Initialize a bn_stats output tile with shape of [128, 6*2] to
# hold outputs of two bn_stats instructions
stats_b = nl.ndarray((128, 6 * 2), dtype=nl.float32)
bn_tile = nl.tile_size.bn_stats_fmax
ix, iy = nl.mgrid[0:128, 0:bn_tile]
iz, iw = nl.mgrid[0:128, 0:6]

for i in range(1024 // bn_tile):
    stats_b[iz, i * 6 + iw] = nisa.bn_stats(b[ix, i * bn_tile + iy], dtype=nl.float32)

mean_var_b = nisa.bn_aggr(stats_b)

# Extract mean and variance
mean_b = mean_var_b[:, 0]
var_b = mean_var_b[:, 1]

nl.store(mean_b_tensor, mean_b)
nl.store(var_b_tensor, var_b)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.isa.bn\_aggr

`nki.isa.bn_aggr(data, *, mask=None, dtype=None, **kwargs)`

Aggregate one or multiple `bn_stats` outputs to generate a mean and variance per partition using Vector Engine.

The input data tile effectively has an array of (count, mean, variance\*count) tuples per partition produced by `bn_stats` instructions. Therefore, the number of elements per partition of data must be a modulo of three.

Note, if you need to aggregate multiple `bn_stats` instruction outputs, it is recommended to declare a SBUF tensor and then make each `bn_stats` instruction write its output into the SBUF tensor at different offsets (see example implementation in Example 2 in `bn_stats`).

Vector Engine performs the statistics aggregation in float32 precision. Therefore, the engine automatically casts the input data tile to float32 before performing float32 computation and is capable of casting the float32 computation results into another data type specified by the `dtype` field, at no additional performance cost. If `dtype` field is not specified, the instruction will cast the float32 results back to the same data type as the input data tile.

### Estimated instruction cost:

$\max(\text{MIN\_II}, 13 * (N/3))$  Vector Engine cycles, where  $N$  is the number of elements per partition in data and  $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

#### Parameters

- **data** – an input tile with results of one or more `bn_stats`
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

#### Returns

an output tile with two elements per partition: a mean followed by a variance

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.isa.local\_gather

`nki.isa.local_gather(src_buffer, index, num_elem_per_idx=1, num_valid_indices=None, *, mask=None)`

Gather SBUF data in `src_buffer` using `index` on GpSimd Engine.

Each of the eight GpSimd cores in GpSimd Engine connects to 16 contiguous SBUF partitions (e.g., core[0] connected to partition[0:16]) and performs gather from the connected 16 SBUF partitions *independently* in parallel. The indices used for gather on each core should also come from the same 16 connected SBUF partitions.

During execution of the instruction, each GpSimd core reads a 16-partition slice from `index`, flattens all indices into a 1D array `indices_1d` (along the partition dimension first). By default with no `num_valid_indices` specified, each GpSimd core will treat all indices from its corresponding 16-partition `index` slice as valid indices. However, when the number of valid indices per core is not a multiple of 16, users can explicitly specify the valid index count per core in `num_valid_indices`. Note, `num_valid_indices` must not exceed the total element count in each 16-partition `index` slice (i.e., `num_valid_indices <= index.size / (index.shape[0] / 16)`).

Next, each GpSimd core uses the flattened `indices_1d` indices as *partition offsets* to gather from the connected 16-partition slice of `src_buffer`. Optionally, this API also allows gathering of multiple contiguous elements starting at each index to improve gather throughput, as indicated by `num_elem_per_idx`. Behavior of out-of-bound index access is undefined.

Even though all eight GpSimd cores can gather with completely different indices, a common use case for this API is to make all cores gather with the same set of indices (i.e., partition offsets). In this case, users can generate indices into 16 partitions, replicate them eight times to 128 partitions and then feed them into `local_gather`.

As an example, if `src_buffer` is (128, 512) in shape and `index` is (128, 4) in shape, where the partition dimension size is 128, `local_gather` effectively performs the following operation:

```
num_gpsimd_cores = 8
num_partitions_per_core = 16

src_buffer = np.random.random_sample([128, 512, 4]).astype(np.float32) * 100
index_per_core = np.random.randint(low=0, high=512, size=(16, 4), dtype=np.uint16)
# replicate 8 times for 8 GpSimd cores
index = np.tile(index_per_core, (num_gpsimd_cores, 1))
num_elem_per_idx = 4
index_hw = index * num_elem_per_idx
num_valid_indices = 64
output_shape = (128, 4, 16, 4)

num_active_cores = index.shape[0] / num_partitions_per_core
num_valid_indices = num_valid_indices if num_valid_indices \
    else index.size / num_active_cores

output_np = np.ndarray(shape=(128, num_valid_indices, num_elem_per_idx),
                        dtype=src_buffer.dtype)

for i_core in range(num_gpsimd_cores):
    start_par = i_core * num_partitions_per_core
    end_par = (i_core + 1) * num_partitions_per_core
    indices_1d = index[start_par:end_par].flatten(order='F')[0: num_valid_indices]

    output_np[start_par:end_par, :, :] = np.take(
        src_buffer[start_par:end_par],
        indices_1d, axis=1)

output_np = output_np.reshape(output_shape)
```

`local_gather` preserves the input data types from `src_buffer` in the gather output. Therefore, no data type casting is allowed in this API. The indices in `index` tile must be uint16 types.

This API has three tile size constraints [subject to future relaxation]:

1. The partition axis size of `src_buffer` must match that of `index` and must be a multiple of 16. In other words, `src_buffer.shape[0] == index.shape[0]` and `src_buffer.shape[0] % 16 == 0`.
2. The number of contiguous elements to gather per index per partition `num_elem_per_idx` must be one of the following values: [1, 2, 4, 8, 16, 32].
3. The number of indices for gather per core must be less than or equal to 4096.

#### Estimated instruction cost:

$150 + (\text{num\_valid\_indices} * \text{num\_elem\_per\_idx}) / C$  GpSimd Engine cycles, where C can be calculated using  $((28 + t * \text{num\_elem\_per\_idx}) / (t * \text{num\_elem\_per\_idx})) / \min(4 / \text{dtype\_size}, \text{num\_elem\_per\_idx})$ . `dtype_size` is the size of `src_buffer.dtype` in bytes. Currently, `t` is a constant 4, but subject to change in future software implementation.

#### Parameters

- **src\_buffer** – an input tile for gathering.
- **index** – an input tile with indices used for gathering.

- **num\_elem\_per\_idx** – an optional integer value to read multiple contiguous elements per index per partition; default is 1.
- **num\_valid\_indices** – an optional integer value to specify the number of valid indices per GpSimd core; default is `index.size / (index.shape[0] / 16)`.
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Returns

an output tile of the gathered data

#### Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: gather src_buffer using index
# Gather input: src_buffer_tile with shape (128, 512, 4)
# Gather indices: index_tile with shape (128, 4)
# We use num_valid_indices indices per core, and read num_elem_per_idx
# contiguous elements per partition.
#####
src_buffer_tile: tensor[128, 512, 4] = nl.load(src_buffer)
index_tile: tensor[128, 4] = nl.load(index)
output_tile: tensor[128, 4, 16, 4] = nisa.local_gather(
    src_buffer_tile, index_tile, num_elem_per_idx, num_valid_indices)

nl.store(output, output_tile)
```

Click [here](#) to download the full NKI code example with equivalent numpy implementation.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### nki.isa.dma\_copy

```
nki.isa.dma_copy(*, dst, src, mask=None, dst_rmw_op=None, oob_mode=oob_mode.error,
                 dge_mode=dge_mode.unknown)
```

Copy data from `src` to `dst` using DMA engine. Both `src` and `dst` tiles can be in device memory (HBM) or SBUF. However, if both `src` and `dst` tiles are in SBUF, consider using [nisa.tensor\\_copy](#) instead for better performance.

#### Parameters

- **src** – the source of copy.
- **dst** – the dst of copy.
- **dst\_rmw\_op** – the read-modify-write operation to be performed at the destination. Currently only `np.add` is supported, which adds the source data to the existing destination data. If `None`, the source data directly overwrites the destination. If `dst_rmw_op` is specified, only `oob_mode=oob_mode.error` is allowed. For best performance with Descriptor Generation Engine (DGE), unique dynamic offsets must be used to access `dst`. Multiple accesses to the same offset will cause a data hazard. If duplicated offsets are present, the compiler automatically adds synchronization to avoid hazards, which slows down computation.

- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see *NKI API Masking* for details)
- **mode** – (optional) Specifies how to handle out-of-bounds (oob) array indices during indirect access operations. Valid modes are:
  - `oob_mode.error`: (Default) Raises an error when encountering out-of-bounds indices.
  - `oob_mode.skip`: Silently skips any operations involving out-of-bounds indices.
 For example, when using indirect gather/scatter operations, out-of-bounds indices can occur if the index array contains values that exceed the dimensions of the target array.
- **dge\_mode** – (optional) specify which Descriptor Generation Engine (DGE) mode to use for copy: `nki.isa.dge_mode.none` (turn off DGE) or `nki.isa.dge_mode.swdge` (software DGE) or `nki.isa.dge_mode.hwdge` (hardware DGE) or `nki.isa.dge_mode.unknown` (by default, let compiler select the best DGE mode). HWDGE is only supported for NeuronCore-v3+.

A cast will happen if the `src` and `dst` have different dtype.

Example:

```
import neuronxcc.nki.isa as nisa

#####
# Example 1: Copy over the tensor to another tensor
#####
nisa.dma_copy(dst=b, src=a)
```

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 2: Load elements from HBM with indirect addressing. If addressing
# results out-of-bound access, the operation will fail.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n//2, 0:m]

expr_arange = 2*nl.arange(n//2)[: , None]
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[64, 512] = nisa.memset(shape=(n//2, m), value=-1, dtype=in_tensor.
↳ dtype)
nisa.dma_copy(src=in_tensor[idx_tile, iy], dst=out_tile[ix, iy], oob_mode=nisa.oob_
↳ mode.error)
```

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 3: Load elements from HBM with indirect addressing. If addressing
```

(continues on next page)



(continued from previous page)

```

# results in out-of-bounds access, the operation will fail.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n//2, 0:m]

# indices are out of range on purpose to demonstrate the error
expr_arange = 3*nl.arange(n//2)[: , None]
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[64, 512] = nisa.memset(shape=(n//2, m), value=-1, dtype=in_tensor.
↪dtype)
nisa.dma_copy(src=in_tensor[idx_tile, iy], dst=out_tile[ix, iy], oob_mode=nisa.oob_
↪mode.error)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 4: Load elements from HBM with indirect addressing. If addressing
# results in out-of-bounds access, the operation will skip indices.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n//2, 0:m]

# indices are out of range on purpose
expr_arange = 3*nl.arange(n//2)[: , None]
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[64, 512] = nisa.memset(shape=(n//2, m), value=-1, dtype=in_tensor.
↪dtype)
nisa.dma_copy(src=in_tensor[idx_tile, iy], dst=out_tile[ix, iy], oob_mode=nisa.oob_
↪mode.skip)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 5: Store elements to HBM with indirect addressing and with
# read-modified-write operation.
#####

...
n, m = in_tensor.shape

```

(continues on next page)

(continued from previous page)

```

ix, iy = nl.mgrid[0:n, 0:m]

expr_arange = 2*nl.arange(n)[: , None]
inp_tile: tensor[64, 512] = nl.load(in_tensor[ix, iy])
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[128, 512] = nisa.memset(shape=(2*n, m), value=1, dtype=in_tensor.
    ↪dtype)
nl.store(out_tensor, value=out_tile)
nisa.dma_copy(dst=out_tensor[idx_tile, iy], src=inp_tile, dst_rmw_op=np.add)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 6: Store elements to HBM with indirect addressing. If indirect
# addressing results out-of-bound access, the operation will fail.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n, 0:m]

expr_arange = 2*nl.arange(n)[: , None]
inp_tile: tensor[64, 512] = nl.load(in_tensor[ix, iy])
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[128, 512] = nisa.memset(shape=(2*n, m), value=-1, dtype=in_tensor.
    ↪dtype)
nl.store(out_tensor, value=out_tile)
nisa.dma_copy(dst=out_tensor[idx_tile, iy], src=inp_tile, oob_mode=nisa.oob_mode.
    ↪error)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 7: Store elements to HBM with indirect addressing. If indirect
# addressing results out-of-bounds access, the operation will skip indices.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n, 0:m]

# indices are out of range on purpose to demonstrate the error
expr_arange = 3*nl.arange(n)[: , None]
inp_tile: tensor[64, 512] = nl.load(in_tensor[ix, iy])
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

```

(continues on next page)

(continued from previous page)

```

out_tile: tensor[128, 512] = nisa.memset(shape=(2*n, m), value=-1, dtype=in_tensor.
↳dtype)
nl.store(out_tensor, value=out_tile)
nisa.dma_copy(dst=out_tensor[idx_tile, iy], src=inp_tile, oob_mode=nisa.oob_mode.
↳error)

```

```

#####
# Example 8: Store elements to HBM with indirect addressing. If indirect
# addressing results out-of-bounds access, the operation will skip indices.
#####

...
n, m = in_tensor.shape
ix, iy = nl.mgrid[0:n, 0:m]

# indices are out of range on purpose
expr_arange = 3*nl.arange(n)[: , None]
inp_tile: tensor[64, 512] = nl.load(in_tensor[ix, iy])
idx_tile: tensor[64, 1] = nisa.iota(expr_arange, dtype=np.int32)

out_tile: tensor[128, 512] = nisa.memset(shape=(2*n, m), value=-1, dtype=in_tensor.
↳dtype)
nl.store(out_tensor, value=out_tile)
nisa.dma_copy(dst=out_tensor[idx_tile, iy], src=inp_tile, oob_mode=nisa.oob_mode.
↳skip)

```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.isa.max8

`nki.isa.max8(*, src, mask=None, dtype=None, **kwargs)`

Find the 8 largest values in each partition of the source tile.

This instruction reads the input elements, converts them to fp32 internally, and outputs the 8 largest values in descending order for each partition. By default, returns the same dtype as the input tensor.

The source tile can be up to 5-dimensional, while the output tile is always 2-dimensional. The number of elements read per partition must be between 8 and 16,384 inclusive. The output will always contain exactly 8 elements per partition. The source and output must have the same partition dimension size:

- source: [par\_dim, ...]
- output: [par\_dim, 8]

### Estimated instruction cost:

N engine cycles, where:

- N is the number of elements per partition in the source tile

### Parameters

- **src** – the source tile to find maximum values from

- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see *NKI API Masking* for details)
- **dtype** – (optional) data type to cast the output type to (see *Supported Data Types* for more information); if not specified, it will default to be the same as the data type of the input tile.

**Returns**

a 2D tile containing the 8 largest values per partition in descending order with shape [par\_dim, 8]

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Generate tile b of 32 * 128 random floating point values
# and get the 8 largest values in each row:
#####
expr_a = nl.rand((32, 128))
a = nisa.max8(src=expr_a)

a_tensor = nl.ndarray([32, 8], dtype=nl.float32, buffer=nl.shared_hbm)
nl.store(a_tensor, value=a)
```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

**nki.isa.nc\_find\_index8**

`nki.isa.nc_find_index8(*, data, vals, mask=None, dtype=None, **kwargs)`

Find indices of the 8 given vals in each partition of the data tensor.

This instruction first loads the 8 values, then loads the data tensor and outputs the indices (starting at 0) of the first occurrence of each value in the data tensor, for each partition.

The data tensor can be up to 5-dimensional, while the vals tensor must be up to 3-dimensional. The data tensor must have between 8 and 16,384 elements per partition. The vals tensor must have exactly 8 elements per partition. The output will contain exactly 8 elements per partition and will be uint16 or uint32 type. Default output type is uint32.

Behavior is undefined if vals tensor contains values that are not in the data tensor.

If provided, a mask is applied only to the data tensor.

**Estimated instruction cost:**

N engine cycles, where:

- N is the number of elements per partition in the data tensor

**Parameters**

- **data** – the data tensor to find indices from
- **vals** – tensor containing the 8 values per partition whose indices will be found
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see *NKI API Masking* for details)
- **dtype** – uint16 or uint32

**Returns**

a 2D tile containing indices (uint16 or uint32) of the 8 values in each partition with shape [par\_dim, 8]

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1: Generate tile b of 32 * 128 random floating point values,
# find the 8 largest values in each row, then find their indices:
#####
# Generate random data
data = nl.rand((32, 128))

# Find max 8 values per row
max_vals = nisa.max8(src=data)

# Create output tensor for indices
indices_tensor = nl.ndarray([32, 8], dtype=nl.uint32, buffer=nl.shared_hbm)

# Find indices of max values
indices = nisa.nc_find_index8(data=data, vals=max_vals)

# Store results
nl.store(indices_tensor, value=indices)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.nc\_match\_replace8**

`nki.isa.nc_match_replace8(*, data, vals, imm, dst_idx=None, mask=None, dtype=None, **kwargs)`

Replace first occurrence of each value in `vals` with `imm` in `data` using the Vector engine and return the replaced tensor. If `dst_idx` tile is provided, the indices of the matched values are written to `dst_idx`.

This instruction reads the input `data`, replaces the first occurrence of each of the given values (from `vals` tensor) with the specified immediate constant and, optionally, output indices of matched values to `dst_idx`. When performing the operation, the free dimensions of both `data` and `vals` are flattened. However, these dimensions are preserved in the replaced output tensor and in `dst_idx` respectively. The partition dimension defines the parallelization boundary. Match, replace, and index generation operations execute independently within each partition.

The `data` tensor can be up to 5-dimensional, while the `vals` tensor can be up to 3-dimensional. The `vals` tensor must have exactly 8 elements per partition. The `data` tensor must have no more than 16,384 elements per partition. The replaced output will have the same shape as the input `data` tensor. `data` and `vals` must have the same number of partitions. Both input tensors can come from SBUF or PSUM.

Behavior is undefined if `vals` tensor contains values that are not in the `data` tensor.

If provided, a mask is applied to the `data` tensor.

**Estimated instruction cost:**

`min(MIN_II, N)` engine cycles, where:

- `N` is the number of elements per partition in the data tensor
- `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

**NumPy equivalent:**

```
# Let's assume we work with NumPy, and ``data``, ``vals`` are 2-dimensional arrays
# (with shape[0] being the partition axis) and imm is a constant float32 value.

import numpy as np

# Get original shapes
data_shape = data.shape
vals_shape = vals.shape

# Reshape to 2D while preserving first dimension
data_2d = data.reshape(data_shape[0], -1)
vals_2d = vals.reshape(vals_shape[0], -1)

# Initialize output array for indices
indices = np.zeros(vals_2d.shape, dtype=np.uint32)

for i in range(data_2d.shape[0]):
    for j in range(vals_2d.shape[1]):
        val = vals_2d[i, j]
        # Find first occurrence of val in data_2d[i, :]
        matches = np.where(data_2d[i, :] == val)[0]
        if matches.size > 0:
            indices[i, j] = matches[0] # Take first match
            data_2d[i, matches[0]] = imm

output = data_2d.reshape(data.shape)
indices = indices.reshape(vals.shape) # Computed only if ``dst_idx`` is specified
```

#### Parameters

- **data** – the data tensor to modify
- **dst\_idx** – (optional) the destination tile to write flattened indices of matched values
- **vals** – tensor containing the 8 values per partition to replace
- **imm** – float32 constant to replace matched values with
- **mask** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)
- **dtype** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.

#### Returns

the modified data tensor

Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import neuronxcc.nki.typing as nt
```

(continues on next page)

(continued from previous page)

```
#####
# Example 1: Generate tile a of random floating point values,
# get the 8 largest values in each row, then replace their first
# occurrences with -inf:
#####
N = 4
M = 16
data_tile = nl.rand((N, M))
max_vals = nisa.max8(src=data_tile)

result = nisa.nc_match_replace8(data=data_tile[:, :], vals=max_vals, imm=float('-inf
→'))
result_tensor = nl.ndarray([N, M], dtype=nl.float32, buffer=nl.shared_hbm)
nl.store(result_tensor, value=result)
```

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import neuronxcc.nki.typing as nt

#####
# Example 2: Read the 8 largest values in each row of the tensor,
# replace the first occurrence with imm, write indices, and return
# the replaced output.
#####
n, m = in_tensor.shape

dst_idx = nl.ndarray((n, 8), dtype=idx_tensor.dtype)

ix, iy = nl.mgrid[0:n, 0:8]

inp_tile: nt.tensor[n, m] = nl.load(in_tensor)
max_vals: nt.tensor[n, 8] = nisa.max8(src=inp_tile)

out_tile = nisa.nc_match_replace8(
    dst_idx=dst_idx[ix, iy], data=inp_tile[:, :], vals=max_vals, imm=imm
)
```

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import neuronxcc.nki.typing as nt

#####
# Example 3: Read the 8 largest values in each row of the tensor,
# after applying the specified mask, replace the first occurrence
# with imm, write indices, and return the replaced output.
#####
n, m = in_tensor.shape

idx_tile = nisa.memset(shape=(n, 8), value=0, dtype=nl.uint32)
```

(continues on next page)

(continued from previous page)

```

ix, iy = nl.mgrid[0:n, 0:m]
inp_tile: nt.tensor[n, m] = nl.load(in_tensor)
max_vals: nt.tensor[n, 8] = nisa.max8(src=inp_tile[ix, iy], mask=(ix < n // 2 and iy
↪ < m // 2))

out_tile = nisa.nc_match_replace8(
    dst_idx=idx_tile[:, :],
    data=inp_tile[ix, iy],
    vals=max_vals,
    imm=imm,
    mask=(ix < n // 2 and iy < m // 2), # mask applies to `data`
)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
import neuronxcc.nki.typing as nt

#####
# Example 4: Read the 8 largest values in each row of the tensor,
# replace the first occurrence with 0.0, write indices, and return
# the replaced output.
#####
n, b, m = data_tensor.shape

n, b, m = data_tensor.shape

out_tensor = nl.ndarray([n, b, m], dtype=data_tensor.dtype, buffer=nl.hbm)
idx_tensor = nl.ndarray([n, 8], dtype=nl.uint32, buffer=nl.hbm)

imm = 0.0
idx_tile = nisa.memset(shape=(n, 8), value=0, dtype=nl.uint32)
out_tile = nisa.memset(shape=(n, b, m), value=0, dtype=data_tensor.dtype)

iq, ir, iw = nl.mgrid[0:n, 0:b, 0:m]
ip, io = nl.mgrid[0:n, 0:8]

inp_tile = nl.load(data_tensor[iq, ir, iw])
max_vals: nt.tensor[n, 8] = nisa.max8(src=inp_tile)

out_tile[iq, ir, iw] = nisa.nc_match_replace8(
    dst_idx=idx_tile[ip, io],
    data=inp_tile[iq, ir, iw],
    vals=max_vals[ip, io],
    imm=imm,
)

```

```

import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl

```

(continues on next page)



(continued from previous page)

```

import neuronxcc.nki.typing as nt

#####
# Example 5: Read the 8 largest values in each row of the tensor,
# replace the first occurrence with 0.0 in-place and write indices.
#####
n, b, m = data_tensor.shape

n, b, m = data_tensor.shape

out_tensor = nl.ndarray([n, b, m], dtype=data_tensor.dtype, buffer=nl.hbm)
idx_tensor = nl.ndarray([n, 8], dtype=nl.uint32, buffer=nl.hbm)

imm = 0.0
idx_tile = nisa.memset(shape=(n, 8), value=0, dtype=nl.uint32)

iq, ir, iw = nl.mgrid[0:n, 0:b, 0:m]
ip, io = nl.mgrid[0:n, 0:8]

inp_tile = nl.load(data_tensor[iq, ir, iw])
max_vals: nt.tensor[n, 8] = nisa.max8(src=inp_tile)

inp_tile[iq, ir, iw] = nisa.nc_match_replace8(
    dst_idx=idx_tile[ip, io],
    data=inp_tile[iq, ir, iw],
    vals=max_vals[ip, io],
    imm=imm,
)

```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### nki.isa.nc\_stream\_shuffle

`nki.isa.nc_stream_shuffle(src, dst, shuffle_mask, *, dtype=None, mask=None, **kwargs)`

Apply cross-partition data movement within a quadrant of 32 partitions from source tile `src` to destination tile `dst` using Vector Engine.

Both source and destination tiles can be in either SBUF or PSUM, and passed in by reference as arguments. In-place shuffle is allowed, i.e., `dst` same as `src`. `shuffle_mask` is a 32-element list. Each mask element must be in data type `int` or affine expression. `shuffle_mask[i]` indicates which input partition the output partition `[i]` copies from within each 32-partition quadrant. The special value `shuffle_mask[i]=255` means the output tensor in partition `[i]` will be unmodified. `nc_stream_shuffle` can be applied to multiple of quadrants. In the case with more than one quadrant, the shuffle is applied to each quadrant independently, and the same `shuffle_mask` is used for each quadrant. `mask` applies to `dst`, meaning that locations masked out by `mask` will be unmodified. For more information about the cross-partition data movement, see [Cross-partition Data Movement](#).

This API has 3 constraints on `src` and `dst`:

1. `dst` must have same data type as `src`.

2. `dst` must have the same number of elements per partition as `src`.
3. The access start partition of `src` (`src_start_partition`), does not have to match or be in the same quadrant as that of `dst` (`dst_start_partition`). However, `src_start_partition/dst_start_partition` needs to follow some special hardware rules with the number of active partitions `num_active_partitions`. `num_active_partitions = ceil(max(src_num_partitions, dst_num_partitions)/32) * 32`, where `src_num_partitions` and `dst_num_partitions` refer to the number of partitions the `src` and `dst` tensors access respectively. `src_start_partition/dst_start_partition` is constrained based on the value of `num_active_partitions`:
  - If `num_active_partitions` is 96/128, `src_start_partition/dst_start_partition` must be 0.
  - If `num_active_partitions` is 64, `src_start_partition/dst_start_partition` must be 0/64.
  - If `num_active_partitions` is 32, `src_start_partition/dst_start_partition` must be 0/32/64/96.

#### Estimated instruction cost:

`max(MIN_II, N)` Vector Engine cycles, where `N` is the number of elements per partition in `src`, and `MIN_II` is the minimum instruction initiation interval for small input tiles. `MIN_II` is roughly 64 engine cycles.

#### Parameters

- **`src`** – the source tile
- **`dst`** – the destination tile
- **`shuffle_mask`** – a 32-element list that specifies the shuffle source and destination partition
- **`dtype`** – (optional) data type to cast the output type to (see [Supported Data Types](#) for more information); if not specified, it will default to be the same as the data type of the input tile.
- **`mask`** – (optional) a compile-time constant predicate that controls whether/how this instruction is executed (see [NKI API Masking](#) for details)

#### Example:

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl
from neuronxcc.nki.typing import tensor

#####
# Example 1:
# Apply cross-partition data movement to a 32-partition tensor,
# in-place shuffling the data in partition[i] to partition[(i+1)%32].
#####

...
a: tensor[32, 128] = nl.load(in_tensor)
a_mgrid = nl.mgrid[0:32, 0:128]
shuffle_mask = [(i - 1) % 32 for i in range(32)]
nisa.nc_stream_shuffle(src=a[a_mgrid.p, a_mgrid.x], dst=a[a_mgrid.p, a_mgrid.x],
↳ shuffle_mask=shuffle_mask)

nl.store(out_tensor, value=a)
```

```
#####
# Example 2:
# Broadcast data in 1 partition to 32 partitions.
#####
```

(continues on next page)

(continued from previous page)

```

...
a: tensor[1, 128] = nl.load(in_tensor)
b = nl.ndarray(shape=(32, 128), dtype=np.float32)
dst_mgrid = nl.mgrid[0:32, 0:128]
src_mgrid = nl.mgrid[0:1, 0:128]
shuffle_mask = [0] * 32
nisa.nc_stream_shuffle(src=a[0, src_mgrid.x], dst=b[dst_mgrid.p, dst_mgrid.x],
↳shuffle_mask=shuffle_mask)

nl.store(out_tensor, value=b)

```

```

#####
# Example 3:
# In the case where src and dst access more than one quadrant (32
# partitions), the shuffle is applied to each quadrant independently,
# and the same shuffle_mask is used for each quadrant.
#####

...
a: tensor[128, 128] = nl.load(in_tensor)
b = nl.ndarray(shape=(128, 128), dtype=np.float32)
mgrid = nl.mgrid[0:128, 0:128]
shuffle_mask = [(i - 1) % 32 for i in range(32)]
nisa.nc_stream_shuffle(src=a[mgrid.p, mgrid.x], dst=b[mgrid.p, mgrid.x], shuffle_
↳mask=shuffle_mask)

nl.store(out_tensor, value=b)

```

This document is relevant for: Inf2, Trn1, Trn2

## NKI ISA Config Enums

|                         |  |
|-------------------------|--|
| <code>engine</code>     | Neuron Device engines                    |
| <code>reduce_cmd</code> | Engine Register Reduce commands          |
| <code>dge_mode</code>   | Neuron Descriptor Generation Engine Mode |

This document is relevant for: Inf2, Trn1, Trn2

## nki.isa.engine

**class** `nki.isa.engine`(value)

Neuron Device engines

**Attributes**

|         |                |
|---------|----------------|
| tensor  | Tensor Engine  |
| vector  | Vector Engine  |
| scalar  | Scalar Engine  |
| gpsimd  | GpSIMD Engine  |
| sync    | Sync Engine    |
| unknown | Unknown Engine |

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.reduce\_cmd**

**class** nki.isa.reduce\_cmd(value)  
Engine Register Reduce commands

**Attributes**

|              |   |
|--------------|---|
| idle         | Not using the accumulator registers   |
| reset        | Resets the accumulator registers to its initial state   |
| reset_reduce | Resets the accumulator registers then immediately accumulate the results of the current instruction into the accumulators |
| reduce       | keeps accumulating over the current value of the accumulator registers  |

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.dge\_mode**

**class** nki.isa.dge\_mode(value)  
Neuron Descriptor Generation Engine Mode

**Attributes**

|         |  |
|---------|--|
| none    | Not using DGE  |
| swdge   | Software DGE   |
| hwdge   | Hardware DGE   |
| unknown | Unknown DGE mode, i.e., let compiler decide the DGE mode |

*This document is relevant for: Inf2, Trn1, Trn2*

**Target**

|                             |  |
|-----------------------------|--|
| <code>nc_version</code>     | NeuronCore version   |
| <code>get_nc_version</code> | Returns the <code>nc_version</code> of the current target context. |

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.nc\_version**

**class** `nki.isa.nc_version(value)`

NeuronCore version

`__init__()`

**Attributes**

|                   |                  |
|-------------------|------------------|
| <code>gen2</code> | Trn1/Inf2 target |
| <code>gen3</code> | Trn2 target      |

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.isa.get\_nc\_version**

`nki.isa.get_nc_version()`

Returns the `nc_version` of the current target context.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

**nki.compiler**

## Allocation Control

|                              |   |
|------------------------------|---|
| <code>sbuf.alloc</code>      | Allocate SBUF memory space for each logical block in a tensor using a customized allocation method. |
| <code>sbuf.mod_alloc</code>  | Allocate SBUF memory space for each logical tile in a tensor through modulo allocation.             |
| <code>sbuf.auto_alloc</code> | Returns a maker to indicate the tensor should be automatically allocated by compiler.               |
| <code>psum.alloc</code>      | Allocate PSUM memory space for each logical block in a tensor using a customized allocation method. |
| <code>psum.mod_alloc</code>  | Allocate PSUM memory space for each logical block in a tensor through modulo allocation.            |
| <code>psum.auto_alloc</code> | Returns a maker to indicate the tensor should be automatically allocated by compiler.               |

This document is relevant for: Inf2, Trn1, Trn2

### nki.compiler.sbuf.alloc

`nki.compiler.sbuf.alloc(func)`

Allocate SBUF memory space for each logical block in a tensor using a customized allocation method.

This is one of the NKI direction allocation APIs. We recommend reading *NKI Direct Allocation Developer Guide* before using these APIs.

In NKI, a SBUF tensor (declared using *NKI tensor creation APIs*) can have three kinds of dimensions, in order: logical block(B), partition(P), and free(F). The partition and free dimensions directly map to the SBUF dimensions. Both B and F can be multi-dimensional, while P must be one-dimensional per Neuron ISA constraints. The block dimension describes how many (P, F) logical tiles this tensor has, but does not reflect the number of physical tiles being allocated.

`ncc.sbuf.alloc` should be assigned to the `buffer` field of a NKI tensor declaration API. For example,

```
nki_tensor = nl.ndarray((4, 8, nl.par_dim(128), 4, 32), dtype=nl.bfloat16,
↪buffer=ncc.sbuf.alloc(...))
```

`ncc.sbuf.alloc` allows programmers to specify the physical location of each logical tile in the tensor. The API accepts a single input `func` parameter, which is a callable object that takes in:

1. a tuple of integers `idx` representing a logical block index,
2. an integer `pdim_size` for the number of partitions the logical tile has, and
3. an integer `fdim_size` for the number of bytes the logical tile has per partition.

The number of integers in `idx` must match the number of B dimensions the SBUF tensor has. For example, for the above `nki_tensor`, we expect the `idx` tuple to have two integers for a 2D block index.

`pdim_size` should match the partition dimension size of the NKI tensor exactly. `fdim_size` should be the total size of F dimension shapes of each logical tile in the tensor, multiplied by the data type size in bytes. For the above `sbuf_tensor`, `pdim_size` should be 128, and `fdim_size` should be `4*32*sizeof(nl.bfloat16) = 256` bytes.

The `func` callable must return a tuple of two integers (`start_partition`, `byte_addr`) indicating the physical tile location for the input logical block index. `start_partition` indicates the lowest partition the physical tile allocation starts from and must follow the these ISA rules:

- If `64 < pdim_size <= 128`, `start_partition` must be 0
- If `32 < pdim_size <= 64`, `start_partition` must be 0 or 64

- If  $0 < \text{pdim\_size} \leq 32$ , `start_partition` must be one of 0/32/64/96

The `byte_addr` indicates the byte offset into each partition the physical tile starts from. On NeuronCore-v2, a valid `byte_addr` can be any integer values from 0 (inclusive) to  $192\text{KiB} - 16\text{KiB} = (192 - 16) * 1024$  (exclusive). 192KiB is the physical size of a SBUF partition (defined in [architecture guide](#)) and 16KiB is allocated for compiler internal usage. In addition, the `base_addr` must be aligned to `nki.language.constants.sbuf_min_align`.

**Note:** In current release, programmers cannot mix NKI tensor declarations using automatic allocation (`ncc.sbuf.auto_alloc()` or the PSUM variant) and direction allocation APIs (`ncc.sbuf.alloc()`, `ncc.sbuf.mod_alloc()` or the PSUM variants) in the same kernel.

### Parameters

**func** – a callable object to specify how to place the logical block in SBUF memory.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.compiler.sbuf.mod\_alloc

`nki.compiler.sbuf.mod_alloc(*, base_addr, base_partition=0, num_par_tiles=(), num_free_tiles=())`

Allocate SBUF memory space for each logical tile in a tensor through modulo allocation.

This is one of the NKI direction allocation APIs. We recommend reading [NKI Direct Allocation Developer Guide](#) before using these APIs.

This API is equivalent to calling `nisa.compiler.alloc()` with a callable `psum_modulo_alloc_func` as defined below.

```

1 from typing import Optional, Tuple
2 from functools import reduce
3 from operator import mul
4 import unittest
5
6 def num_elms(shape):
7     return reduce(mul, shape, 1)
8
9 def linearize(shape, indices):
10    return sum(i * num_elms(shape[dim+1:]) for dim, i in enumerate(indices))
11
12 def modulo_allocate_func(base, allocate_shape, scale):
13     def func(indices):
14         if not allocate_shape:
15             # default shape is always (1, 1, ...)
16             allocate_shape_ = (1, ) * len(indices)
17         else:
18             allocate_shape_ = allocate_shape
19         mod_idx = tuple(i % s for i, s in zip(indices, allocate_shape_))
20         return linearize(shape=allocate_shape_, indices=mod_idx) * scale + base
21     return func
22
23 def mod_alloc(base_addr: int, *,
24               base_partition: Optional[int] = 0,
```

(continues on next page)

(continued from previous page)

```

25         num_par_tiles: Optional[Tuple[int, ...]] = (),
26         num_free_tiles: Optional[Tuple[int, ...]] = ()):
27     def sbuf_modulo_alloc_func(idx, pdim_size, fdim_size):
28         return (modulo_allocate_func(base_partition, num_par_tiles, pdim_size)(idx),
29               modulo_allocate_func(base_addr, num_free_tiles, fdim_size)(idx))
30     return sbuf_modulo_alloc_func
31

```

Here's an example usage of this API:

```

nki_tensor = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
                        buffer=nki.compiler.sbuf.mod_alloc(base_addr=0, num_free_
→tiles=(2, )))

for i_block in nl.affine_range(4):
    nki_tensor[i_block, :, :] = nl.load(...)
    ...                               = nl.exp(nki_tensor[i_block, :, :])

```

This produces the following allocation:

Table 7.3: Modulo Allocation Example

| Logical Tile Index | Physical Tile start_partition | Physical Tile byte_addr   |
|--------------------|-------------------------------|---|
| (0, )              | 0                             | $0 + (0 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 0$    |
| (1, )              | 0                             | $0 + (1 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 1024$ |
| (2, )              | 0                             | $0 + (2 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 0$    |
| (3, )              | 0                             | $0 + (3 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 1024$ |

With above scheme, we are able to implement double buffering in `nki_tensor`, such that `nl.load` in one iteration can write to one physical tile while `nl.exp` of the previous iteration can read from the other physical tile simultaneously.

---

**Note:** In current release, programmers cannot mix NKI tensor declarations using automatic allocation (`ncc.sbuf.auto_alloc()` or the PSUM variant) and direction allocation APIs (`ncc.sbuf.alloc()`, `ncc.sbuf.mod_alloc()` or the PSUM variants).

---

#### Parameters

- **base\_addr** – the base address in the free(F) dimension of the SBUF in bytes.
- **base\_partition** – the partition where the physical tile starts from. Must be 0 in the current version.
- **num\_par\_tiles** – the number of physical tiles on the partition dimension of SBUF allocated for the tensor. The length of the tuple must be empty or equal to the length of block dimension for the tensor.
- **num\_free\_tiles** – the number of physical tiles on the free dimension of SBUF allocated for the tensor. The length of the tuple must be empty or equal to the length of block dimension for the tensor.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2



## nki.compiler.sbuf.auto\_alloc

### nki.compiler.sbuf.auto\_alloc()

Returns a maker to indicate the tensor should be automatically allocated by compiler. All SBUF tensors in a kernel must either all be marked as `auto_alloc()`, or all be allocated with `alloc` or `mod_alloc`.

Initialize a tensor with `buffer=nl.sbuf` is equivalent to `buffer=ncc.sbuf.auto_alloc()`.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.compiler.psum.alloc

### nki.compiler.psum.alloc(func)

Allocate PSUM memory space for each logical block in a tensor using a customized allocation method.

This is one of the NKI direction allocation APIs. We recommend reading [NKI Direct Allocation Developer Guide](#) before using these APIs.

In NKI, a PSUM tensor (declared using [NKI tensor creation APIs](#)) can have three kinds of dimensions, in order: logical block(B), partition(P), and free(F). The partition and free dimensions directly map to the PSUM dimensions. Both B and F can be multi-dimensional, while P must be one-dimensional per Neuron ISA constraints. The block dimension describes how many (P, F) logical tiles this tensor has, but does not reflect the number of physical tiles being allocated.

`ncc.psum.alloc` should be assigned to the `buffer` field of a NKI tensor declaration API. For example,

```
nki_tensor = nl.ndarray((2, 4, nl.par_dim(128), 512), dtype=nl.float32, buffer=ncc.  
    ↪ psum.alloc(...))
```

`ncc.psum.alloc` allows programmers to specify the physical location of each logical tile in the tensor. The API accepts a single input `func` parameter, which is a callable object that takes in:

1. a tuple of integers `idx` representing a logical block index,
2. an integer `pdim_size` for the number of partitions the logical tile has, and
3. an integer `fdim_size` for the number of bytes the logical tile has per partition.

The number of integers in `idx` must match the number of B dimensions the PSUM tensor has. For example, for the above `nki_tensor`, we expect the `idx` tuple to have two integers for a 2D block index.

`pdim_size` should match the partition dimension size of the NKI tensor exactly. `fdim_size` should be the total size of F dimension shapes of each logical tile in the tensor, multiplied by the data type size in bytes. For the above `nki_tensor`, `pdim_size` should be 128, and `fdim_size` should be `512*sizeof(nl.float32) = 2048` bytes.

---

**Note:** In current release, `fdim_size` cannot exceed 2KiB, which is the size of a single PSUM bank per partition. Therefore, a physical PSUM tile cannot span multiple PSUM banks. Check out [Trainium/Inferentia2 Architecture Guide for NKI](#) for more information on PSUM banks.

---

The `func` returns a tuple of three integers (`bank_id`, `start_partition`, `byte_addr`) indicating the physical tile location for the input logical block index.

`bank_id` indicates the PSUM bank ID of the physical tile. `start_partition` indicates the lowest partition the physical tile allocation starts from. The `byte_addr` indicates the byte offset into each PSUM bank per partition the physical tile starts from.

---

**Note:** In current release, `start_partition` and `byte_addr` must both be 0.

---



---

**Note:** In current release, programmers cannot mix NKI tensor declarations using automatic allocation (`ncc.psum.auto_alloc()` or the SBUF variant) and direction allocation APIs (`ncc.psum.alloc()`, `ncc.psum.mod_alloc()` or the SBUF variants) in the same kernel.

---

### Parameters

**func** – a callable object to specify how to place the logical block in PSUM memory.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## nki.compiler.psum.mod\_alloc

```
nki.compiler.psum.mod_alloc(*, base_bank, base_addr=0, base_partition=0, num_bank_tiles=(),
                             num_par_tiles=(), num_free_tiles=())
```

Allocate PSUM memory space for each logical block in a tensor through modulo allocation.

This is one of the NKI direction allocation APIs. We recommend reading *NKI Direct Allocation Developer Guide* before using these APIs.

This API is equivalent to calling `nki.compiler.psum.alloc()` with a callable `psum_modulo_alloc_func` as defined below.

```
1 from typing import Optional, Tuple
2 from functools import reduce
3 from operator import mul
4 import unittest
5
6 def num_elems(shape):
7     return reduce(mul, shape, 1)
8
9 def linearize(shape, indices):
10     return sum(i * num_elems(shape[dim+1:]) for dim, i in enumerate(indices))
11
12 def modulo_allocate_func(base, allocate_shape, scale):
13     def func(indices):
14         if not allocate_shape:
15             # default shape is always (1, 1, ...)
16             allocate_shape_ = (1, ) * len(indices)
17         else:
18             allocate_shape_ = allocate_shape
19         mod_idx = tuple(i % s for i, s in zip(indices, allocate_shape_))
20         return linearize(shape=allocate_shape_, indices=mod_idx) * scale + base
21     return func
22
23 def mod_alloc(base_addr: int, *,
24               base_bank: Optional[int] = 0,
25               num_bank_tiles: Optional[Tuple[int]] = (),
26               base_partition: Optional[int] = 0,
```

(continues on next page)

(continued from previous page)

```

27         num_par_tiles: Optional[Tuple[int]] = (),
28         num_free_tiles: Optional[Tuple[int]] = ()):
29     def psum_modulo_alloc_func(idx, pdim_size, fdim_size):
30         # partial bank allocation is not allowed
31         return (modulo_allocate_func(base_bank, num_bank_tiles, 1)(idx),
32               modulo_allocate_func(base_partition, num_par_tiles, pdim_size)(idx),
33               modulo_allocate_func(base_addr, num_free_tiles, fdim_size)(idx))
34     return psum_modulo_alloc_func
35

```

Here's an example usage of this API:

```

psum_tensor = nl.ndarray((4, nl.par_dim(128), 512), dtype=nl.float32,
                        buffer=ncc.psum.mod_alloc(base_bank=0,
                                                  base_addr=0,
                                                  num_bank_tiles=(2,)))

for i_block in nl.affine_range(4):
    psum[i_block, :, :] = nisa.nc_matmul(...)
    ...                      = nl.exp(psum[i_block, :, :])

```

This produces the following allocation:

Table 7.4: Modulo Allocation Example

| Logical Tile Index | Physical Tile bank_id | Physical Tile start_partition | Physical Tile byte_addr |
|--------------------|-----------------------|-------------------------------|-------------------------|
| (0,)               | 0                     | 0                             | 0                       |
| (1,)               | 1                     | 0                             | 0                       |
| (2,)               | 0                     | 0                             | 0                       |
| (3,)               | 1                     | 0                             | 0                       |

With above scheme, we are able to implement double buffering in `nki_tensor`, such that `nisa.nc_matmul` in one iteration can write to one physical tile while `nl.exp` of the previous iteration can read from the other physical tile simultaneously.

**Note:** In current release, programmers cannot mix NKI tensor declarations using automatic allocation (`ncc.psum.auto_alloc()` or the SBUF variant) and direction allocation APIs (`ncc.psum.alloc()`, `ncc.psum.mod_alloc()` or the SBUF variants).

#### Parameters

- **base\_addr** – the base address in bytes along the free(F) dimension of the PSUM bank. Must be 0 in the current version.
- **base\_bank** – the base bank ID that the physical tiles start from.
- **num\_bank\_tiles** – the number of PSUM banks allocated for the tensor.
- **base\_partition** – the partition ID the physical tiles start from. Must be 0 in the current version.
- **num\_par\_tiles** – the number of physical tiles along the partition dimension allocated for the tensor. The length of the tuple must be empty or equal to the length of block dimension for the tensor. Currently must be an empty tuple or (1, 1, ...).
- **num\_free\_tiles** – the number of physical tiles on the free dimension per PSUM bank allocated for the tensor. The length of the tuple must be empty or equal to the length of

block dimension for the tensor. Currently must be an empty tuple or (1, 1, ...).

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### **nki.compiler.psum.auto\_alloc**

**nki.compiler.psum.auto\_alloc()**

Returns a maker to indicate the tensor should be automatically allocated by compiler. All PSUM tensors in a kernel must either all be marked as `auto_alloc()`, or all be allocated with `alloc` or `mod_alloc`.

Initialize a tensor with `buffer=n1.psum` is equivalent to `buffer=ncc.psum.auto_alloc()`.

*This document is relevant for:* Inf2, Trn1, Trn2

### **Kernel Decorators**

|  |  |
|--|--|
| <code>skip_middle_end_transformations</code> | Skip all middle end transformations on the kernel                        |
| <code>enable_stack_allocator</code>          | Use stack allocator to allocate the psum and sbuf tensors in the kernel. |
| <code>force_auto_alloc</code>                | Force automatic allocation to be turned on in the kernel.                |

*This document is relevant for:* Inf2, Trn1, Trn2

### **nki.compiler.skip\_middle\_end\_transformations**

**nki.compiler.skip\_middle\_end\_transformations**(*func=None*)

Skip all middle end transformations on the kernel

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### **nki.compiler.enable\_stack\_allocator**

**nki.compiler.enable\_stack\_allocator**(*func=None, log\_level=50*)

Use stack allocator to allocate the psum and sbuf tensors in the kernel.

Must use together with `skip_middle_end_transformations`.

```
from neuronxcc import nki

@nki.compiler.enable_stack_allocator
@nki.compiler.skip_middle_end_transformations
@nki.jit
def kernel(...):
    ...
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## nki.compiler.force\_auto\_alloc

`nki.compiler.force_auto_alloc(func=None)`

Force automatic allocation to be turned on in the kernel.

This will ignore any direct allocation inside the kernel

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI API Common Fields

### Supported Data Types

*Supported Data Types by NKI* below lists all supported data types by NKI. Almost all the NKI APIs accept a data type field, `dtype`, which can either be a NumPy equivalent type or a `nki.language` data type.

Table 7.5: Supported Data Types by NKI

|         | Data Type                                       | Accepted dtype Field by NKI APIs                                |
|---------|---|---|
| Integer | 8-bit unsigned integer                          | <code>nki.language.uint8</code> or <code>numpy.uint8</code>     |
|         | 8-bit signed integer                            | <code>nki.language.int8</code> or <code>numpy.int8</code>       |
|         | 16-bit unsigned integer                         | <code>nki.language.uint16</code> or <code>numpy.uint16</code>   |
|         | 16-bit signed integer                           | <code>nki.language.int16</code> or <code>numpy.int16</code>     |
|         | 32-bit unsigned integer                         | <code>nki.language.uint32</code> or <code>numpy.uint32</code>   |
|         | 32-bit signed integer                           | <code>nki.language.int32</code> or <code>numpy.int32</code>     |
| Float   | <code>float8_e4m3 (1S,4E,3M)<sup>2</sup></code> | <code>nki.language.float8_e4m3</code>                           |
|         | <code>float8_e5m2 (1S,5E,2M)</code>             | <code>nki.language.float8_e5m2</code>                           |
|         | <code>float16 (1S,5E,10M)</code>                | <code>nki.language.float16</code> or <code>numpy.float16</code> |
|         | <code>bfloat16 (1S,8E,7M)</code>                | <code>nki.language.bfloat16</code>                              |
|         | <code>tfloat32 (1S,8E,10M)</code>               | <code>nki.language.tfloat32</code>                              |
|         | <code>float32 (1S,8E,23M)</code>                | <code>nki.language.float32</code> or <code>numpy.float32</code> |
| Boolean | boolean stored as uint8                         | <code>nki.language.bool_</code> or <code>numpy.bool</code>      |

### Supported Math Operators for NKI ISA

*Supported Math Operators by NKI ISA* below lists all the mathematical operator primitives supported by NKI. Many *nki.isa* APIs (instructions) allow programmable operators through the `op` field. The supported operators fall into two categories: *bitvec* and *arithmetic*. In general, instructions using *bitvec* operators expect integer data types and treat input elements as bit patterns. On the other hand, instructions using *arithmetic* operators accept any valid NKI data types and convert input elements into float32 before performing the operators.

<sup>2</sup> S: sign bits, E: exponent bits, M: mantissa bits

Table 7.6: Supported Math Operators by NKI ISA

|            | Operator                    | op                                      | Legal Reduction op | Supported Engine      |
|------------|-----------------------------|---|--------------------|-----------------------|
| Bitvec     | Bitwise Not                 | <code>nki.language.invert</code>        | N                  | Vector                |
|            | Bitwise And                 | <code>nki.language.bitwise_and</code>   | Y                  | Vector                |
|            | Bitwise Or                  | <code>nki.language.bitwise_or</code>    | Y                  | Vector                |
|            | Bitwise Xor                 | <code>nki.language.bitwise_xor</code>   | Y                  | Vector                |
|            | Arithmetic Shift Left       | <code>nki.language.left_shift</code>    | N                  | Vector                |
|            | Arithmetic Shift Right      | Not supported                           | N                  | Vector                |
|            | Logical Shift Left          | <code>nki.language.left_shift</code>    | N                  | Vector                |
|            | Logical Shift Right         | <code>nki.language.right_shift</code>   | N                  | Vector                |
| Arithmetic | Add                         | <code>nki.language.add</code>           | Y                  | Vec-tor/GpSIMD/Scalar |
|            | Subtract                    | <code>nki.language.subtract</code>      | Y                  | Vector                |
|            | Multiply                    | <code>nki.language.multiply</code>      | Y                  | Vec-tor/GpSIMD/Scalar |
|            | Max                         | <code>nki.language.maximum</code>       | Y                  | Vector                |
|            | Min                         | <code>nki.language.minimum</code>       | Y                  | Vector                |
|            | Is Equal to                 | <code>nki.language.equal</code>         | N                  | Vector                |
|            | Is Not Equal to             | <code>nki.language.not_equal</code>     | N                  | Vector                |
|            | Is Greater than or Equal to | <code>nki.language.greater_equal</code> | N                  | Vector                |
|            | Is Greater than to          | <code>nki.language.greater</code>       | N                  | Vector                |
|            | Is Less than or Equal to    | <code>nki.language.less_equal</code>    | N                  | Vector                |
|            | Is Less than                | <code>nki.language.less</code>          | N                  | Vector                |
|            | Logical Not                 | <code>nki.language.logical_not</code>   | N                  | Vector                |
|            | Logical And                 | <code>nki.language.logical_and</code>   | Y                  | Vector                |
|            | Logical Or                  | <code>nki.language.logical_or</code>    | Y                  | Vector                |
|            | Logical Xor                 | <code>nki.language.logical_xor</code>   | Y                  | Vector                |
|            | Reverse Square Root         | <code>nki.language.rsqrt</code>         | N                  | GpSIMD/Scalar         |
|            | Reciprocal                  | <code>nki.language.reciprocal</code>    | N                  | Vector/Scalar         |
|            | Absolute                    | <code>nki.language.abs</code>           | N                  | Vector/Scalar         |
|            | Power                       | <code>nki.language.power</code>         | N                  | GpSIMD                |

**Note** Add and Multiply are supported on Scalar Engine only from NeuronCore-v3. 32-bit integer Add and Multiply are only supported on GpSIMD Engine.

## Supported Activation Functions for NKI ISA

*Supported Activation Functions by NKI ISA* below lists all the activation function supported by the `nki.isa`. `activation` API. These activation functions are approximated with piece-wise polynomials on Scalar Engine. *NOTE:* if input values fall outside the supported **Valid Input Range** listed below, the Scalar Engine will generate invalid output results.

Table 7.7: Supported Activation Functions by NKI ISA

| Function Name                | Accepted op by Scalar Engine  | Valid Input Range                                |
|------------------------------|---|--|
| Identity                     | <code>nki.language.copy</code> or <code>numpy.copy</code>             | <code>[-inf, inf]</code>                         |
| Square                       | <code>nki.language.square</code> or <code>numpy.square</code>         | <code>[-inf, inf]</code>                         |
| Sigmoid                      | <code>nki.language.sigmoid</code>                                     | <code>[-inf, inf]</code>                         |
| Relu                         | <code>nki.language.relu</code>  | <code>[-inf, inf]</code>                         |
| Gelu                         | <code>nki.language.gelu</code>  | <code>[-inf, inf]</code>                         |
| Gelu Derivative              | <code>nki.language.gelu_dx</code>                                     | <code>[-inf, inf]</code>                         |
| Gelu with Tanh Approximation | <code>nki.language.gelu_apprx_tanh</code>                             | <code>[-inf, inf]</code>                         |
| Silu                         | <code>nki.language.silu</code>  | <code>[-inf, inf]</code>                         |
| Silu Derivative              | <code>nki.language.silu_dx</code>                                     | <code>[-inf, inf]</code>                         |
| Tanh                         | <code>nki.language.tanh</code> or <code>numpy.tanh</code>             | <code>[-inf, inf]</code>                         |
| Softplus                     | <code>nki.language.softplus</code>                                    | <code>[-inf, inf]</code>                         |
| Mish                         | <code>nki.language.mish</code>  | <code>[-inf, inf]</code>                         |
| Erf                          | <code>nki.language.erf</code>   | <code>[-inf, inf]</code>                         |
| Erf Derivative               | <code>nki.language.erf_dx</code>                                      | <code>[-inf, inf]</code>                         |
| Exponential                  | <code>nki.language.exp</code> or <code>numpy.exp</code>               | <code>[-inf, inf]</code>                         |
| Natural Log                  | <code>nki.language.log</code> or <code>numpy.log</code>               | <code>[2<sup>-64</sup>, 2<sup>64</sup>]</code>   |
| Sine                         | <code>nki.language.sin</code> or <code>numpy.sin</code>               | <code>[-PI, PI]</code>                           |
| Arctan                       | <code>nki.language.arctan</code> or <code>numpy.arctan</code>         | <code>[-PI/2, PI/2]</code>                       |
| Square Root                  | <code>nki.language.sqrt</code> or <code>numpy.sqrt</code>             | <code>[2<sup>-116</sup>, 2<sup>118</sup>]</code> |
| Reverse Square Root          | <code>nki.language.rsqrt</code>                                       | <code>[2<sup>-87</sup>, 2<sup>97</sup>]</code>   |
| Reciprocal                   | <code>nki.language.reciprocal</code> or <code>numpy.reciprocal</code> | <code>±[2<sup>-42</sup>, 2<sup>42</sup>]</code>  |
| Sign                         | <code>nki.language.sign</code> or <code>numpy.sign</code>             | <code>[-inf, inf]</code>                         |
| Absolute                     | <code>nki.language.abs</code> or <code>numpy.abs</code>               | <code>[-inf, inf]</code>                         |

## NKI API Masking

All `nki.language` and `nki.isa` APIs accept an optional input field, `mask`. The `mask` field is an execution predicate known at compile-time, which informs the compiler to skip generating the instruction or generate the instruction with a smaller input tile shape. Masking is handled completely by Neuron compiler and hence does not incur any performance overhead in the generated instructions.

The `mask` can be created using comparison expressions (e.g., `a < b`) or multiple comparison expressions concatenated with `&` (e.g., `(a < b) & (c > d)`). The left- or right-hand side expression of each comparator must be an affine expression of `nki.language.arange()`, `nki.language.affine_range()` or `nki.language.program_id()`. Each comparison expression should indicate which range of indices along one of the input tile axes should be valid

for the computation. For example, assume we have an input tile `in_tile` of shape (128, 512), and we would like to perform a square operation on this tile for elements in `[0:64, 0:256]`, we can invoke the `nki.language.square()` API using the following:

```
import neuronxcc.nki.language as nl

...
i_p = nl.arange(128)[: , None]
i_f = nl.arange(512)[None, :]

out_tile = nl.square(in_tile, mask=((i_p<64) & (i_f<256)))
```

The above example will be lowered into a hardware ISA instruction that only processes 64x256 elements by Neuron Compiler.

The above mask definition works for most APIs where there is only one input tile or both input tiles share the same axes. One exception is the `nki.language.matmul` and similarly `nki.isa.nc_matmul` API, where the two input tiles `lhs` and `rhs` contain three unique axes:

1. The contraction axis: both `lhs` and `rhs` partition axis (`lhs_rhs_p`)
2. The first axis of matmul output: `lhs` free axis (`lhs_f`)
3. The second axis of matmul output: `rhs` free axis (`rhs_f`)

As an example, let's assume we have `lhs` tile of shape (`sz_p`, `sz_m`) and `rhs` tile of shape (`sz_p`, `sz_n`), and we call `nki.language.matmul` to calculate an output tile of shape (`sz_m`, `sz_n`):

```
import neuronxcc.nki.language as nl

i_p = nl.arange(sz_p)[: , None]

i_lhs_f = nl.arange(sz_m)[None, :]
i_rhs_f = nl.arange(sz_n)[None, :] # same as `i_rhs_f = i_lhs_f`

result = nl.matmul(lhs[i_p, i_lhs_f], rhs[i_p, i_rhs_f], transpose_x=True)
```

Since both `i_lhs_f` and `i_rhs_f` are identical to the Neuron Compiler, the Neuron Compiler cannot distinguish the two input axes if they were to be passed into the `mask` field directly.

Therefore, we introduce “operand masking” syntax for matmult APIs to let users to precisely define the masking on the inputs to the matmult APIs (currently only matmult APIs support operand masking, subject to changes in future releases). Let's assume we need to constraint `sz_m <= 64` and `sz_n <= 256`:

```
import neuronxcc.nki.language as nl

i_p = nl.arange(sz_p)[: , None]

i_lhs_f = nl.arange(sz_m)[None, :]
i_rhs_f = nl.arange(sz_n)[None, :] # same as `i_rhs_f = i_lhs_f`

i_lhs_f_virtual = nl.arange(sz_m)[None, :, None]

result = nl.matmul(lhs_T[i_lhs_f <= 64], rhs[i_rhs_f <= 256], transpose_x=True)
```

There are two notable use cases for masking:

1. When the tiling factor doesn't divide the tensor dimension sizes



## 2. Skip ineffectual instructions that compute known output values

We will present an example of the first use case below. Let's assume we would like to evaluate the exponential function on an input tensor of shape `[sz_p, sz_f]` from HBM. Since the input to `nki.language.load/nki.language.store/nki.language.exp` expects a tile with a partition axis size not exceeding `nki.language.tile_size.pmax == 128`, we should loop over the input tensor using a tile size of `[nki.language.tile_size.pmax, sz_f]`.

However, `sz_p` is not guaranteed to be an integer multiple of `nki.language.tile_size.pmax`. In this case, one option is to write a loop with trip count of `sz_p // nki.language.tile_size.pmax` followed by a single invocation of `nki.language.exp` with an input tile of shape `[sz_p % nki.language.tile_size.pmax, sz_f]`. This effectively “unrolls” the last instance of tile computation, which could lead to messy code in a complex kernel. Using masking here will allow us to avoid such unrolling, as illustrated in the example below:

```
import neuronxcc.nki.language as nl
from torch_neuronx import nki_jit

@nki_jit
def tensor_exp_kernel_(in_tensor, out_tensor):

    sz_p, sz_f = in_tensor.shape

    i_f = nl.arange(sz_f)[None, :]

    trip_count = math.ceil(sz_p/nl.tile_size.pmax)

    for p in nl.affine_range(trip_count):
        # Generate tensor indices for the input/output tensors
        # pad index to pmax, for simplicity
        i_p = p * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]

        # Load input data from external memory to on-chip memory
        # only read up to sz_p
        in_tile = nl.load(in_tensor[i_p, i_f], mask=(i_p < sz_p))

        # perform the computation
        out_tile = nl.exp(in_tile, mask=(i_p < sz_p))

        # store the results back to external memory
        # only write up to sz_p
        nl.store(out_tensor[i_p, i_f], value=out_tile, mask=(i_p<sz_p))
```

## NKI Type Promotion

When the data types (dtypes) of inputs to an arithmetic operation (i.e., add, multiply, tensor\_tensor, etc.) differ, we promote the dtypes following the rules below:

**(float, integer):** Pick the float type. Example:

- `(np.int32, np.float16) -> np.float16`
- `(np.uint16, nl.tfloat32) -> nl.tfloat32`

**(float, float):** Pick the wider float type or a new widened type that fits the values range. Example:

- `(np.float32, nl.tfloat32) -> np.float32`
- `(np.float32, nl.bfloat16) -> np.float32`

- (np.float16, nl.bfloat16) -> np.float32 (new widened type)
- (nl.float8\_e4m3, np.float16) -> np.float16
- (nl.float8\_e4m3, nl.bfloat16) -> nl.bfloat16
- (nl.float8\_e4m3, nl.float8\_e5m2) -> nl.bfloat16 (new widened type)

(int, int): Pick the wider type or a new widened type that fits the values range. Example:

- (np.int16, np.int32) -> np.int32
- (np.uint8, np.uint16) -> np.uint16
- (np.uint16, np.int32) -> np.int32
- (np.int8, np.uint8) -> np.int16 (new widened type)
- (np.int8, np.uint16) -> np.int32 (new widened type)
- (np.int32, np.uint32) -> np.float32 (new widened type is float32, since int64 isn't supported on the hardware)

The output of the arithmetic operation will get the promoted type by default.

**Note:** The Vector Engine internally performs most of the computation in FP32 (see [Vector Engine](#)) and casts the output back to the specific type.

```
x = np.ndarray((N, M), dtype=nl.float8_e4m3)
y = np.ndarray((N, M), dtype=np.float16)
z = nl.add(x, y) # calculation done in FP32, output cast to np.float16
assert z.dtype == np.float16
```

To prevent the compiler from automatically widening output dtype based on mismatching input dtypes, you may explicitly set the output dtype in the arithmetic operation API. This would be useful if the output is passed into another operation that benefits from a smaller dtype.

```
x = np.ndarray((N, M), dtype=nl.bfloat16)
y = np.ndarray((N, M), dtype=np.float16)
z = nl.add(x, y, dtype=nl.bfloat16) # without explicit `dtype`, `z.dtype` would have
↳ been np.float32
assert z.dtype == nl.bfloat16
```

## Weakly Typed Scalar Type Inference

Weakly typed scalars (scalar values where the type wasn't explicitly specified) will be inferred as the widest dtype supported by hardware:

- bool --> uint8
- integer --> int32
- floating --> float32

Doing an arithmetic operation with a scalar may result in a larger output type than expected, for example:

- (np.int8, 2) -> np.int32
- (np.float16, 1.2) -> np.float32

To prevent larger dtypes from being inferred from weak scalar types, do either of:

1. Explicitly set the datatype of the scalar, like np.int8(2), so that the output type is what you desire:

```
x = np.ndarray((N, M), dtype=np.float16)
y = np.float16(2)
z = nl.add(x, y)
assert z.dtype == np.float16
```

2. Explicitly set the output dtype of the arithmetic operation:

```
x = np.ndarray((N, M), dtype=np.int16)
y = 2
z = nl.add(x, y, dtype=nl.bfloat16)
assert z.dtype == nl.bfloat16
```

**Note:** The Vector Engine internally performs most of the computation in FP32 (see [Vector Engine](#)) and casts the output back to the specific type.

## NKI Engine Selection for Operators Supported on Multiple Engines

There is a tradeoff between precision and speed on different engines for operators with multiple engine options. Users can select which engine to map to based on their needs. We take reciprocal and reverse square root as two examples and explain the tradeoff below.

1. Reciprocal can run on Scalar Engine or Vector Engine:

Reciprocal can run on Vector Engine with `nki.isa.reciprocal` or on Scalar Engine with `nki.isa.activation(nl.reciprocal)`. Vector Engine performs reciprocal at a higher precision compared to Scalar Engine; however, the computation throughput of reciprocal on Vector Engine is about 8x lower than Scalar Engine for large input tiles. For input tiles with a small number of elements per partition (less than 64, processed one per cycle), instruction initiation interval (roughly 64 cycles) dominates performance so Scalar Engine and Vector Engine have comparable performance. In this case, we suggest using Vector Engine to achieve better precision.

**Estimated cycles on different engines:**

| Cost ( <i>Engine Cycles</i> ) | Condition  |
|-------------------------------|--|
| $\max(\text{MIN\_II}, N)$     | mapped to Scalar Engine <code>nki.isa.scalar_engine</code> |
| $\max(\text{MIN\_II}, 8*N)$   | mapped to Vector Engine <code>nki.isa.vector_engine</code> |

where,

- $N$  is the number of elements per partition in the input tile.
- $\text{MIN\_II}$  is the minimum instruction initiation interval for small input tiles.  $\text{MIN\_II}$  is roughly 64 engine cycles.

**Note** `nki.isa.activation(op=nl.reciprocal)` doesn't support setting bias on NeuronCore-v2.

2. Reverse square root can run on GpSIMD Engine or Scalar Engine:

Reverse square root can run on GpSIMD Engine with `nki.isa.tensor_scalar(op0=nl.rsqrt, operand0=0.0)` or on Scalar Engine with `nki.isa.activation(nl.rsqrt)`. GpSIMD Engine performs reverse square root at a higher precision compared to Scalar Engine; however, the computation throughput of reverse square root on GpSIMD Engine is 4x lower than Scalar Engine.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI API Errors

### err\_1d\_arange\_not\_supported

Indexing a NKI tensor with 1D arange is not supported.

NKI expects tile indices to have at least two dimensions to match the underlying memory (SBUF or PSUM)

```
tmp = nl.zeros((128, 1), dtype=nl.float32, buffer=nl.sbuf)
i = nl.arange(64)
c = nl.exp(tmp[i, 0]) # Error: indexing tensor `tmp` with 1d arange is not supported,
```

You can workaround the problem by introducing new axes like the following code:

```
tmp = nl.zeros((128, 1), dtype=nl.float32, buffer=nl.sbuf)
i = nl.arange(64)[:, None]
c = nl.exp(tmp[i, 0])
```

Or using simple slicing:

```
tmp = nl.zeros((128, 1), dtype=nl.float32, buffer=nl.sbuf)
c = nl.exp(tmp[0:64, 0])
```

### err\_activation\_bias\_invalid\_type

Bias parameter of activation or activation\_reduce must be a vector of type float32, float16, or bfloat16.

```
nisa.activation(op=nl.exp, data=data[...], bias=nisa.memset((128, 1), 1.2, dtype=np.
↪float32)) # ok
nisa.activation(op=nl.exp, data=data[...], bias=nisa.memset((128, 1), 1.2, dtype=nl.
↪bfloat16)) # ok
nisa.activation(op=nl.exp, data=data[...], bias=nisa.memset((128, 1), 1.2, dtype=np.
↪int8)) # not supported
```

### err\_activation\_scale\_invalid\_type

Scale parameter of activation or activation\_reduce must be a scalar or vector of type float32.

```
nisa.activation(op=nl.exp, data=data[...], scale=1.2) # ok
nisa.activation(op=nl.exp, data=data[...], scale=nisa.memset((128, 1), 1.2, dtype=np.
↪float32)) # ok
nisa.activation(op=nl.exp, data=data[...], scale=nisa.memset((128, 1), 1.2, dtype=np.
↪float16)) # not supported
```

## err\_activation\_scale\_scalar\_or\_vector

Scale parameter of activation must be either a scalar value or a 1D vector spanning the partition dimension.

```
nisa.activation(op=nl.exp, data=data[...], scale=1.2) # ok
nisa.activation(op=nl.exp, data=data[...], scale=nisa.memset((128, 1), 1.2, dtype=np.
↪float32)) # ok
nisa.activation(op=nl.exp, data=data[...], scale=nisa.memset((1, 128), 1.2, dtype=np.
↪float32)) # not supported
nisa.activation(op=nl.exp, data=data[...], scale=nisa.memset((128, 128), 1.2, dtype=np.
↪float32)) # not supported
```

## err\_ambiguous\_tensor\_truth\_value

ValueError: Cannot evaluate truth value of a multi-element tensor/array.

This error occurs in two common scenarios: 1. Using logical operators (and, or, not) on multi-element tensors 2. Using tensor objects in conditional statements without explicit None checks

Example of problematic code:

```
def func(a, b: Optional[tensor]):
    ix, iy = nl.mgrid[0:128, 0:128]
    a_tile: tensor[128, 128] = nl.load(a[ix, iy])

    not_a_tile = not (a_tile > 0) # The truth value of an array with more than one_
↪element is ambiguous

    if b: # The truth value of an array with more than one element is ambiguous
        pass
```

Correct usage:

For tensor operations, use appropriate element-wise operators: - Use ~ instead of not for boolean negation - Use & instead of and for logical AND - Use | instead of or for logical OR

For None checks, use explicit 'is' comparisons:

```
def func(a, b: Optional[tensor]):
    ix, iy = nl.mgrid[0:128, 0:128]
    a_tile: tensor[128, 128] = nl.load(a[ix, iy])

    not_a_tile = ~(a_tile > 0) # Element-wise negation

    if b is not None: # Explicit None check
        pass
```

### Note:

This error is similar to NumPy's ValueError for ambiguous truth value of arrays, as tensors/arrays can contain multiple boolean values and cannot be automatically reduced to a single boolean.

## err\_annotation\_shape\_mismatch

Tensor shape and the annotated shape mismatch

NKI check the object shape based on python type annotation in the *target: type = value* syntax, NKI will throw an error if the expected shape and the object shape mismatch.

For example:

```
import neuronxcc.nki.typing as nt
data: nt.tensor[128, 512] = nl.zeros((par_dim(128), 128), dtype=np.float32) # Error:
↳ shape of `data[128, 128]` does not match the expected shape of [128, 512]
```

## err\_bias\_tensor\_must\_be\_specified\_in\_allocation

Bias tensor of an activation op must be specified in allocated NKI kernels.

```
data = .... # assume data is of shape (128, 128)
exp = nl.ndarray((par_dim(128), 512), dtype=nl.bfloat16, buffer=ncc.sbuf.mod_alloc(base_
↳ addr=0))
exp[...] = nisa.activation(np.exp,
                           data=data[...]) # Error, bias argument must also be specified

exp[...] = nl.exp(data=data[...])
# Error, nl.exp maps to the the instruction as nisa.activation, must use nisa.activation,
↳ and specify bias tensor in allocation kernels
```

## err\_cannot\_assign\_to\_index

An *index* or *mask* tensor does not support item assignment.

**You may explicitly call *iota* to convert an *index* tensor to a normal *tile* before any assignments.**

```
_, x = nl.mgrid[0:1, 0:8]
x[0, 5] = 1024 # Error: 'index' tensor does not support item assignment
y = nisa.iota(x, dtype=nl.uint32)
y[0, 5] = 1024 # works
```

## err\_cannot\_update\_immutable\_parameter

Cannot update immutable parameter

By default, all parameters to the top level nki kernels are immutable, updating immutable parameters in the kernel is not allowed.

```
def kernel(in_tensor):
    x = nl.load(in_tensor)
    y = x + 1
    # Parameter `in_tensor` is immutable by default, cannot modify immutable parameter
    nl.store(in_tensor, value=y) # Error: Cannot update immutable parameter
```

(continues on next page)

(continued from previous page)

```
return in_tensor
```

You could explicitly annotate the parameter as mutable:

```
import neuronxcc.nki.typing as nt

def kernel(in_tensor: nt.mutable_tensor):
    x = nl.load(in_tensor)
    y = x + 1
    nl.store(in_tensor, value=y) # ok, in_tensor is mutable based on the annotation

    return in_tensor
```

Alternatively, you could return a copy of the input parameter if you don't intend to modify the input parameter:

```
import neuronxcc.nki.isa as nisa import neuronxcc.nki.language as nl

def kernel(in_tensor):
    out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
                             buffer=nl.shared_hbm)

    nisa.dma_copy(dst=out_tensor, src=in_tensor)

    x = nl.load(out_tensor) y = x + 1 nl.store(out_tensor, value=y) # ok
    return out_tensor
```

### err\_control\_flow\_condition\_dependent\_on\_arange

Control-flow depending on `nl.arange` or `nl.mgrid` is not supported.

```
for j0 in nl.affine_range(4096):
    i1 = nl.arange(512)[None, :]
    j = j0 * 512 + i1
    if j > 2048: # Error: Control-flow depending on `nl.arange` or `nl.mgrid` is not
    ↪ supported
        y = nl.add(x[0, j], x[0, j - 2048])
```

In the above example, `j` depends on the value of `i1`, which is `nl.arange(512)[None, :]`. NKI does not support using `nl.arange` or `nl.mgrid` in control-flow condition. To workaround this error, you can use the `mask` parameter:

```
for j0 in nl.affine_range(4096):
    i1 = nl.arange(512)[None, :]
    j = j0 * 512 + i1
    y = nl.add(x[0, j], x[0, j - 2048], mask=j > 2048)
```

### err\_copy\_dynamic\_indirect\_indices\_not\_natively\_supported

If indices are dynamic (i.e. known only at runtime), copying indirect memory references is not natively supported.

For example, when overloading the *assignment* operator as a copy operation, the following code will throw an error:

```
data_tensor: tensor[128, 8, 4] = nl.load(data_tensor)
idx_tile: tensor[8, 1] = nl.load(idx_tensor)
out_sbuf: tensor[8, 4] = nl.ndarray([8, 4], dtype=data_tensor.dtype,
                                   buffer=nl.sbuf)

iy, iz = nl.mgrid[0:8, 0:4]
out_sbuf[iy, iz] = data_tensor[0, idx_tile, iz] # idx_tile only known at runtime
```

To fix this error, consider using *nisa.tensor\_copy\_dynamic\_src* as follows.

```
out_sbuf[iy, iz] = nisa.tensor_copy_dynamic_src(data_tensor[0, idx_tile, iz])
```

### err\_dynamic\_control\_flow\_not\_supported

Dynamic control-flow depending on tensor value is currently not supported by NKI.

```
cnd = nl.load(a) # a have shape of [1, 1]
if cnd:          # Error: dynamic control-flow depending on tensor value is not
    ↪supported.
    nl.store(b, 1)
```

### err\_exceed\_max\_supported\_dimension

NKI API tensor parameter exceeds max supported number of dimensions.

Certain NKI APIs have restrictions on how many dimensions the tensor parameter can have:

```
x = nl.zeros(shape=[64, 32, 2], dtype=np.float32, buffer=nl.sbuf)
b = nl.transpose(x) # Error: parameter 'x[64, 32, 2]' of 'transpose' exceed max supported
    ↪number of dimensions of 2.

x = nl.zeros(shape=[64, 64], dtype=np.float32, buffer=nl.sbuf)
b = nl.transpose(x) # Works if input `x` only have 2 dimensions (i.e. rank=2)
```

### err\_failed\_to\_infer\_tile\_from\_local\_tensor

NKI requires inputs of all compute APIs to be valid tiles with the first dimension being the partition dimension.

```
# We mark the second dimension as the partition dimension
a = nl.zeros((4, nl.par_dim(8), 8), dtype=nl.float32, buffer=nl.sbuf)
c = nl.add(a, 32) # Error: Failed to infer tile from tensor 'a',
```

To fix the problem you can use index tensor *a* to generate a tile whose first dimension is the partition dimension



```
# We mark the second dimension of tensor a as the partition dimension
a = nl.zeros((4, nl.par_dim(8), 8), dtype=nl.float32, buffer=nl.sbuf)
c = nl.ndarray((4, nl.par_dim(8), 8), dtype=nl.float32, buffer=nl.sbuf)
for i in range(4):
    # result of `a[i]` is a tile with shape (8, 8) and the first dimension is the partition_
    ↪ dimension
    c[i] = nl.add(a[i], 32) # works
    # Or explicitly generate a tile with `nl.arange`
    ix = nl.arange(8)[: , None]
    iy = nl.arange(8)[None, :]
    # result of `a[i, ix, iy]` is a tile with shape (8, 8) and the first dimension is the_
    ↪ partition dimension
    c[i, ix, iy] = nl.add(a[i, ix, iy], 32) # also works
```

### err\_hbm\_tensor\_with\_init\_value\_not\_supported

Creating HBM tensor with init value is not supported.

```
t = nl.full((3, 128, 512), fill_value=1.0, buffer=nl.shared_hbm) # t on hbm and has an_
    ↪ init value
# Error: Creating HBM tensor with init value is not supported.
```

To work around the limitation you need to explicitly initialize the tensor with `nl.store`:

```
t = nl.ndarray((3, 128, 512), buffer=nl.shared_hbm)
for i in range(3):
    nl.store(dst=t[i, :, :], value=1.0)
```

### err\_indirect\_indices\_free\_dim

Dynamic indexing for load/store only supports the indirect indexing to be on the partition or block dimension. Refer to the code examples in [nl.load](#) and [nl.store](#).

Also, if you're using `nl.mgrid` you may get this error even though your indirect indexing was on the partition dimension, use `nl.arange` instead.

```
i_p, i_f = nl.mgrid[0:64, 0:512] # this won't work for dynamic access

i_p = nl.arange(64)[: , None] # this works for dynamic access
i_f = nl.arange(512)[None, :]

data_tile = nl.load(data_tensor[idx_tile[i_p, 0], i_f])
```

### err\_leading\_dimension\_of\_tensor\_must\_be\_partition

The leading dimension of SBUF/PSUM tensors must be the partition dimension.

NKI used to allow block dimensions in front of the partition dimension, but block dimension has been deprecated. Please refer to [Migration Guide for Block Dimension](#)

### err\_local\_variable\_used\_out\_of\_scope

Tensors in NKI are not allowed to be used outside of their parent scope.

Tensors in NKI have a stricter scope rules than Python. In NKI, control blocks in if/else/for statements will introduce their own scope for tensors. A tensor defined in if/else/for control blocks are not allowed to be used outside of the scope.

```
for i in range(4):
    if i < 2:
        tmp = nl.load(a)
    else:
        tmp = nl.load(b)

    nl.store(c, tmp) # Error: Local variable 'tmp' is referenced outside of its parent_
↪scope ...
```

To fix the problem, you can rewrite the above code as:

```
for i in range(4):
    tmp = nl.ndarray(shape=a.shape, dtype=a.dtype)
    if i < 2:
        tmp[...] = nl.load(a)
    else:
        tmp[...] = nl.load(b)

    nl.store(c, tmp)
```

This stricter scope rules may also introduce unexpected error like the following:

```
data = nl.zeros((par_dim(128), 128), dtype=np.float32)

for i in nl.sequential_range(4):
    i_tile = nisa.iota(i, dtype=nl.uint32).broadcast_to(data.shape)
    data = data + i_tile # Warning: shadowing local tensor 'float32 data[128, 128]' with a_
↪new object, use 'data[...]' if you want to update the existing object

nl.store(ptr, value=data) # # Error: Local variable 'tmp' is referenced outside of its_
↪parent scope ...
```

To fix the problem you can follow the suggestion from the warning

```
data = nl.zeros((par_dim(128), 128), dtype=np.float32)

for i in nl.sequential_range(4):
    i_tile = nisa.iota(i, dtype=nl.uint32).broadcast_to(data.shape)
    data[...] = data + i_tile
```

(continues on next page)

(continued from previous page)

```
nl.store(ptr, value=data)
```

### err\_mutable\_parameter\_not\_returned

A mutable kernel parameter must be returned by *return*

```
import neuronxcc.nki.typing as nt

def kernel(in_tensor: nt.mutable_tensor):
    out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
                            buffer=nl.shared_hbm)

    x = nl.load(in_tensor)
    y = x + 1
    nl.store(out_tensor, value=y)
    # Also update mutable parameter `in_tensor`
    nl.store(in_tensor, value=y)

    # But didnt return the mutable parameter `in_tensor`
    return out_tensor # Error: Mutable kernel parameter not returned by `return` statement
```

To fix this error, you need to return the mutable parameter to the returned list:

```
import neuronxcc.nki.typing as nt

def kernel(in_tensor: nt.mutable_tensor):
    out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
                            buffer=nl.shared_hbm)

    x = nl.load(in_tensor)
    y = x + 1
    nl.store(out_tensor, value=y)
    nl.store(in_tensor, value=y)

    return out_tensor, in_tensor # ok
```

### err\_nested\_kernel\_with\_spmd\_grid

Calling a NKI kernel with a SPMD grid from another NKI kernel is not supported.

```
@nki.trace
def kernel0(...):
    ...

@nki.trace
def kernel1(...):
    ...

@nki_jit
def kernel_top():
    kernel0(...) # works
```

(continues on next page)

(continued from previous page)

```
kernel1[4, 4](...) # Error: Calling kernel with spmd grid (kernel1[4,4]) inside_
↳ another kernel is not supported
```

### err\_nki\_api\_outside\_of\_nki\_kernel

Calling NKI API outside of NKI kernels is not supported.

Make sure the NKI kernel function decorated with *nki.jit*.

### err\_num\_partition\_exceed\_arch\_limit

Number of partitions exceeds architecture limitation.

NKI requires the number of partitions of a tile to not exceed the architecture limitation of 128

For example in Trainium:

```
x = nl.zeros(shape=[256, 1024], dtype=np.float32, buffer=nl.sbuf) # Error: number of_
↳ partitions 256 exceed architecture limitation of 128.
x = nl.zeros(shape=[128, 1024], dtype=np.float32, buffer=nl.sbuf) # Works
```

### err\_num\_partition\_mismatch

Number of partitions mismatch.

Most of the APIs in the *nki.isa* module require all operands to have the same number of partitions. For example, the *nki.isa.tensor\_tensor()* requires all operands to have the same number of partitions.

```
x = nl.zeros(shape=[128, 512], dtype=np.float32, buffer=nl.sbuf)
y0 = nl.zeros(shape=[1, 512], dtype=np.float32, buffer=nl.sbuf)
z = nisa.tensor_tensor(x, y0, op=nl.add) # Error: number of partitions (dimension 0 size_
↳ of a tile) mismatch in parameters (data1[128, 512], data2[1, 512]) of 'tensor_tensor'

y1 = y0.broadcast_to([128, 512]) # Call 'broadcast_to' to explicitly broadcast on_
↳ the partition dimension
z = nisa.tensor_tensor(x, y1, op=nl.add) # works because x and y1 has the same number of_
↳ partitions
```

### err\_shared\_hbm\_must\_in\_kernel\_level

shared\_hbm tensor can only be created in top level kernel scope

Creating shared\_hbm tensors inside a loop, under if condition or inside another function called by the top-level nki kernel is not supported.

Consider hoist the creation of shared\_hbm tensors to the top level kernel scope.

```
@nki.jit
def kernel(...):
    a = nl.ndarray((128, 512), dtype=nl.float32,
```

(continues on next page)

(continued from previous page)

```

        buffer=nl.shared_hbm) # works

for i in range(8):
    b = nl.ndarray((128, 512), dtype=nl.float32,
                   buffer=nl.shared_hbm) # Error: shared_hbm buffer can only be created_
↳top level kernel scope

    if nl.program_id(0) >= 1:
        c = nl.ndarray((128, 512), dtype=nl.float32,
                       buffer=nl.shared_hbm) # Error: shared_hbm buffer can only be created_
↳top level kernel scope

    # Call another function
    func(...)

def func(...):
    d = nl.ndarray((128, 512), dtype=nl.float32,
                   buffer=nl.shared_hbm) # Error: shared_hbm buffer can only be created top_
↳level kernel scope

```

### err\_size\_of\_dimension\_exceed\_arch\_limit

Size of dimension exceeds architecture limitation.

Certain NKI APIs have restrictions on dimension sizes of the parameter tensor:

```

x = nl.zeros(shape=[128, 512], dtype=np.float32, buffer=nl.sbuf)
b = nl.transpose(x) # Error: size of dimension 1 in 'x[128, 512]' of 'transpose' exceed_
↳architecture limitation of 128.

x = nl.zeros(shape=[128, 128], dtype=np.float32, buffer=nl.sbuf)
b = nl.transpose(x) # Works size of dimension 1 < 128

```

### err\_store\_dst\_shape\_smaller\_than\_other\_shape

Illegal shape in assignment destination.

The destination of assignment must have the same or bigger shape than the source of assignment. Assigning multiple values to the same element in the assignment destination from a single NKI API is not supported

```

x = nl.zeros(shape=(128, 512), dtype=nl.float32, buffer=nl.sbuf)
y = nl.zeros(shape=(128, 1), dtype=nl.float32, buffer=nl.sbuf)

y[...] = x # Error: Illegal assignment destination shape in 'a = b': shape [128, 1] of_
↳parameter 'a' is smaller than other parameter shapes b[128, 512].
x[...] = y # ok, if we are broadcasting from source to the destination of the assignment

```

## err\_tensor\_access\_out\_of\_bound

Tensor access out-of-bound.

Out-of-bound access is considered illegal in NKI. When the indices are calculated from nki indexing APIs, out-of-bound access results in a compile-time error. When the indices are calculated dynamically at run-time, such as indirect memory accesses, out-of-bound access results in run-time exceptions during execution of the kernel.

```
x = nl.ndarray([128, 4000], dtype=np.float32, buffer=nl.hbm)
for i in nl.affine_range((4000 + 512 - 1) // 512):
    tile = nl.mgrid[0:128, 0:512]
    nl.store(x[tile.p, i * 512 + tile.x], value=0) # Error: Out-of-bound access for
    ↪ tensor `x` on dimension 1: index range [0, 4095] exceed dimension size of 4000
```

You could carefully check the corresponding indices and make necessary correction. If the indices are correct and intentional, out-of-bound access can be avoided by providing a proper mask:

```
x = nl.ndarray([128, 4000], dtype=np.float32, buffer=nl.hbm)
for i in nl.affine_range((4000 + 512 - 1) // 512):
    tile = nl.mgrid[0:128, 0:512]
    nl.store(x[tile.p, i * 512 + tile.x], value=0,
            mask=i * 512 + tile.x < 4000) # Ok
```

## err\_tensor\_creation\_on\_scratchpad\_with\_init\_value\_not\_supported

Creating SBUF/PSUM tensor with init value is not supported in allocated NKI kernels.

```
t = nl.full((3, par_dim(128), 512), fill_value=1.0, buffer=ncc.sbuf.mod_alloc(base_
    ↪ addr=0)) # t is allocated and has an init value
# Error: Creating SBUF/PSUM tensor with init value is not supported in allocated NKI
    ↪ kernels.
```

## err\_tensor\_output\_not\_written\_to

A tensor was either passed as an output parameter to kernel but never written to, or no output parameter was passed to the kernel at all. At least one output parameter must be provided to kernels.

If you did pass an output parameter to your kernel, and this still occurred, this means the tensor was never written to. The most common cause for this is a dead-loop, such as when a range expression evaluates to 0 and the loop performing the store operation is not actually being entered. But this can occur in any situation in which a loop is never entered, regardless of flow-control construct (for, if, while, etc..)

```
def incorrect(tensor_in, tensor_out):
    M = 128
    N = M + 1

    for i in nl.affine_range( M // N ): # This is the cause of the error, as N > M, M // N
    ↪ will evaluate to 0
        a = nl.load(tensor_in)
        nl.store(tensor_out, value=a) # This store will never be called.

def also_incorrect_in_the_same_way(tensor_in, tensor_out, cnd):
```

(continues on next page)

(continued from previous page)

```
# This will cause the error if the value of `cnd` is False
while cnd:
    a = nl.load(tensor_in)
    nl.store(tensor_out, value=a) # This store will never be called.
```

Consider doing the following:

1. Evaluate your range expressions and conditionals to make sure they're what you intended. If you were trying to perform a computation on tiles smaller than your numerator (M in this case), use `math.ceil()` around your range expression. e.g. `nl.affine_range(math.ceil(M / N))`. You will likely need to pass a mask to your load and store operations as well to account for this.
2. If the possible dead-loop is intentional, you need to issue a store that writes to the entire tensor somewhere in the kernel outside of the dead loop. One good way to do this is to invoke `store()` on your output tensor with a default value.

For example:

```
def memset_output(input, output, cnd):
    # Initialize the output if we cannot guarantee the output are always written later
    nl.store(output[i_p, i_f], value=0)

    while cnd: # Ok even if the value of `cnd` is False
        a = nl.load(tensor_in)
        nl.store(tensor_out, value=a)
```

## err\_transpose\_on\_tensor\_engine\_not\_allowed\_in\_allocated\_kernel

Unsupported transpose case in allocated NKI kernels:

- `nisa.nc_transpose()` with `TensorEngine`, or
- `nl.matmul()` without setting `transpose_x=True`.

User must use their own allocated identity matrix, and call `nisa.nc_matmul()` explicitly to perform transpose on `TensorEngine`.

```
a = .... # assume a has shape [128, 128]
result_a = nl.ndarray((par_dim(128), 128), dtype=nl.bfloat16, buffer=ncc.psum.mod_
↳alloc(byte_addr=0))
result_a[...] = nisa.nc_transpose(a[...]) # Error, calling nc_transpose() with_
↳TensorEngine is not allowed in allocated kernels

b = ... # assume b has shape [32, 32]
result_b = nl.ndarray((par_dim(32), 32), dtype=nl.bfloat16, buffer=ncc.psum.mod_
↳alloc(byte_addr=0))
result_b[...] = nisa.nc_transpose(b[...]) # Error, must specify engine=NeuronEngine.
↳Vector
result_b[...] = nisa.nc_transpose(b[...], engine=NeuronEngine.Vector) # pass
```

## err\_unexpected\_output\_dependencies

Unexpected output dependencies.

NKI assume kernel instances in the spmd grid and iteration between `affine_range` can be executed in parallel require synchronization on the output. As a result, each iteration of the loop will write to a different memory location.

```
a = nl.ndarray((4, 128, 512), dtype=nl.float32, buffer=nl.sbuf)

for i in nl.affine_range(4):
    a[0] = 0 # Unexpected output dependencies, different iterations of i loop write to
    ↪ `a[0]`
```

To fix the problem, you could either index the destination with the missing indices:

```
a = nl.ndarray((4, 128, 512), dtype=nl.float32, buffer=nl.sbuf)

for i in nl.affine_range(4):
    a[i] = 0 # Ok
```

Or if you want to write to the same memory location, you could use `sequential_range` which allows writing to the same memory location:

```
a = nl.ndarray((4, 128, 512), dtype=nl.float32, buffer=nl.sbuf)

for i in nl.sequential_range(4):
    a[0] = 0 # Also ok, we dont expect the sequential_range to execute in parallel
```

## err\_unsupported\_memory

NKI API parameters are in the wrong memory.

NKI enforces API-specific requirements on which memory the parameters are allocated, that is, HBM, SBUF or PSUM. NKI will throw this error when the operands of a NKI API call are not placed in the correct memory.

```
tmp = nl.ndarray((4, 4), dtype=nl.float32, buffer=nl.sbuf)
x = nl.load(tmp) # Error: Expected operand 'src' of 'load' to be in address space 'hbm', but
    ↪ got a tile in 'sbuf' instead.

tmp = nl.ndarray((4, 4), dtype=nl.float32, buffer=nl.hbm)
x = nl.exp(tmp) # Error: Expected operand 'x' of 'exp' to be in address space 'psum/sbuf',
    ↪ but got a tile in 'hbm' instead.
```

## err\_unsupported\_mixing\_basic\_advanced\_tensor\_indexing

Mixing basic tensor indexing and advanced tensor indexing is not supported

```
a = nl.zeros((4, 4), dtype=nl.float32, buffer=nl.sbuf)
i = nl.arange(4)[: , None]
c = nl.exp(a[i, :]) # Error: Mixing basic tensor indexing and advanced tensor indexing
    ↪ is not supported.
```

You could avoid the error by either use basic indexing or advanced indexing but not both:



```
c = nl.exp(a[:, :]) # ok

i = nl.arange(4)[: , None]
j = nl.arange(4)[None, :]
c = nl.exp(a[i, j]) # also ok
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI Developer Guides

*This document is relevant for: Inf2, Trn1, Trn2*

## Getting Started with NKI

In this guide, we will implement a simple “Hello World” style NKI kernel and run it on a NeuronDevice (Trainium/Inferentia2 or beyond device). We will showcase how to invoke a NKI kernel standalone through NKI baremetal mode and also through ML frameworks (PyTorch and JAX). Before diving into kernel implementation, let’s make sure you have the correct environment setup for running NKI kernels.

## Environment Setup

You need a [Trn1](#) or [Inf2](#) instance set up on AWS to run NKI kernels on a NeuronDevice. Once logged into the instance, follow steps below to ensure you have all the required packages installed in your Python environment.

NKI is shipped as part of the Neuron compiler package. To make sure you have the latest compiler package, see [Setup Guide](#) for an installation guide.

You can verify that NKI is available in your compiler installation by running the following command:

```
python -c 'import neuronxcc.nki'
```

This attempts to import the NKI package. It will error out if NKI is not included in your Neuron compiler version or if the Neuron compiler is not installed. The import might take about a minute the first time you run it. Whenever possible, we recommend using local instance NVMe volumes instead of EBS for executable code.

If you intend to run NKI kernels without any ML framework for quick prototyping, you will also need [NumPy](#) installed.

To call NKI kernels from PyTorch, you also need to have `torch_neuronx` installed. For an installation guide, see [PyTorch Neuron Setup](#). You can verify that you have `torch_neuronx` installed by running the following command:

```
python -c 'import torch_neuronx'
```

To call NKI kernels from JAX, you need to have `jax_neuronx` installed. For an installation guide, see [JAX Neuron Setup](#). You can verify that you have `jax_neuronx` installed by running the following command:

```
python -c 'import jax_neuronx'
```

## Implementing your first NKI kernel

In current NKI release, all input tensors must be passed into the kernel as device memory (HBM) tensors on a Neuron-Device. Similarly, output tensors returned from the kernel must also reside in device memory. The body of the kernel typically consists of three main phases:

1. Load the inputs from device memory to on-chip memory (SBUF).
2. Perform the desired computation.
3. Store the outputs from on-chip memory to device memory.

For more details on the above terms, see *NKI Programming Model*.

Below is a small NKI kernel example. In this example, we take two tensors and add them element-wise to produce an output tensor of the same shape.

```

1  from neuronxcc import nki
2  import neuronxcc.nki.language as nl
3
4
5  @nki.jit
6  def nki_tensor_add_kernel(a_input, b_input):
7
8      """NKI kernel to compute element-wise addition of two input tensors
9      """
10
11     # Check all input/output tensor shapes are the same for element-wise operation
12     assert a_input.shape == b_input.shape
13
14     # Check size of the first dimension does not exceed on-chip memory tile size limit,
15     # so that we don't need to tile the input to keep this example simple
16     assert a_input.shape[0] <= nl.tile_size.pmax
17
18     # Load the inputs from device memory to on-chip memory
19     a_tile = nl.load(a_input)
20     b_tile = nl.load(b_input)
21
22     # Specify the computation (in our case: a + b)
23     c_tile = nl.add(a_tile, b_tile)
24
25     # Create a HBM tensor as the kernel output
26     c_output = nl.ndarray(a_input.shape, dtype=a_input.dtype, buffer=nl.shared_hbm)
27
28     # Store the result to c_output from on-chip memory to device memory
29     nl.store(c_output, value=c_tile)
30
31     # Return kernel output as function output
32     return c_output

```

Now let us walk through the above code:

## Importing NKI

We start by importing `neuronxcc.nki` which includes function decorators to compile NKI kernels and also `neuronxcc.nki.language` which implements the NKI language. We will go into more detail regarding the NKI language in *NKI Programming Model*, but for now you can think of it as a tile-level domain-specific language.

```
from neuronxcc import nki
import neuronxcc.nki.language as nl
```

## Defining a kernel

Next we define the `nki_tensor_add_kernel` Python function, which contains the NKI kernel code. The kernel is decorated with `nki.jit`, which allows Neuron compiler to recognize this is NKI kernel code and trace it correctly. Input tensors (`a_input` and `b_input`) are passed by reference into the kernel, just like any other Python function input parameters.

```
@nki.jit
def nki_tensor_add_kernel(a_input, b_input):
```

## Checking input shapes

To keep this getting started guide simple, this kernel example expects all input and output tensors have the same shapes for an element-wise addition operation. We further restrict the first dimension of the input/output tensors to not exceed `nl.tile_size.pmax == 128`. More detailed discussion on tile size limitation is available in *NKI Programming Model*. Note, all of these restrictions *can* be lifted with tensor broadcasting/reshape and tensor tiling with loops in NKI. For more kernel examples, check out *NKI tutorials*.

```
# Check all input/output tensor shapes are the same for element-wise operation
assert a_input.shape == b_input.shape

# Check size of the first dimension does not exceed on-chip memory tile size limit,
# so that we don't need to tile the input to keep this example simple
assert a_input.shape[0] <= nl.tile_size.pmax
```

## Loading inputs

Most NKI kernels start by loading inputs from device memory to on-chip memory. We need to do that because computation can only be performed on data in the on-chip memory.

```
a_tile = nl.load(a_input)
b_tile = nl.load(b_input)
```

## Defining the desired computation

After loading the two input tiles, it is time to define the desired computation. In this case, we perform a simple element-wise addition between two tiles:

```
c_tile = nl.add(a_tile, b_tile)
```

Note that `c_tile = a_tile + b_tile` will also work, as NKI overloads simple Python operators such as `+`, `-`, `*`, and `/`. For a complete set of available NKI APIs, refer to [NKI API Reference Manual](#).

## Storing and returning outputs

To return the output tensor of the kernel, we first declare a NKI tensor `c_output` in device memory (HBM) and then store the output tile `c_tile` from on-chip memory to `c_output` using `nl.store`. We end the kernel execution by returning `c_output` using a standard Python return call. This will allow the host to access the output tensor.

```
# Create a HBM tensor as the kernel output
c_output = nl.ndarray(a_input.shape, dtype=a_input.dtype, buffer=nl.shared_hbm)

# Store the result to c_output from on-chip memory to device memory
nl.store(c_output, value=c_tile)

# Return kernel output as function output
return c_output
```

## Running the kernel

Next, we will cover three unique ways to run the above NKI kernel on a NeuronDevice:

1. NKI baremetal: run NKI kernel with no ML framework involvement
2. PyTorch: run NKI kernel as a PyTorch operator
3. JAX: run NKI kernel as a JAX operator

All three run modes can call the same kernel function decorated with the `nki.jit` decorator as discussed above:

```
1 @nki.jit
2 def nki_tensor_add_kernel(a_input, b_input):
```

The `nki.jit` decorator automatically chooses the correct run mode by checking the incoming tensor type:

1. NumPy arrays as input: run in NKI baremetal mode
2. PyTorch tensors as input: run in PyTorch mode
3. JAX tensors: run in JAX mode

See [nki.jit](#) API doc for more details.

**Note:** NKI baremetal mode is the most convenient way to prototype and optimize performance a NKI kernel alone. For production ML workloads, we highly recommend invoking NKI kernels through a ML framework (PyTorch or JAX). This allows you to integrate NKI kernels in your regular compute graph to accelerate certain operators (see [NKI](#)

*Kernel as a Framework Custom Operator* for details) and leverage the more optimized host-to-device data transfer handling available in ML frameworks.

## NKI baremetal

Baremetal mode expects input tensors of the NKI kernel to be **NumPy arrays**. The kernel also converts its NKI output tensors to **NumPy arrays**. To invoke the kernel, we first initialize the two input tensors *a* and *b* as NumPy arrays. Finally, we call the NKI kernel just like any other Python function:

```
1 import numpy as np
2
3 a = np.ones((4, 3), dtype=np.float16)
4 b = np.ones((4, 3), dtype=np.float16)
5
6 # Run NKI kernel on a NeuronDevice
7 c = nki_tensor_add_kernel(a, b)
8
9 print(c)
```

**Note:** Alternatively, we can decorate the kernel with [nki.baremetal](#) or pass the mode parameter to the `nki.jit` decorator, `@nki.jit(mode='baremetal')`, to bypass the dynamic mode detection. See [nki.baremetal](#) API doc for more available input arguments for the baremetal mode.

## PyTorch

To run the above `nki_tensor_add_kernel` kernel using PyTorch, we initialize the input and output tensors as PyTorch device tensors instead.

```
1 import torch
2 from torch_xla.core import xla_model as xm
3
4 device = xm.xla_device()
5
6 a = torch.ones((4, 3), dtype=torch.float16).to(device=device)
7 b = torch.ones((4, 3), dtype=torch.float16).to(device=device)
8
9 c = nki_tensor_add_kernel(a, b)
10
11 print(c) # an implicit XLA barrier/mark-step (triggers XLA compilation)
```

Running the above code for the first time will trigger compilation of the NKI kernel, which might take a few minutes before printing any output. The printed output should be as follows:

```
tensor([[2., 2., 2.],
        [2., 2., 2.],
        [2., 2., 2.],
        [2., 2., 2.]], device='xla:1', dtype=torch.float16)
```

---

**Note:** Alternatively, we can pass the `mode='torchxla'` parameter into the `nki.jit` decorator to bypass the dynamic mode detection.

---

## JAX

To run the above `nki_tensor_add_kernel` kernel using JAX, we initialize the input tensors as JAX tensors:

```
1 import jax.numpy as jnp
2
3 a = jnp.ones((4, 3), dtype=jnp.float16)
4 b = jnp.ones((4, 3), dtype=jnp.float16)
5
6 c = nki_tensor_add_kernel(a, b)
7
8 print(c)
```

---

**Note:** Alternatively, we can pass the `mode='jax'` parameter into the `nki.jit` decorator to bypass the dynamic mode detection.

---

## Download links

- NKI baremetal script: `getting_started_baremetal.py`
- PyTorch script: `getting_started_torch.py`
- JAX script: `getting_started_jax.py`

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI Programming Model

The NKI programming model enables developers to create custom kernels to program NeuronCores, where every kernel consists of three main stages:

1. **Loading of inputs** from device memory (High Bandwidth Memory, or HBM) to the on-chip SRAM (State Buffer, or SBUF).
2. **Computation definition**, to be executed on the NeuronCore compute engines.
3. **Storing of outputs** from on-chip SRAM (SBUF) back to device memory (HBM).

Fig. 7.2 below is a simplified diagram of a NeuronCore along with its attached HBM device memory. NKI kernels in current release can target a single NeuronCore-v2 or up to two NeuronCore-v3.

As shown in Fig. 7.2, a single NeuronCore consists of two on-chip SRAMs (SBUF and PSUM) and four heterogenous compute engines: the Tensor Engine, Vector Engine, Scalar Engine, and GpSimd Engine. For more information about the compute engine capabilities, see *NeuronDevice Architecture Guide*. Next, let's dive into the memory hierarchy design of NeuronCore, which provides the necessary architecture knowledge to understand the NKI programming model.

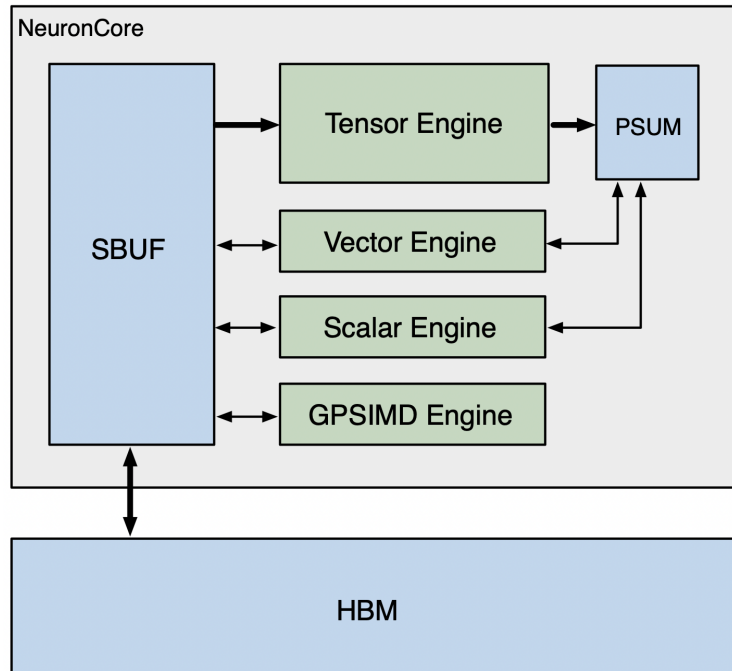


Fig. 7.2: NeuronCore Architecture (multiple NeuronCores available per NeuronDevice)

## Memory hierarchy

Fig. 7.3 below shows the four-level memory hierarchy available to a single NeuronCore. The ranges provided in the figure are intended to calibrate the programmer's mental model. See [NeuronDevice Architecture Guide](#) for the exact values.

Similar to standard memory hierarchy in other devices, memories near the top of the hierarchy are the closest to the compute engines; therefore, they are designed to provide the highest bandwidth and lowest latency. However, the faster memories have smaller capacities compared to memories near the bottom. Unlike memory hierarchy for traditional processors (e.g., CPU, GPU), all the memories available to a NeuronCore are software-managed. They are managed either directly by the programmers or the Neuron SDK. In other words, NeuronCore does not have a hardware cache system to perform any data movement across memories that is opaque to the program. Next, let's discuss the different memories bottom-up.

## NeuronCore external memory

The two memories at the bottom of the hierarchy, host memory and device memory, are both considered *external* memory for a NeuronCore. These memories are **linear memory**, where multi-dimensional tensors must be stored in a flattened manner.

The **host memory** is the CPU-attached DRAM, which is accessible by the host CPUs and all the NeuronCores attached to the instance. NKI kernels currently do not provide APIs to move data in and out of the host memory directly, but we can rely on ML frameworks such as PyTorch or JAX to send input data from host memory into NeuronDevice and vice versa. For an example of this, see [Getting Started with NKI](#).

The **device memory** resides within a NeuronDevice and uses High Bandwidth Memory (HBM) technologies starting from NeuronDevice v2. This means that device memory and HBM refer to the same thing within NKI. Currently, the input and output parameters to NKI kernels must be HBM tensor references. Input tensors in HBM must be loaded into memory within a NeuronCore before any computation can take place.

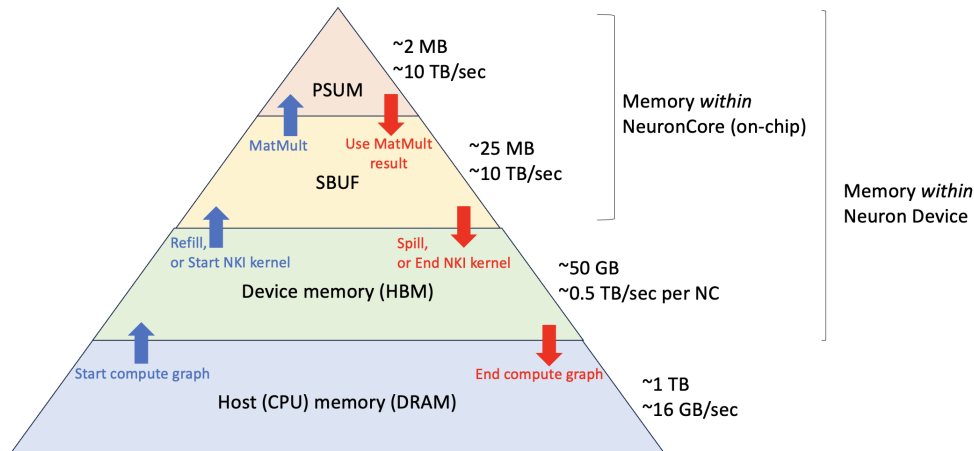


Fig. 7.3: NeuronCore Memory Hierarchy with Capacity and Bandwidth Ranges

### NeuronCore internal memory

The two memories at the top of the hierarchy, SBUF and PSUM, are both considered *internal*, on-chip memory for a NeuronCore. Both memories are **two-dimensional** memory, organized in **128 partitions**. The partitions size of PSUM is typically much smaller than SBUF, and PSUM/SBUF partition sizes vary with NeuronCore generations.

State Buffer (SBUF) memory is the main software-managed on-chip SRAM. The SBUF is accessible by all the compute engines within a NeuronCore. NKI kernel input tensors from HBM must be loaded into the SBUF for computation using `nki.language.load`, and computed output tensors of the kernel must be stored back into the HBM from SBUF using `nki.language.store` before the host can access them. In addition, SBUF is used for storing intermediate data within the kernel, generated by the compute engines. Note, SBUF has **~20x higher bandwidth** than HBM, but needs to be carefully managed to minimize HBM accesses for better performance.

Lastly, Partial Sum Buffer (PSUM) memory is a small, dedicated memory designed for storing matrix multiplication (MatMult) results computed by the tensor engine. Tensor Engine is able to read-add-write to every address in PSUM. Therefore, PSUM is useful for performing large MatMult calculations using multiple tiles where multiple MatMult instructions need to accumulate into the same output tile. As is shown in Fig. 7.2, PSUM memory can also be read and written by the vector and scalar engines. However, due to the limited capacity of PSUM, we recommend that you reserve PSUM space for the tensor engine to write MatMult outputs and to use the vector and scalar engines to evict MatMult results back to SBUF as soon as possible.

Note that to optimize kernel performance, it is a good practice for NKI programmers to be mindful of SBUF and PSUM usage through careful *tiling* and loop fusion. However, ultimately the Neuron compiler performs memory allocation for SBUF and PSUM and assigns memory addresses to kernel intermediate data. When the cumulative size of live data defined by the NKI kernel overflows the capacity of any on-chip memory, the Neuron compiler inserts the necessary spills or refills between that memory and the next-tier memory in the hierarchy.



## Representing data in NKI

NKI represents data in NeuronCore's memory hierarchy with built-in type `Tensor` and its subclasses.

A `Tensor` is a multi-dimensional array which contains elements with the same data type. Programmers can pass `Tensor` in and out of NKI kernels, and declare or initialize `Tensor` in any memory within the `NeuronDevice` (PSUM, SBUF, HBM) using APIs such as `nki.language.ndarray`, `nki.language.zeros`, and `nki.language.full`. Input and output tensors from ML frameworks to NKI kernels can be reinterpreted as NKI `Tensor` of `hbm` buffer type in the same underlying memory buffer.

`Tensor` in NeuronCore's internal memories (SBUF and PSUM) also have a dimension mapped to the partitions of the internal memories. We call this dimension the `partition dimension`. By default, NKI infers the first dimension (that is, the left most dimension) as the `partition dimension` of `Tensor`. Users could also explicitly annotate the `partition dimension` with `par_dim` from `nki.language`. For example:

```
# NKI infers the left most dimension as the partition dimension (size 128 below)
x = nl.ndarray((128, 32, 512), dtype=nl.float32, buffer=nl.sbuf)

# Same as above but more verbose
y = nl.ndarray((nl.par_dim(128), 32, 512), dtype=nl.float32, buffer=nl.sbuf)

# We can also explicitly annotate the partition dimension if we want the partition
# dimension
# to be on the other dimensions. In the following code we are creating a tensor whose
# partition
# dimension is the second from the left most dimension
z = nl.ndarray((128, nl.par_dim(32), 512), dtype=nl.float32, buffer=nl.sbuf)
```

There is a special subclass of `Tensor` called `Index`. `Index` represents the result of the affine expression over variables produced by index-generating APIs, such as loop variables, `nki.language.program_id`, `nki.language.affine_range`, and `nki.language.arange`.

A `Tensor` whose `partition dimension` is the first dimension is also called a `Tile` in NKI. In the above code example, `x` and `y` is a `Tile`, `z` is not a `Tile`. All NKI APIs take `Tile` as input and return a `Tile` as output. We will give more explanation in [Tile-based operations](#).

## Tile-based operations

All NKI APIs operate on `Tile`, which aligns with NeuronCore instruction set architecture (NeuronCore ISA).

```
x = nl.ndarray((128, 32, 512), dtype=nl.float32, buffer=nl.sbuf)
xx = nl.exp(x) # works

z = nl.ndarray((128, nl.par_dim(32), 512), dtype=nl.float32, buffer=nl.sbuf)
zz = nl.exp(z) # not supported
```

To call NKI APIs to process data in a `Tensor` whose `partition dimension` is not the first dimension, users need to generate `Tiles` from the `Tensor`. This can be done by indexing the `Tensor` with a tuple of `Index`, following standard Python syntax `Tensor[Index, Index, ...]`. For example:

```
z = nl.ndarray((128, nl.par_dim(32), 512), dtype=nl.float32, buffer=nl.sbuf)
for i in range(128):
    zz = nl.exp(z[i, :, :]) # works
```

We will provide more discussion of the indexing in [Tensor Indexing](#). Next, let's discuss two important considerations when working with tile-based operations in NKI: [data layout](#) and [tile size](#) constraints.

## Layout considerations

When working with multi-dimensional arrays in any platform, it is important to consider the physical memory layout of the arrays, or how data is stored in memory. For example, in the context of 1D linear memory, we can store a 2D array in a row-major layout or a column-major layout. Row-major layouts place elements within each row in contiguous memory, and column-major layouts place elements within each column in contiguous memory.

As discussed in the [Memory hierarchy](#) section, the on-chip memories, SBUF and PSUM, are arranged as 2D memory arrays. The first dimension is the partition dimension P with 128 memory partitions that can be read and written in parallel by compute engines. The second dimension is the free dimension F where elements are read and written sequentially. A tensor is placed in SBUF and PSUM across both P and F, with the same start offset across all P partitions used by the tensor. [Fig. 7.4](#) below illustrates a default tensor layout. Note that a tile in NKI must map `shape[0]` to the partition dimension.

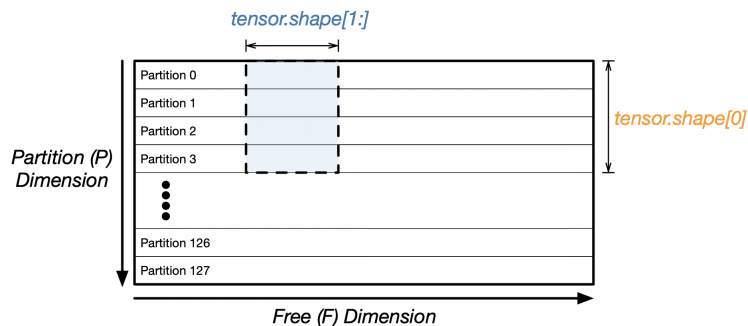


Fig. 7.4: Tensor mapped to partition and free dimensions of SBUF and PSUM

Similar to other domain-specific languages that operate on tensors, NKI defines a *contraction axis* of a tensor as the axis over which reduction is performed, for example the summation axis in a dot product. NKI also defines a *parallel axis* as an axis over which the same operation is performed on all elements. For example, if we take a `[100, 200]` matrix and sum each row independently to get an output of shape `[100, 1]`, then the row-axis (`axis[0]`, left-most) is the parallel axis, and the column-axis (`axis[1]`, right-most) is the contraction axis.

To summarize, the partition and free dimensions of a NKI tensor dictate how the tensor is stored in the 2D on-chip memories physically, while the parallel and contraction axes of a tensor are logical axes that are determined by the computation to be done on the tensor.

The NeuronCore compute engines impose two layout constraints:

- **[LC#1]** For matrix multiplication operations, the contraction axis of both input tiles must be mapped to the P dimension.
- **[LC#2]** For operations that are not matrix multiplication operations, such as scalar or vector operations, the parallel axis should be mapped to the P dimension.

LC#1 means that to perform a matrix multiplication of shapes `[M, K]` and `[K, N]`, Tensor Engine (the engine performing this operation) requires the K dimension to be mapped to the partition dimension in SBUF for both input matrices. Therefore, you need to pass shapes `[K, M]` and `[K, N]` into the `nki.isa.nc_matmul` API, as the partition dimension is always the left-most dimension for an input tile to any NKI compute API.

To help developers get started with NKI quickly, NKI also provides a high-level API `nki.language.matmul` that can take `[M, K]` and `[K, N]` input shapes and invoke the necessary layout shuffling on the input data before sending it to the Tensor Engine `matmul` instruction.

LC#2, on the other hand, is applicable to many instructions supported on Vector, Scalar and GpSimd Engines. See [nki.isa.tensor\\_reduce](#) API as an example.

## Tile size considerations

Besides layout constraints, NeuronCore hardware further imposes three tile-size constraints in NKI:

- [TC#1] The P dimension size of a tile in both SBUF and PSUM must never exceed `nki.tile_size.pmax == 128`.
- [TC#2] For tiles in PSUM, the F dimension size must not exceed `nki.tile_size.psum_fmax == 512`.
- [TC#3] Matrix multiplication input tiles F dimension size must not exceed `nki.tile_size.gemm_stationary_fmax == 128` on the left-hand side (LHS), or `nki.tile_size.gemm_moving_fmax == 512` on the right-hand side (RHS).

You are responsible for breaking your tensors according to these tile-size constraints. If the constraints are not met properly, the NKI kernel compilation throws a `SyntaxError` indicating which constraint is violated. For example, below we show a simple kernel that applies the exponential function to every element of an input tensor. To start, let's write a kernel that expects a hard-coded shape of (128, 512) for both input and output tensors:

```

1  import neuronxcc.nki.language as nl
2  from neuronxcc import nki
3
4  @nki.jit
5  def tensor_exp_kernel(in_tensor):
6      """NKI kernel to compute elementwise exponential of an input tensor
7
8      Args:
9          in_tensor: an input tensor of shape [128,512]
10     Returns:
11         out_tensor: an output tensor of shape [128,512]
12     """
13
14     out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
15                             buffer=nl.shared_hbm)
16
17     # Generate indices for the input/output tensors
18     i_p = nl.arange(128)[: , None]
19     i_f = nl.arange(512)[None, :]
20
21     # Load input data from HBM to on-chip memory
22     in_tile = nl.load(in_tensor[i_p, i_f])
23
24     # perform the computation:
25     out_tile = nl.exp(in_tile)
26
27     # store the results back to HBM
28     nl.store(out_tensor[i_p, i_f], value=out_tile)
29
30     return out_tensor
31
32
33 if __name__ == "__main__":

```

(continues on next page)

(continued from previous page)

```

34 import torch
35 from torch_xla.core import xla_model as xm
36
37 device = xm.xla_device()
38
39 shape = (128, 512)
40 in_tensor = torch.ones(shape, dtype=torch.bfloat16).to(device=device)
41 out_tensor = tensor_exp_kernel_(in_tensor)
42
43 print(out_tensor) # an implicit XLA barrier/mark-step

```

As expected, the output tensor is an element-wise exponentiation of the input-tensor (a tensor of ones):

```

tensor([[2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        ...,
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188]],
        device='xla:1', dtype=torch.bfloat16)

```

Now let's examine what happens if the input/output tensor shapes do not match the shape of the compute kernel. As an example, we can change the input and output tensor shape from [128, 512] to [256, 512]:

```

1  if __name__ == "__main__":
2      import torch
3      from torch_xla.core import xla_model as xm
4
5      device = xm.xla_device()
6
7      shape = (256, 512) # Previously (128, 512)
8      in_tensor = torch.ones(shape, dtype=torch.bfloat16).to(device=device)
9      out_tensor = tensor_exp_kernel_(in_tensor)
10
11     print(out_tensor) # an implicit XLA barrier/mark-step

```

Since the compute kernel is expecting (128, 512) input/output tensors, but we used a (256, 512) input/output tensor instead, the bottom half of the output tensor becomes garbage data:

```

tensor([[2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        [2.7188, 2.7188, 2.7188, ..., 2.7188, 2.7188, 2.7188],
        ...,
        [0.5273, 0.6055, 0.4336, ..., 0.9648, 0.9414, 0.4062],
        [0.7109, 0.2539, 0.7227, ..., 0.7344, 0.2539, 0.1211],
        [0.8867, 0.2109, 0.8789, ..., 0.8477, 0.2227, 0.1406]],
        device='xla:1', dtype=torch.bfloat16)

```

We could try to fix this by changing the tile size inside the compute kernel to (256, 512) as well, and see what happens: (NOTE: This violates tile-size constraint #1!):

```

1 import neuronxcc.nki.language as nl
2 from neuronxcc import nki
3
4
5 @nki.jit
6 def tensor_exp_kernel_(in_tensor):
7     """NKI kernel to compute elementwise exponential of an input tensor
8
9     Args:
10         in_tensor: an input tensor of shape [128,512]
11     Returns:
12         out_tensor: an output tensor of shape [128,512]
13     """
14     out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
15                             buffer=nl.shared_hbm)
16
17     # Generate indices for the input/output tensors
18     i_p = nl.arange(256)[: , None] # Previously nl.arange(128)
19     i_f = nl.arange(512)[None, :]
20
21     # Load input data from HBM to on-chip memory
22     in_tile = nl.load(in_tensor[i_p, i_f])
23
24     # perform the computation:
25     out_tile = nl.exp(in_tile)
26
27     # store the results back to HBM
28     nl.store(out_tensor[i_p, i_f], value=out_tile)
29
30
31 if __name__ == "__main__":
32     import torch
33     from torch_xla.core import xla_model as xm
34
35     device = xm.xla_device()
36
37     shape = (256, 512) # Previously (128, 512)
38     in_tensor = torch.ones(shape, dtype=torch.bfloat16).to(device=device)
39     out_tensor = tensor_exp_kernel_(in_tensor)
40
41     print(out_tensor) # an implicit XLA barrier/mark-step

```

Here, Neuron compiler identifies the tile-size constraint violation and fails compilation with the following exception:

```
SyntaxError: Size of partition dimension 256 exceeds architecture limitation of 128.
```

Now, let's see how NKI developers can build a kernel that properly handles (256, 512) input/output tensors with a simple loop. We can use the `nki.language.tile_size.pmax` constant defined in NKI as the maximum partition dimension size in a tile.

```

1 import neuronxcc.nki.language as nl
2 from torch_neuronx import nki_jit
3

```

(continues on next page)

(continued from previous page)

```

4 @nki_jit
5 def tensor_exp_kernel(in_tensor):
6     """NKI kernel to compute elementwise exponential of an input tensor
7
8     Args:
9         in_tensor: an input tensor of shape [256,512]
10    Returns:
11        out_tensor: an output tensor of shape [256,512]
12    """
13    out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
14                            buffer=nl.shared_hbm)
15
16    i_f = nl.arange(512)[None, :]
17
18    for k in nl.affine_range(2):
19        # Generate tensor indices for the input/output tensors
20        i_p = k * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]
21
22        # Load input data from HBM to on-chip memory
23        in_tile = nl.load(in_tensor[i_p, i_f])
24
25        # perform the computation
26        out_tile = nl.exp(in_tile)
27
28        # store the results back to HBM
29        nl.store(out_tensor[i_p, i_f], value=out_tile)
30
31    return out_tensor

```

The `nl.affine_range(2)` API call returns a list of integers `[0, 1]`. `nl.affine_range` should be the default loop iterator choice in NKI, when the loop has no loop-carried dependency. Note, associative reductions are not considered loop carried dependencies in this context. One such example is accumulating results of multiple matrix multiplication calls into the same output buffer using `+=` (see [Matmul Tutorial](#) for an example). Otherwise, `nl.sequential_range` should be used to handle loop-carried dependency. Note, Neuron compiler transforms any usage of Python `range()` API into `nl.sequential_range()` under the hood. See [NKI iterator API](#) for a detailed discussion of various loop iterator options in NKI.

While the code above does handle (256, 512) tensors correctly, it is rather inflexible since it only supports input shape of (256, 512). Therefore, as a last step, we extend this kernel to handle varying input/output sizes:

```

1 import neuronxcc.nki.language as nl
2 from neuronxcc import nki
3 import math
4
5 @nki_jit
6 def tensor_exp_kernel(in_tensor):
7     """NKI kernel to compute elementwise exponential of an input tensor
8
9     Args:
10        in_tensor: an input tensor of ANY 2D shape (up to SBUF size)
11    Returns:
12        out_tensor: an output tensor of ANY 2D shape (up to SBUF size)
13    """

```

(continues on next page)

(continued from previous page)

```

14  sz_p, sz_f = in_tensor.shape
15  out_tensor = nl.ndarray((sz_p, sz_f), dtype=in_tensor.dtype,
16                        buffer=nl.shared_hbm)
17
18
19  i_f = nl.arange(sz_f)[None, :]
20
21  for p in nl.affine_range(math.ceil(sz_p / nl.tile_size.pmax)):
22      # Generate tensor indices for the input/output tensors
23      # pad index to pmax, for simplicity
24      i_p = p * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]
25
26      # Load input data from external memory to on-chip memory
27      # only read up to sz_p
28      in_tile = nl.load(in_tensor[i_p, i_f], mask=(i_p<sz_p))
29
30      # perform the computation
31      out_tile = nl.exp(in_tile, mask=(i_p<sz_p))
32
33      # store the results back to external memory
34      # only write up to sz_p
35      nl.store(out_tensor[i_p, i_f], value=out_tile, mask=(i_p<sz_p))
36
37  return out_tensor

```

The above example handles cases where `in_tensor.shape[0]` is not a multiple of 128 by passing a `mask` field into the `nl.load` and `nl.store` API calls. For more information, refer to [NKI API Masking](#).

Later in this guide, we'll explore another way to launch a kernel with varying input/output shapes, with a single program multiple data programming model, or [SPMD](#). The SPMD programming model removes the need for explicit looping over different tiles with variable trip counts, which could lead to cleaner and more readable code.

## Tensor Indexing

As mentioned above, we can index Tensor with standard Python syntax to produce Tiles. There are two styles of indexing: Basic and Advanced Tensor Indexing. Note that currently NKI does not support mixing Basic and Advanced Tensor Indexing in the same Index tuple.

### Basic Tensor Indexing

We can index a Tensor with fewer indices than dimensions, we get a *view* of the original tensor as a sub-dimensional tensor. For example:

```

x = nl.ndarray((2, 2, 2), dtype=nl.float32, buffer=nl.hbm)

# `x[1]` return a view of x with shape of [2, 2]
# [[x[1, 0, 0], x[1, 0, 1]], [x[1, 1, 0], x[1, 1, 1]]]
assert x[1].shape == [2, 2]

```

By indexing a Tensor like this, we can generate a Tile with the partition dimension in the first dimension and feed the Tile to NKI compute APIs:

```
# Not a tile, cannot directly feed to a NKI compute API
x = nl.ndarray((2, nl.par_dim(2), 2), dtype=nl.float32)
# Error
y = nl.exp(x)

# `x[1]` have shape [2, 2], and the first dimension is the partition dimension of the
↳ original
# tensor. We can feed it to a NKI compute API.
y = nl.exp(x[1])
```

NKI also supports **slicing** in basic tensor indexing:

```
x = nl.ndarray((2, 128, 1024), dtype=nl.float32, buffer=nl.hbm)

# `x[1, :, :]` is the same as `x[1]`
assert x[1, :, :].shape == [128, 1024]

# Get a smaller view of the third dimension
assert x[1, :, 0:512].shape == [128, 512]

# `x[:, 1, 0:2]` returns a view of x with shape of [2, 2]
# [[x[0, 1, 0], x[0, 1, 1]], [x[1, 1, 0], x[1, 1, 1]]]
assert x[:, 1, 0:2].shape == [2, 2]
```

## Advanced Tensor Indexing

So far we have only shown basic indexing in tensors. However, NeuronCore offers much more flexible tensorized memory access in its on-chip SRAMs along the free dimension. You can use this to efficiently stride the SBUF/PSUM memories at high performance for all NKI APIs that access on-chip memories. However, such flexible indexing is not supported along the partition dimension. That being said, device memory (HBM) is always more performant when accessed sequentially.

In this section, we share several use cases that benefit from advanced memory access patterns and demonstrate how to implement them in NKI.

Advanced Tensor Indexing in NKI leverages the `nl.arange` API.

### Case #1 - Tensor split to even and odd columns

Here we split an input tensor into two output tensors, where the first output tensor gathers all the even columns from the input tensor, and the second output tensor gathers all the odd columns from the input tensor. We assume the rows of the input tensors are mapped to SBUF partitions. Therefore, we are effectively gathering elements along the free dimension of the input tensor. Fig. 7.5 below visualizes the input and output tensors.

```
1 from neuronxcc import nki
2 import neuronxcc.nki.language as nl
3 import math
4
5 @nki.jit
6 def tensor_split_kernel(in_tensor):
7     """NKI kernel to split an input tensor into two output tensors, along the column axis.
```

(continues on next page)



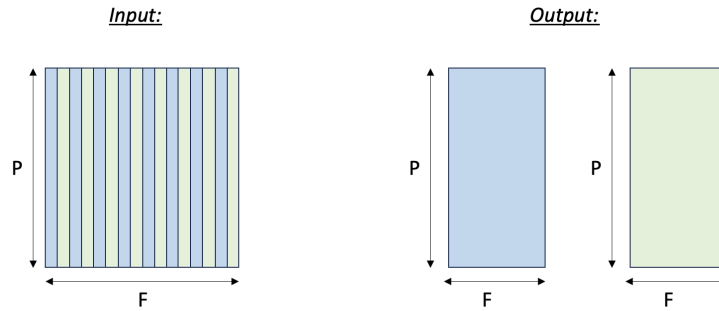


Fig. 7.5: Tensor split to even and odd columns

(continued from previous page)

The even columns of the input tensor will be gathered into the first output tensor, and the odd columns of the input tensor will be gathered into the second output tensor.

Args:

`in_tensor`: an input tensor

Returns:

`out_tensor_even`: a first output tensor (will hold the even columns of the input tensor)

`out_tensor_odd`: a second output tensor (will hold the odd columns of the input tensor)

"""

# Extract tile sizes.

`sz_p, sz_f = in_tensor.shape`

`sz_fout_even = sz_f - sz_f // 2`

`sz_fout_odd = sz_f // 2`

`out_tensor_even = nl.ndarray((sz_p, sz_fout_even), dtype=in_tensor.dtype, buffer=nl.shared_hbm)`

`out_tensor_odd = nl.ndarray((sz_p, sz_fout_odd), dtype=in_tensor.dtype, buffer=nl.shared_hbm)`

# We assume that all three tensors have the same partition dimension size

# and it does not exceed pmax

`assert in_tensor.shape[0] == out_tensor_even.shape[0] == out_tensor_odd.shape[0]`

`assert in_tensor.shape[0] <= nl.tile_size.pmax`

# Make sure even/odd output tensors have correct free dimension size

`assert sz_fout_even == math.ceil(sz_f / 2)`

`assert sz_fout_odd == math.floor(sz_f / 2)`

# Generate tensor indices for the input/output tensors

`i_p = nl.arange(sz_p)[: , None]`

`i_f = nl.arange(sz_f)[None, :]`

`i_fout_even = nl.arange(sz_fout_even)[None, :]`

`i_fout_odd = nl.arange(sz_fout_odd)[None, :]`

# Split pattern:

(continues on next page)

(continued from previous page)

```

42 i_f_even = (2 * i_fout_even)
43 i_f_odd = (2 * i_fout_odd + 1)
44
45 # Load input data from external memory to on-chip memory
46 in_tile = nl.load(in_tensor[i_p, i_f])
47
48 # Perform the split
49 # these assignments invoke copy instructions under the hood
50 # which can execute on either Scalar or Vector Engine
51 # (decided by compiler instruction scheduler)
52 out_tile_even = in_tile[i_p, i_f_even]
53 out_tile_odd = in_tile[i_p, i_f_odd]
54
55 # Store the results back to external memory
56 nl.store(out_tensor_even[i_p, i_fout_even], value=out_tile_even)
57 nl.store(out_tensor_odd[i_p, i_fout_odd], value=out_tile_odd)
58
59 return out_tensor_even, out_tensor_odd
60
61
62 if __name__ == "__main__":
63     import torch
64     from torch_xla.core import xla_model as xm
65
66     device = xm.xla_device()
67
68     X, Y = 4, 5
69     in_tensor = torch.arange(X * Y, dtype=torch.bfloat16).reshape(X, Y).to(device=device)
70
71     out1_tensor, out2_tensor = tensor_split_kernel_(in_tensor)
72     print(in_tensor, out1_tensor, out2_tensor)

```

The main concept in this example is that we introduced the even (`i_f_even`) and odd (`i_f_odd`) indices. Note that both indices are affine expressions of the form `start + stride * nl.arange(size)` with a specific start offset (0/1 respectively) and stride (2 for both cases). This allows us to stride through the `in_tile` memory and copy it to both output tiles (`out_tile_even` and `out_tile_odd`), according to the desired pattern.

## Case #2 - Transpose tensor along the f axis

In this example we transpose a tensor along two of its axes. Note, there are two main types of transposition in NKI:

1. Transpose between the partition-dimension axis and one of the free-dimension axes, which is achieved via the [`nki.isa.nc\_transpose`](#) API.
2. Transpose between two free-dimension axes, which is achieved via a [`nki.language.copy`](#) API, with indexing manipulation in the transposed axes to re-arrange the data.

In this example, we'll focus on the second case: consider a three-dimensional input tensor `[P, F1, F2]`, where the `P` axis is mapped to the different SBUF partitions and the `F1` and `F2` axes are flattened and placed in each partition, with `F1` being the major dimension. Our goal in this example is to transpose the `F1` and `F2` axes with a parallel dimension `P`, which would re-arrange the data within each partition. [Fig. 7.6](#) below illustrates the input and output tensor layouts.

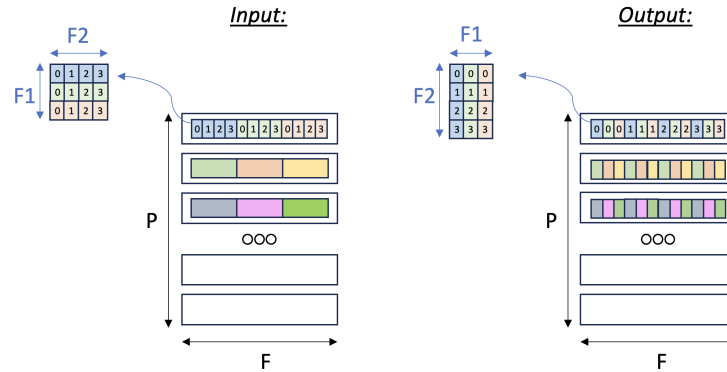


Fig. 7.6: Tensor F1:F2 Transpose

```

1  import neuronxcc.nki as nki
2  import neuronxcc.nki.language as nl
3
4
5  @nki.jit
6  def tensor_transpose2D_kernel_(in_tensor, shape2D):
7      """
8      NKI kernel to reorder the elements on axis[1] of the input tensor.
9
10     Every row of the input tensor is a flattened row-major 2D matrix.
11     The shape2D argument defines the dimensions of the flattened matrices (#rows,#cols).
12     Our goal in this kernel is to transpose these flattened 2D matrices, i.e. make them (
13     ↪ #cols,#rows).
14
15     Example:
16         in_tensor = [a0,a1,a2,a3,b0,b1,b2,b3,c0,c1,c2,c3]
17         shape2D = (3,4)
18     this means that in_tensor has 3 rows and 4 columns, i.e. can be represented as:
19         [a0,a1,a2,a3]
20         [b0,b1,b2,b3]
21         [c0,c1,c2,c3]
22     after transpose, we expect to get:
23         [a0,b0,c0]
24         [a1,b1,c1]
25         [a2,b2,c2]
26         [a3,b3,c3]
27     Thus, out_tensor is expected to be [a0,b0,c0,a1,b1,c1,a2,b2,c2,a3,b3,c3]
28
29     Args:
30         in_tensor: an input tensor
31         shape2D: tuple representing the dimensions to be transposed: (#rows, #cols)
32         out_tensor: an output (transposed) tensor
33     """
34     out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
35                             buffer=nl.shared_hbm)
36     # Gather input shapes
37     sz_p, _ = in_tensor.shape

```

(continues on next page)

(continued from previous page)

```

37
38 # Load input data from external memory to on-chip memory
39 in_tile = nl.load(in_tensor)
40
41 # Performing f1/f2 transpose
42 # =====
43 # The desired transpose pattern is provided as an input:
44 sz_f1, sz_f2 = shape2D
45
46 # We're going to need 3 indices to perform f1:f2 transpose.
47 # - i_p0 is the parallel index
48 # - i_f1 and i_f2 are both free-dim indices, and will be used to transpose between the
49 ↪ f1/f2 axes
49 i_p0 = nl.arange(sz_p)[: , None, None]
50 i_f1 = nl.arange(sz_f1)[None, :, None]
51 i_f2 = nl.arange(sz_f2)[None, None, :]
52
53 # Perform the transposition via a SBUF-to-SBUF copy, with access-pattern manipulation
54 # Note that we have 2D tensors and 3 indices, since we need to represent a 2D access
55 ↪ pattern *per partition*
55 # RHS traverses an F1 x F2 matrix in a row major manner
56 # LHS traverses an F2 x F1 (new) matrix in a row major manner
57 out_tile = nl.ndarray(shape=(sz_p, sz_f2*sz_f1), dtype=out_tensor.dtype)
58 out_tile[i_p0, i_f2*sz_f1+i_f1] = nl.copy(in_tile[i_p0, i_f1*sz_f2+i_f2])
59
60 # Finally, we store out_tile to external memory
61 nl.store(out_tensor, value=out_tile)
62
63 return out_tensor

```

The main concept introduced in this example is a 2D memory access pattern per partition, via additional indices. We copy `in_tile` into `out_tile`, while traversing the memory in different access patterns between the source and destination, thus achieving the desired transposition.

You may download the full runnable script from [Transpose2d tutorial](#).

### Case #3 - 2D pooling operation

Lastly, we examine a case of dimensionality reduction. We implement a 2D MaxPool operation, which is used in many vision neural networks. This operation takes  $C \times [H, W]$  matrices and reduces each matrix along the H and W axes. To leverage free-dimension flexible indexing, we can map the C (parallel) axis to the P dimension and H/W (contraction) axes to the F dimension. Performing such a 2D pooling operation requires a 4D memory access pattern in the F dimension, with reduction along two axes. [Fig. 7.7](#) below illustrates the input and output tensor layouts.

```

1 from neuronxcc import nki
2 import neuronxcc.nki.language as nl
3
4 @nki.jit
5 def tensor_maxpool_kernel(in_tensor, pool_size):
6     """NKI kernel to compute a 2D max-pool operation
7
8     Args:

```

(continues on next page)

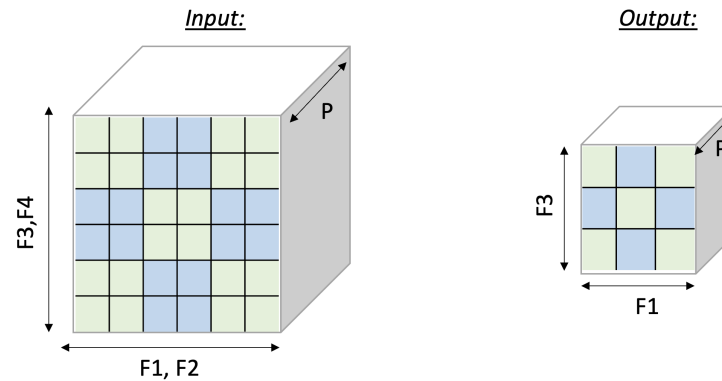


Fig. 7.7: 2D-Pooling Operation (reducing on axes F2 and F4)

(continued from previous page)

```

9     in_tensor: an input tensor, of dimensions C x H x W
10    pool_size: integer P representing a (square) pool-window size
11 Returns:
12     out_tensor: the resulting output tensor, of dimensions C x (H/P) x (W/P)
13 """
14
15 # Get input/output dimensions
16 sz_cin, sz_hin, sz_win = in_tensor.shape
17 sz_hout, sz_wout = sz_hin // pool_size, sz_win // pool_size
18 out_tensor = nl.ndarray((sz_cin, sz_hout, sz_wout), dtype=in_tensor.dtype,
19                        buffer=nl.shared_hbm)
20
21 # Set relevant sizes
22 sz_p = sz_cin
23 sz_pool = pool_size
24
25 # Generate tensor h/w index patterns
26 # 3D indexing according to [C, H, W]
27 i_p = nl.arange(sz_p)[: , None, None] # 3D for
28 i_win = nl.arange(sz_win)[None, None, :]
29 i_hin = nl.arange(sz_hin)[None, :, None]
30
31 i_wout = nl.arange(sz_wout)[None, None, :]
32 i_hout = nl.arange(sz_hout)[None, :, None]
33
34 # Generate pool index patterns (requires two extra dimensions, for the pool window)
35 i_0 = nl.arange(sz_p)[: , None, None, None, None] #
36 i_1 = nl.arange(sz_hin//sz_pool)[None, :, None, None, None] # y_outer
37 i_2 = nl.arange(sz_pool)[None, None, :, None, None] # y_inner
38 i_3 = nl.arange(sz_win//sz_pool)[None, None, None, :, None] # x_outer
39 i_4 = nl.arange(sz_pool)[None, None, None, None, :] # x_inner
40
41 # Load input data from external memory to on-chip memory
42 # Declare ndarray to force a 3D tensor (temporary requirement)
43 in_tile = nl.ndarray([sz_p, sz_hin, sz_win], dtype=in_tensor.dtype)
44 in_tile[:, :, :] = nl.load(in_tensor[i_p, i_hin, i_win])

```

(continues on next page)

(continued from previous page)

```

45
46     # Perform the pooling operation:
47     # We use numpy's advanced indexing, in order to extend in_tile to 5D, and then reduce-
48     ↪max two dimension.
49     # axis[0] is the index for p_dim, and thus doesn't participate in the reduction_
49     ↪operation.
50     # axis[1] and axis[2] together index the rows, with axis[2] responsible for inner_
50     ↪strides
51     # (i.e. inside a pooling window), and axis[1] responsible for the outer strides. As_
51     ↪such, we reduce over axis[2].
52     # Similarly, axis[3] and axis[4] together index the columns, and we thus reduce over_
52     ↪axis[4].
53     out_tile = nl.max(in_tile[i_0, sz_pool*i_1+i_2, sz_pool*i_3+i_4], axis=[2,4])
54
55     # Store the results back to external memory
56     nl.store(out_tensor[i_p, i_hout, i_wout], value=out_tile)
57
58     return out_tensor
59
60 if __name__ == "__main__":
61     import torch
62     from torch_xla.core import xla_model as xm
63
64     device = xm.xla_device()
65
66     # Now let's run the kernel
67     POOL_SIZE = 2
68     C, HIN, WIN = 2, 6, 6
69     HOUT, WOUT = HIN//POOL_SIZE, WIN//POOL_SIZE
70
71     in_tensor = torch.arange(C * HIN * WIN, dtype=torch.bfloat16).reshape(C, HIN, WIN).
71     ↪to(device=device)
72     out_tensor = tensor_maxpool_kernel_(in_tensor, POOL_SIZE)
73
74     print(in_tensor, out_tensor) # an implicit XLA barrier/mark-step

```

## SPMD: Launching multiple instances of a kernel

So far we have discussed how to launch a single NKI kernel instance, in which the full input tensor is processed. In this section, we discuss how to launch multiple instances of the same kernel and slice the full input tensor across kernel instances using a single program multiple data programming model (SPMD).

**Note:** In current NKI release, adopting the SPMD programming model has **no** impact on performance of NKI kernel, and therefore is considered **optional**. A SPMD program is compiled into an executable that targets one NeuronCore, and the different instances of the SPMD program are executed serially on a single NeuronCore. This is subject to changes in future releases.

NKI allows users to launch multiple instances of a kernel, which are organized in a user-defined multi-dimensional grid. The grid indices are then used by the different kernel instances to select which input and output data to access. There is

no restriction on the number of dimensions in an SPMD grid, nor on the size of each dimension. Each kernel instance can find its coordinates within the launch grid using the `nki.language.program_id` API. Neuron compiler translates the SPMD launch grid into nested loops of compute-kernel invocations, which are then executed on the NeuronCore.

As an example, we'll perform a  $C=A @ B$  matrix multiplication, where A and B are of shape (512, 128) and (128, 1024) respectively. We partition the output tensor C of shape (512, 1024) into 4x2 tiles and assign the task of computing each output tile to a different kernel instance. A 4x2 launch-grid is chosen in this case, in order to make each compute kernel instance operate on a single tile in A and a single tile in B, while adhering to the *tile-size constraints*.

With a 2D 4x2 launch grid, the (i, j) kernel instance is responsible for computing the (i, j) tile of C. The computation of the (i, j) tile requires the corresponding rows of A and columns of B. This induces a four-way row-wise partitioning of A and a two-way column-wise partitioning of B, as shown in Fig. 7.8.

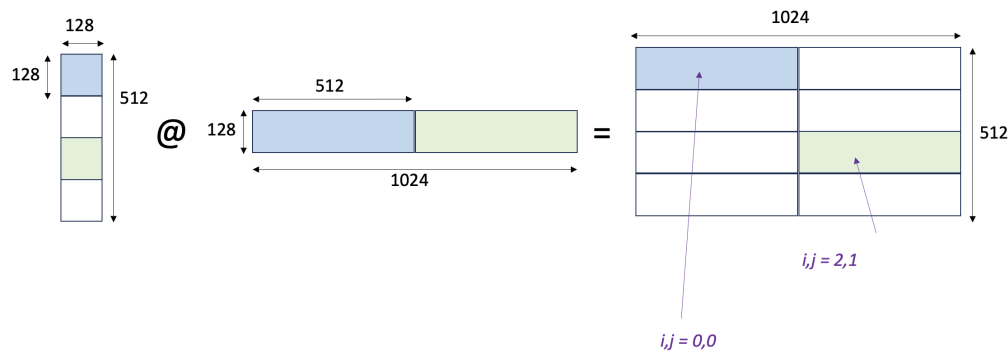


Fig. 7.8: Visualization of 512x128x1024 matrix multiplication using SPMD

In this SPMD kernel example, we will use the high-level `nki.language.matmul` API, so that we can focus on the concept of SPMD without worrying about the layout requirement of Tensor Engine (*LC#1*). To achieve the best performance, we suggest transposing input A and invoking another NKI kernel instead, which solely performs matmul operations on Tensor Engine using `nki.isa.nc_matmul` without extra overhead in changing input layouts to meet *LC#1*.

```

1  import neuronxcc.nki.language as nl
2  from neuronxcc import nki
3
4
5  @nki.jit
6  def matmul_128x128x512_spmd(A, B):
7      """NKI kernel to compute a 128x128x512 matrix multiplication operation.
8          Use SPMD program IDs to index into the full A and B input tensor to get tiles
9          for 128x128x512 matrix multiplication.
10
11      Args:
12          A: an input tensor of shape [M=512,K=128],
13             a left hand side argument of the matrix multiplication,
14          B: an input tensor of shape [K=128,N=1024],
15             a right hand side argument of the matrix multiplication
16          result: the resulting output tensor of shape [M=512,N=1024]
17      """
18      N, K = A.shape
19      K_, M = B.shape
20      assert K == K_
21      # Create output tensor shared between all SPMD instances as result tensor
22      result = nl.ndarray((N, M), dtype=A.dtype, buffer=nl.shared_hbm)

```

(continues on next page)

(continued from previous page)

```

23
24 # Defining starting indexes for input A and B
25 i_A_row = nl.program_id(0) * 128
26 i_B_col = nl.program_id(1) * 512
27
28 # Loading the inputs (HBM->SBUF)
29 A_tile = nl.load(A[i_A_row:i_A_row+128, 0:128])
30 B_tile = nl.load(B[0:128, i_B_col:i_B_col+512])
31
32 # Perform the matrix-multiplication
33 # Note1: nl.matmul will invoke a transpose on A_tile before performing the actual
↪matmul operation
34 # Note2: A NKI matmul instruction always writes to PSUM in float32 data-type
35 result_psum = nl.matmul(A_tile, B_tile)
36
37 # Copy the result from PSUM back to SBUF, and cast to expected output data-type
38 result_sbuf = nl.copy(result_psum, dtype=result.dtype)
39
40 # The result of a [128,128] x [128,512] matrix multiplication has a shape of [128,
↪512].
41 # This dictates which indices to use to address the result tile.
42 nl.store(result[i_A_row:i_A_row+128, i_B_col:i_B_col+512], value=result_sbuf)
43
44 return result
45
46 if __name__ == "__main__":
47     from torch_xla.core import xla_model as xm
48     import torch
49
50     device = xm.xla_device()
51
52     A = torch.ones((512, 128), dtype=torch.bfloat16).to(device=device)
53     B = torch.ones((128, 1024), dtype=torch.bfloat16).to(device=device)
54
55     # Launch kernel with a 2D grid
56     result = matmul_128x128x512_spmd[4, 2](A, B)
57
58     print(result) # an implicit XLA barrier/mark-step

```

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2



## NKI Kernel as a Framework Custom Operator

This document demonstrates how to insert a NKI kernel as a custom operator into a PyTorch or JAX model using simple code examples.

### Using NKI kernels

To register a NKI kernel registration, you need to call a decorated NKI function.

Let's examine a guiding example below where we randomly initialize two inputs, add them together, and then multiply the result by the two input tensors element-wise. This effectively calculates:  $a * b * (a + b)$ .

We define a common NKI kernel for addition. For more information on the kernel, see [SPMD Tensor Addition](#).

```

1  import neuronxcc.nki as nki
2  import neuronxcc.nki.language as nl
3
4
5  @nki.jit
6  def nki_tensor_add_kernel_(a_input, b_input):
7      """NKI kernel to compute element-wise addition of two input tensors
8
9      This kernel assumes strict input/output sizes can be uniformly tiled to [128,512]
10
11     Args:
12         a_input: a first input tensor
13         b_input: a second input tensor
14
15     Returns:
16         c_output: an output tensor
17     """
18     # Create output tensor shared between all SPMD instances as result tensor
19     c_output = nl.ndarray(a_input.shape, dtype=a_input.dtype, buffer=nl.shared_hbm)
20
21     # Calculate tile offsets based on current 'program'
22     offset_i_x = nl.program_id(0) * 128
23     offset_i_y = nl.program_id(1) * 512
24
25     # Generate tensor indices to index tensors a and b
26     ix = offset_i_x + nl.arange(128)[: , None]
27     iy = offset_i_y + nl.arange(512)[None, :]
28
29     # Load input data from device memory (HBM) to on-chip memory (SBUF)
30     # We refer to an indexed portion of a tensor as an intermediate tensor
31     a_tile = nl.load(a_input[ix, iy])
32     b_tile = nl.load(b_input[ix, iy])
33
34     # compute a + b
35     c_tile = a_tile + b_tile
36
37     # store the addition results back to device memory (c_output)
38     nl.store(c_output[ix, iy], value=c_tile)
39

```

(continues on next page)

(continued from previous page)

```

40 # Transfer the ownership of `c_output` to the caller
41 return c_output

```

## PyTorch

We can perform  $(a + b) * a * b$  using native PyTorch code.

```

import torch
from torch_xla.core import xla_model as xm

device = xm.xla_device()

a = torch.randn(256, 1024, dtype=torch.float32).to(device)
b = torch.randn(256, 1024, dtype=torch.float32).to(device)
c = a + b
out = a * b * c

print(out)

```

Now let's replace the tensor addition ( $c = a + b$ ) with a NKI kernel. To do this we replace the  $+$  operator with a call to the NKI kernel caller (`nki_tensor_add`), and everything else works as before.

```

1 def nki_tensor_add(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the
5     ↪ inputs/outputs
6
7     Args:
8         a_input: a first input tensor, of shape [N*128, M*512]
9         b_input: a second input tensor, of shape [N*128, M*512]
10
11     Returns:
12         a tensor of shape [N*128, M*512], the result of a_input + b_input
13     """
14     # The SPMD launch grid denotes the number of kernel instances.
15     # In this case, we use a 2D grid where the size of each invocation is 128x512
16     grid_x = a_input.shape[0] // 128
17     grid_y = a_input.shape[1] // 512
18
19     return nki_tensor_add_kernel_[grid_x, grid_y](a_input, b_input)

```

```

device = xm.xla_device()
a = torch.randn(256, 1024, dtype=torch.float32).to(device)
b = torch.randn(256, 1024, dtype=torch.float32).to(device)
c = nki_tensor_add(a, b) # calling a NKI kernel, instead of the built-in torch op
out = a * b * c
print(out)

```

To understand what happens under the hood when we compile the above code, we can print HLO IR graph generated

by XLA by setting the `NEURON_FRAMEWORK_DEBUG` environment variable. For example, you may add the following lines to your code:

```
import os
os.environ['NEURON_FRAMEWORK_DEBUG'] = "1"
```

A `.pbtxt` file is then written in your run directory that has the corresponding human-readable HLO IR.

Let's examine the XLA output of this example. In line #5 we can identify that the tensor addition is now mapped to an HLO custom-call instruction, with `AwsNeuronCustomNativeKernel` as `custom_call_target`. The output of that custom-call is then consumed by the next instruction in line #6 as usual.

```
1 ENTRY %SyncTensorsGraph.22 (p0.2: f32[256,1024], p1.2: f32[256,1024]) -> (f32[256,1024])
  ↳ {
2   %p1.2 = f32[256,1024]{1,0} parameter(1), frontend_attributes={neff_input_name="input1"}
3   %p0.2 = f32[256,1024]{1,0} parameter(0), frontend_attributes={neff_input_name="input0"}
4   %multiply = f32[256,1024]{1,0} multiply(f32[256,1024]{1,0} %p1.2, f32[256,1024]{1,0}
  ↳ %p0.2)
5   %custom-call.2 = f32[256,1024]{1,0} custom-call(f32[256,1024]{1,0} %p1.2, f32[256,1024]
  ↳ {1,0} %p0.2), custom_call_target="AwsNeuronCustomNativeKernel", api_version=API_
  ↳ VERSION_UNSPECIFIED, backend_config="...")
6   %multiply.1 = f32[256,1024]{1,0} multiply(f32[256,1024]{1,0} %multiply, f32[256,1024]{1,
  ↳ 0} %custom-call.2)
7   ROOT %tuple = (f32[256,1024]{1,0}) tuple(f32[256,1024]{1,0} %multiply.1), frontend_
  ↳ attributes={neff_output_names="output0"}
8 }
```

The Neuron compiler replaces the above custom-call with the corresponding NKI kernel implementation while optimizing the rest of the compute graph as usual. At the end of the compilation process, a single compiled binary NEFF file is generated representing the entire graph including the NKI kernel. For more information about NEFF files, see [Neuron Compiler](#).

## JAX

We can perform  $(a + b) * a * b$  using native JAX code.

```
import jax
import jax.numpy as jnp

@jax.jit
def jax_customop_tutorial(a, b):
    c = a + b
    out = a * b * c
    return out

seed = jax.random.PRNGKey(0)
seed_a, seed_b = jax.random.split(seed)
a = jax.random.normal(seed_a, (256, 1024), dtype=jnp.float32)
b = jax.random.normal(seed_b, (256, 1024), dtype=jnp.float32)

print(jax_customop_tutorial(a, b))
```

Similar to the PyTorch example above, let's replace the tensor addition ( $c = a + b$ ) with the addition NKI kernel. To do this we replace the `+` operator with a call to the NKI kernel caller (`nki_tensor_add`), and everything else works

as before.

```

1 def nki_tensor_add(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the_
    ↪ inputs/outputs
5
6     Args:
7         a_input: a first input tensor, of shape [N*128, M*512]
8         b_input: a second input tensor, of shape [N*128, M*512]
9
10    Returns:
11        a tensor of shape [N*128, M*512], the result of a_input + b_input
12    """
13
14    # The SPMD launch grid denotes the number of kernel instances.
15    # In this case, we use a 2D grid where the size of each invocation is 128x512
16    grid_x = a_input.shape[0] // 128
17    grid_y = a_input.shape[1] // 512
18
19    return nki_tensor_add_kernel_[grid_x, grid_y](a_input, b_input)

```

```

import jax
import jax.numpy as jnp

@jax.jit
def jax_customop_tutorial(a, b):
    c = nki_tensor_add(a, b) # calling a NKI kernel, instead of the built-in jax op
    out = a * b * c
    return out

seed = jax.random.PRNGKey(0)
seed_a, seed_b = jax.random.split(seed)
a = jax.random.normal(seed_a, (256, 1024), dtype=jnp.float32)
b = jax.random.normal(seed_b, (256, 1024), dtype=jnp.float32)
print(jax_customop_tutorial(a, b))

```

To understand what happens under the hood when we compile the above code, we can print the HLO IR graph by adding the following snippet to your code:

```

print(jax.jit(jax_customop_tutorial)
      .lower(a, b)
      .compile()
      .runtime_executable()
      .hlo_modules()[0].to_string()
)

```

Let's examine the XLA output of this example. In line #7 we can identify that the tensor addition is now mapped to an HLO custom-call instruction, similar to PyTorch. The output of that custom-call is then consumed by the next instruction in line #8 as usual.

```

1 HloModule jit_add, entry_computation_layout= {(f32[256,1024]{1,0}, f32[256,1024]{1,0})->
    ↪ (f32[256,1024]{1,0})}, allow_spmd_sharding_propagation_to_output={true}

```

(continues on next page)

(continued from previous page)

```

2 ENTRY %main.11 (Arg_0.1: f32[256,1024], Arg_1.2: f32[256,1024]) -> (f32[256,1024]) {
3   %Arg_0.1 = f32[256,1024]{1,0} parameter(0), sharding={replicated}
4   %Arg_1.2 = f32[256,1024]{1,0} parameter(1), sharding={replicated}
5   %multiply.0 = f32[256,1024]{1,0} multiply(f32[256,1024]{1,0} %Arg_0.1, f32[256,1024]{1,
6   ↪ 0} %Arg_1.2), metadata={op_name="jit(add)/jit(main)/jit(jax_customop_tutorial)/mul"
7   ↪ source_file="/tmp/ipykernel_3935360/2333914945.py" source_line=61}
8   %custom-call.0 = f32[256,1024]{1,0} custom-call(f32[256,1024]{1,0} %Arg_0.1, f32[256,
9   ↪ 1024]{1,0} %Arg_1.2), custom_call_target="AwsNeuronCustomNativeKernel", api_
10  ↪ version=API_VERSION_STATUS_RETURNING, metadata={op_name="jit(add)/jit(main)/jit(jax_
11  ↪ customop_tutorial)/nki_call[func=<function nki_tensor_add_kernel_ at 0xf6be28f6f80>
12  ↪ grid=(2, 2) out_shape=(ShapeDtypeStruct(shape=(256, 1024), dtype=float32),)]" source_
13  ↪ file="/home/ubuntu/nki/src/jax_neuronx/core.py" source_line=34}, backend_config="..."
14  %multiply.1 = f32[256,1024]{1,0} multiply(f32[256,1024]{1,0} %multiply.0, f32[256,1024]
15  ↪ {1,0} %custom-call.0), metadata={op_name="jit(add)/jit(main)/jit(jax_customop_
16  ↪ tutorial)/mul" source_file="/tmp/ipykernel_3935360/2333914945.py" source_line=61}
17  ROOT %tuple.10 = (f32[256,1024]{1,0}) tuple(f32[256,1024]{1,0} %multiply.1)
18 }

```

The Neuron compiler replaces the above custom-call with the corresponding NKI kernel implementation while optimizing the rest of the compute graph as usual. At the end of the compilation process, a single compiled binary NEFF file is generated representing the entire graph including the NKI kernel. For more information about NEFF files, see [Neuron Compiler](#).

## Using NKI in training graphs

If you are using NKI to implement a new operator in a training graph, you might need to make the new operator interplay with the autograd engine in the framework. To do this, in PyTorch, you can subclass the framework's base operator class and implement both the `forward()` and `backward()` methods. The autograd engine then uses the `backward()` method when performing auto-differentiation. See [Extending torch.autograd](#) in the PyTorch Docs for instructions on doing this in PyTorch. To do this in JAX, you can create a `custom_vjp` rule (vjp stands for Vector-Jacobian product), which binds the `forward()` and `backward()` calls. See [Autodiff Cookbook](#) in the JAX Docs for instructions on doing this.

Let's reuse the `nki_tensor_add` kernels from before and demonstrate how to train a simple compute graph  $(a+b) * a * b$  in both PyTorch and JAX.

## PyTorch

We define a `NkiAddFunc` class, which leverages the `nki_tensor_add` kernel in its `forward()` function. The gradients of both input tensors in  $y = a + b$  are ones, so the `backward()` function propagates the  $dy$  gradients from the previous backward function.

```

import torch
import torch_xla.core.xla_model as xm
device = xm.xla_device()

class NkiAddFunc(torch.autograd.Function):
    @staticmethod

```

(continues on next page)

(continued from previous page)

```

def forward(ctx, a, b):
    return nki_tensor_add(a, b)

@staticmethod
def backward(ctx, dy, *args):
    # gradients for a and b
    return dy, dy

# now, let's define the compute graph
a = torch.randn(256, 1024, dtype=torch.float32).to(device).detach().requires_grad_()
b = torch.randn(256, 1024, dtype=torch.float32).to(device).detach().requires_grad_()
c = NkiAddFunc.apply(a, b)
out = a * b * c

# here we define a (dummy) loss-function, in prep for backward propagation
loss = out.sum()

# lastly, let's invoke the auto-grad engine
loss.backward()

xm.mark_step()

```

## JAX

We define a custom\_vjp function `nki_add_func` by using the `@jax.custom_vjp` decorator which directly calls the `nki_tensor_add` kernel. We then define and register the `forward()` and `backward()` implementations of the `nki_add_func` function via `defvjp()`. Just like the PyTorch example before, the `backward()` implementation simply passes the gradients through. Finally, to start training, we execute the forward pass by calling `nki_add_func(a, b) * x * y`. To get the gradients, we call `jax.grad` directly with a loss function.

```

@jax.custom_vjp
def nki_add_func(a, b):
    return nki_tensor_add(a, b)

def f_forward(a, b):
    # operator output and residual (same as input here)
    return nki_add_func(a, b), (a, b)

def f_backward(res, grad):
    # gradients for a and b
    return grad, grad

nki_add_func.defvjp(f_forward, f_backward) # line 11

@jax.jit
def jax_customop_tutorial_and_grad(a, b):
    out = nki_add_func(a, b) * x * y

    # use the same dummy loss function (output sum) as PyTorch example above
    grad = jax.grad(lambda x, y: (nki_add_func(x, y) * x * y).sum(), argnums=(0, 1))(a, b)
    return out, *grad

```

(continues on next page)

(continued from previous page)

```
c, grad_a, grad_b = jax_customop_tutorial_and_grad(a, b)
```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## NeuronDevice Architecture Guide for NKI

NKI currently supports the following NeuronDevice generations:

- Trainium/Inferentia2, available on AWS `trn1`, `trn1n` and `inf2` instances
- Trainium2, available on AWS `trn2` instances and UltraServers

The documents below provide an architecture deep dive of each NeuronDevice generation, with a focus on areas that NKI developers can directly control through kernel implementation. [Trainium/Inferentia2 Architecture Guide](#) serves as a foundational architecture guide for understanding basics of any NeuronDevice generation, while [Trainium2 Architecture Guide](#) walks through architecture enhancements compared to the previous generation in details. Therefore, we suggest new NKI developers start with [Trainium/Inferentia2 Architecture Guide](#) before exploring newer NeuronDevice architecture.

Trainium/Inferentia2 Architecture Guide

Trainium2 Architecture Guide

*This document is relevant for:* Inf2, Trn1, Trn2

## Trainium/Inferentia2 Architecture Guide for NKI

In this guide, we will dive into hardware architecture of second-generation NeuronDevices: Trainium/Inferentia2. Our goal is to equip advanced Neuron users with sufficient architectural knowledge to write performant NKI kernels and troubleshoot performance issues on NeuronDevices using [neuron-profile](#), a profiler tool designed specifically for NeuronDevices. This guide is also written assuming readers have read through [NKI Programming Model](#) and familiarized themselves with key NKI concepts.

[Fig. 7.9](#) shows a block diagram of a Trainium and Inferentia2 device. At a high level, both Trainium and Inferentia2 devices consist of:

- 2 NeuronCores (v2).
- 2 HBM stacks with a total device memory capacity of 32GiB and bandwidth of 820 GB/s.
- 32 DMA (Direct Memory Access) engines to move data within and across devices.
- 6 CC-Cores for collective communication.
- 2 (Inferentia2) or 4 (Trainium) NeuronLink-v2 for device-to-device collective communication.

The rest of this guide will go into details of each compute engine in NeuronCore-v2 and supported data movement patterns across the memory hierarchy.

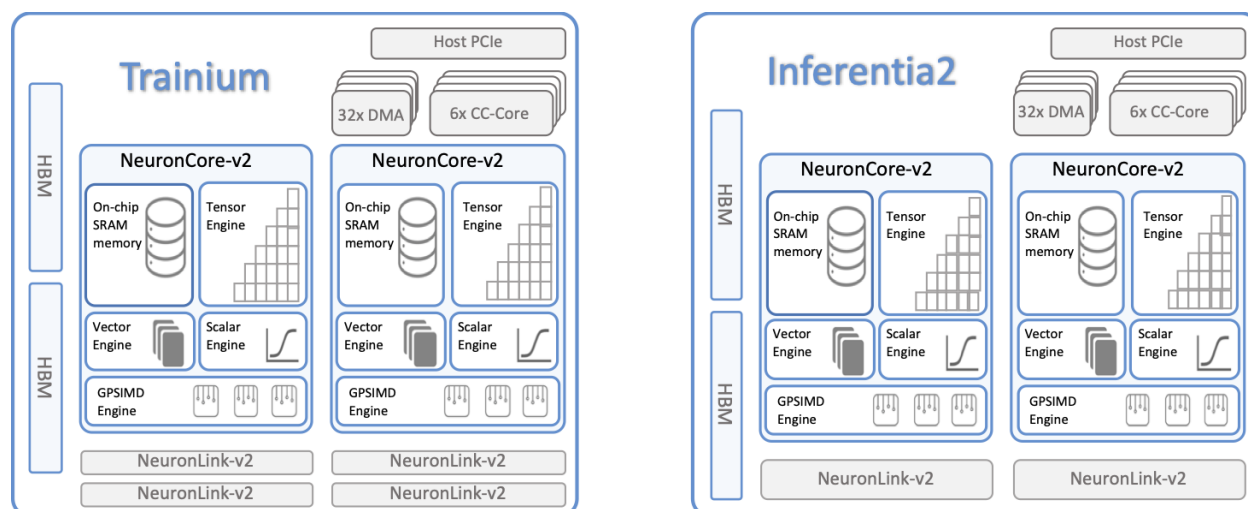


Fig. 7.9: Trainium/Inferentia2 Device Diagrams.

## NeuronCore-v2 Compute Engines

In this section, we will describe the architectural details within a NeuronCore-v2. The figure below is a simplified diagram of the compute engines and their connectivity to the two on-chip SRAMs: state buffer (SBUF) and partial sum buffer (PSUM).

A NeuronCore-v2 consists of four heterogeneous compute engines (Tensor, Vector, Scalar and GpSimd), each of which is designed to accelerate different types of operators in modern machine learning models. These engines execute their own instruction sequences *asynchronously* in parallel, but they can perform explicit synchronization to meet data and resource dependency requirements through atomic semaphores in hardware. In NKI, programmers are not required to program such engine synchronization manually. If the synchronization is not explicitly specified, the Neuron Compiler will insert the required synchronizations during compilation, based on data dependencies identified in the NKI kernel. NKI API calls without data dependencies can run in parallel if they have different target engines.

In addition, it is often useful to take engine data-path width and frequency into account when optimizing performance for a multi-engine operator:

| Device Architecture  | Compute Engine | Data-path Width (elements/cycle) | Frequency (GHz) |
|----------------------|----------------|----------------------------------|-----------------|
| Trainium/Inferentia2 | Tensor         | 2x128 (input); 1x128 (output)    | 2.8             |
|                      | Vector         | 128 input/output                 | 1.12            |
|                      | Scalar         |                                  | 1.4             |
|                      | GpSimd         |                                  | 1.4             |

Memory-wise, a NeuronCore-v2 consists of two software-managed on-chip SRAMs, a 24MiB SBUF as the main data storage and a 2MiB PSUM as a dedicated accumulation buffer for Tensor Engine. Both SBUF and PSUM are considered two-dimensional memories with 128 partitions each, i.e., one SBUF partitions has 192KiB of memory while one PSUM partition has 16KiB. We will cover more details on data movements with SBUF/PSUM later [here](#).

The rest of this section will cover the following topics for each compute engine:

- Key functionalities.
- Layout and tile size requirement for input and output tensors.
- Best practices to achieve good performance on the engine.



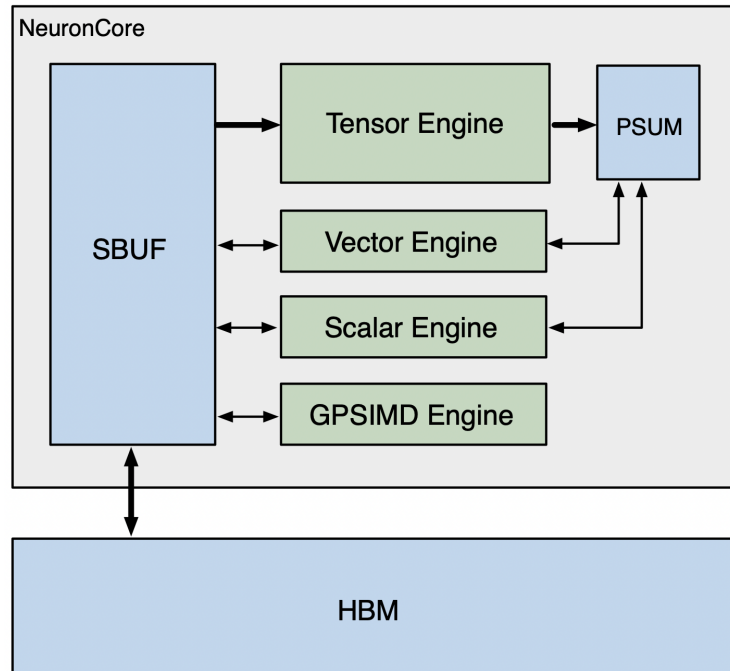


Fig. 7.10: NeuronCore-v2 and its device memory (HBM).

## Tensor Engine

Tensor Engine (TensorE from now on) is specially designed to accelerate matrix-multiplications (matmuls), as well as other operators that can be executed using matrix multiplications such as 2D convolutions. We also note that TensorE can be used for advanced data movement from SBUF to PSUM, including transposition and broadcast (more discussion below [here](#)). Architecturally, the engine is built around a [systolic array](#) with 128 rows and 128 columns of processing elements, which streams input data from SBUF and writes output to PSUM.

**Data Types.** TensorE supports [BF16](#), FP16, [TF32](#), and cFP8 input matrix data types at a maximum throughput of 92 TFLOPS, as well as 23 TFLOPS for FP32 inputs. TensorE performs mixed-precision calculations, with accumulations at FP32 precision. Therefore, the output data of a TensorE calculation is always in FP32.

**Layout.** To understand the layout and tiling constraints of TensorE, let's visualize its connection to SBUF and PSUM as below. Note, PSUM partition dimension is purposely rotated 90 degrees compared to SBUF partition dimension due to systolic array data flow.

As shown in the diagram above, TensorE must **read** input matrices from **SBUF** and **write** output matrices to **PSUM**. PSUM also allows near-memory accumulation of multiple matrix multiplication output tiles (detailed usage discussed [here](#)).

In NKI, to perform a multiplication of two matrices,  $x[M, K]$  and  $y[K, N]$ , you may invoke the NKI language API `nki.language.matmul(x, y)` directly. The returned tile has a shape of  $[M, N]$  as expected. At the hardware level, TensorE requires both input tiles to have the **contraction dimension K** in the SBUF partition dimension, that is, the first dimension of input shapes (LC #1 as discussed in [NKI Programming Model](#)). This ISA requirement is reflected in the low-level API `nki.isa.nc_matmul`, which takes `stationary` and `moving` matrices as input parameters. Therefore, `nki.language.matmul(x, y)` is a two-step computation: invoking `nki.isa.nc_transpose(x)` to get `stationary` and then `nki.isa.nc_matmul(stationary, moving)` to get the final result. In other words, `nki.isa.nc_matmul(stationary[K,M], moving[K,N])` performs a `stationary.T @ moving` calculation, which will result in an output with dimensions  $[M, N]$ .

For every `nki.isa.nc_matmul(stationary, moving)` call, TensorE executes two distinct Neuron ISA instructions:

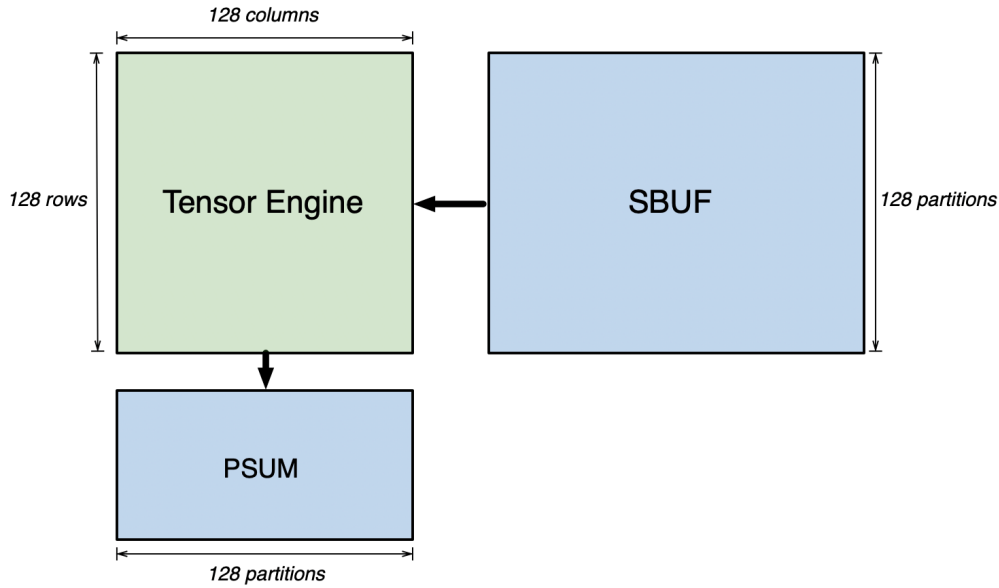


Fig. 7.11: Tensor Engine and SRAM Connectivity.

- **LoadStationary** (short for LS): This instruction loads the **stationary** from SBUF and caches it in internal storage of TensorE
- **MultiplyMoving** (short for MM): This instruction loads the **moving** from SBUF and multiplies **moving** across the pre-loaded **stationary** matrix from the previous LoadStationary instruction. The output of this instruction is the output of the `nki.isa.nc_matmul` call written to PSUM.

With the above instruction sequence, we as NKI programmers effectively map input tile **stationary** as the stationary tensor and input tile **moving** as the moving tensor for TensorE. As a rule-of-thumb for layout analysis, the **free** axis of the **stationary** tensor always becomes the partition (first) axis of the output tile, while the **free** axis of the **moving** tensor becomes the free axis of the output. Fig 7.12 below visualizes this concept by showing a matrix multiplication in both mathematical and TensorE views.

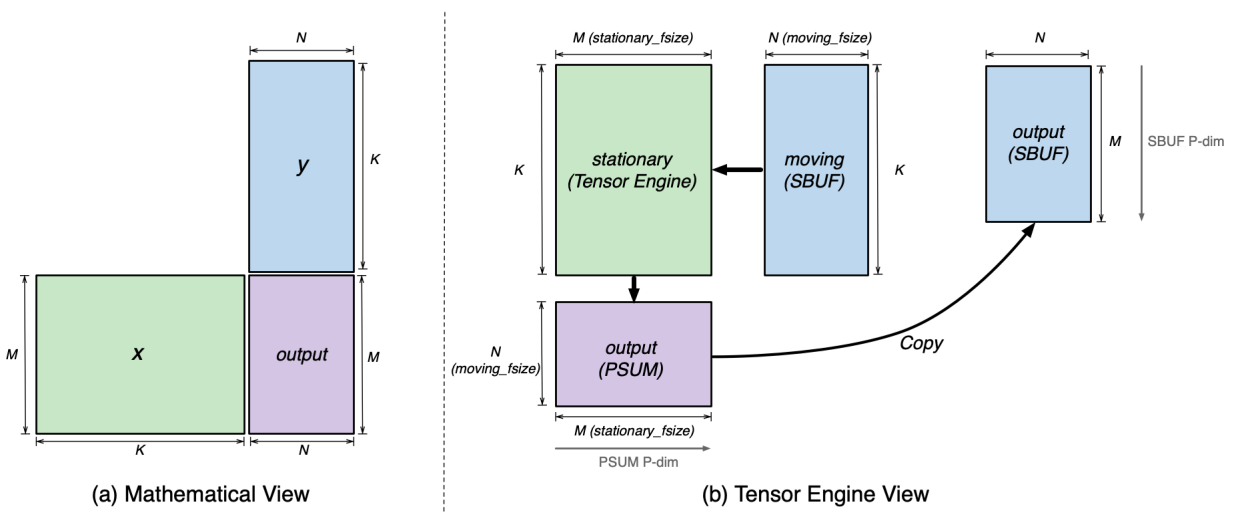


Fig. 7.12: MxKxN Matrix Multiplication Visualization.

However, programmers are also free to map stationary tile to the moving tensor instead, which would lead to the same output tile but transposed: `nki.isa.nc_matmul(moving[K,N], stationary[K,M]) = moving.T @ stationary = outputT[N, M]`. In fact, mapping high-level input tiles to the low-level stationary/moving tensors in TensorE is an important layout decision that NKI programmers should consider to minimize data transposes. Programmers should make this decision based on layout requirements imposed by the compute engine that is going to consume the matrix multiplication output. See NKI Performance Guide for more discussion.

**Tile Size.** The `nki.isa.nc_matmul` API enforces the following constraints on the input/output tile sizes:

1. stationary tensor free axis size (`stationary_fsize`) must never exceed 128, due to the number of PE columns in TensorE.
2. stationary/moving tensor partition axis size (`stationary_psize/moving_psize`) must never exceed 128, due to the number of PE rows and also the number of SBUF partitions.
3. moving tensor free axis size (`moving_fsize`) must never exceed 512, due to the fact that each `nc_matmul` can only write to a single PSUM bank, which can only hold 512 FP32 elements per PSUM partition.

When the shapes of the input matrices defined in the user-level operator exceed any of the above tile size limitation, we must tile the input matrices and invoke multiple `nki.isa.nc_matmul` calls to perform the matrix multiplication. Exceeding the `stationary_fsize` (#1) or `moving_fsize` (#3) tile limitations for M or N should lead to fully independent `nki.isa.nc_matmul` with disjoint output tiles. However, when K exceeds the `stationary_psize/moving_psize` limit, we need to tile the input matrices in the contraction dimension and invoke multiple `nki.isa.nc_matmul` to accumulate into the *same* output buffer in PSUM. Refer to the [Tiling Matrix Multiplications](#) tutorial for a NKI code example.

## Alternative Use Case

One interesting use case of TensorE is low-latency data reshape within NeuronCore, which typically involves multiplying a matrix to be reshaped with a compile-time constant matrix filled with zeros and ones.

As an example, we can perform a 128x128 matrix transposition (i.e., swap the free and partition axis of the matrix) using `nki.isa.nc_matmul(transpose_input, identity)`, where `transpose_input` is the matrix to be transposed and `identity` is a 128x128 identity matrix. In fact, this is exactly what `nki.isa.nc_transpose()` does, when TensorE is chosen as the compute engine.

Similarly, we can broadcast a vector occupying a single partition to M ( $M \leq 128$ ) partitions using `nki.isa.nc_matmul(ones, broadcast_input, is_stationary_onezero=True)`, where `ones` is a 1xM vector filled with ones and `broadcast_input` is the vector to be broadcast. In fact, NKI invokes such `matmul` under the hood when `broadcast_input.broadcast_to((M, broadcast_input.shape[1]))` is called.

In general, we can achieve many more complex data reshapes in TensorE, such as shuffling partitions of a SBUF tensor, by constructing appropriate zero/one patterns as one of the `matmul` inputs.

Finally, we can also leverage TensorE for data summation across SBUF partitions (P-dim summation). For example, a vector laid out across SBUF partitions can be reduced into a single sum using TensorE as shown in the diagram below. Note, this utilizes only a single PE column of the TensorE; therefore, depending on the surrounding operators, this may not be the best use of TensorE. If you can do summation within each partition (F-dim summation), see [nki.isa.tensor\\_reduce](#) for an alternative reduction implementation on Vector Engine. It is recommended to choose the engine based on the natural layout of your input data to avoid any transpositions.

As TensorE is the most performant compute engine of the NeuronCore in terms of FLOPS, the goal is to have it execute meaningful computation at high utilization as much as possible. The above “alternative use cases” stop TensorE from performing *useful* computations at *high* throughput and therefore, should generally be avoided. However, there are situations where it is advisable to use them:

- Operators that do not require heavy `matmuls` anyhow, e.g. normalization, softmax.

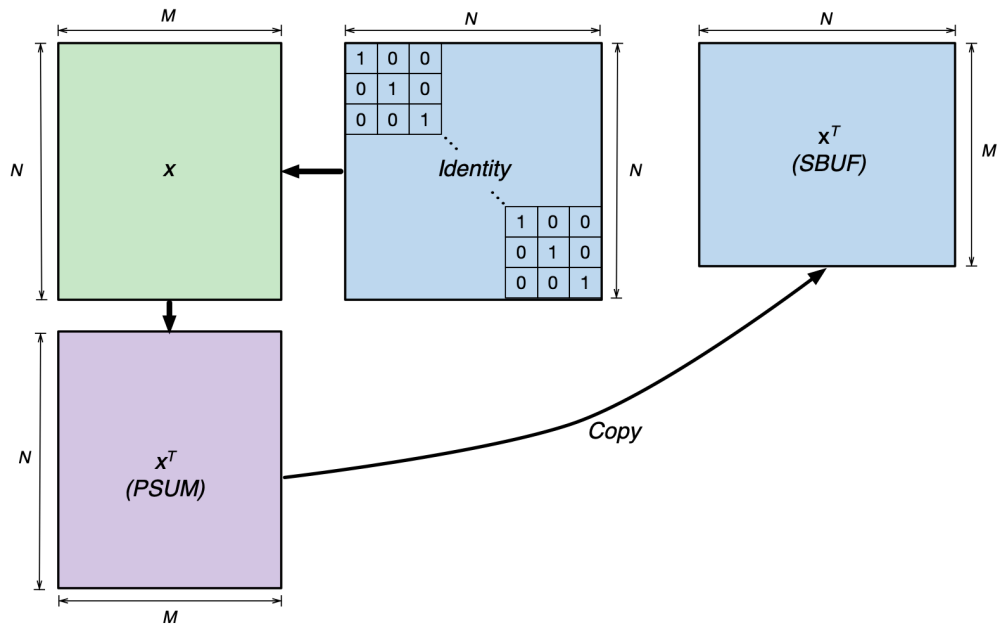


Fig. 7.13: Transposition.

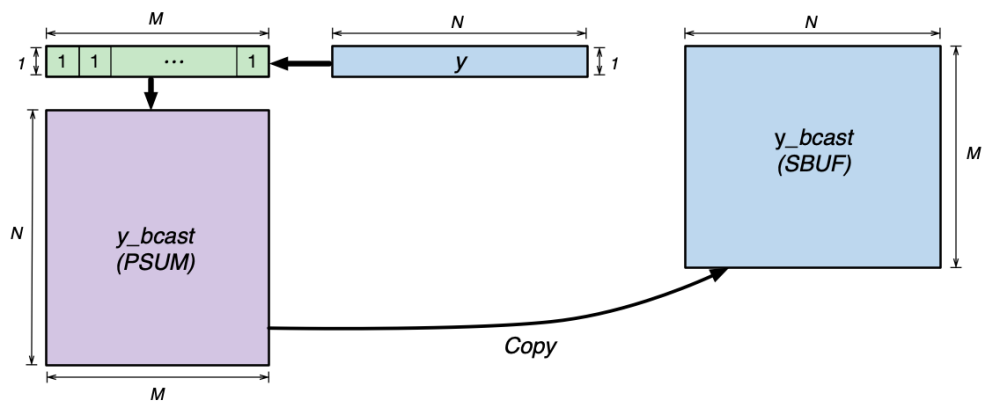


Fig. 7.14: Partition Broadcast.

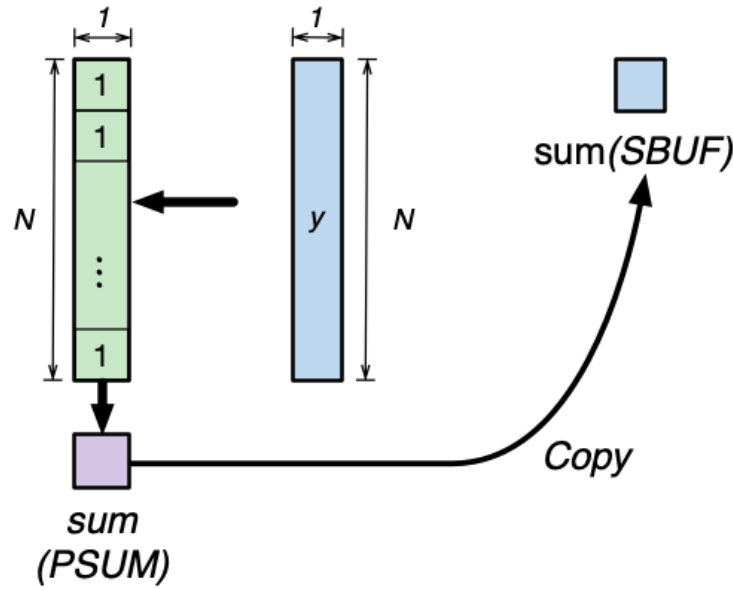


Fig. 7.15: Cross-Partition Accumulation

- Layout conflicts between producer and consumer engines where broadcast/transpose are absolutely unavoidable (see example in fused attention tutorial).

### Performance Consideration

As a rule of thumb, TensorE can achieve the best throughput when it runs many back-to-back `nki.isa.nc_matmul` with both input matrices at the largest possible tiles sizes (stationary is  $128 \times 128$  and moving is  $128 \times 512$ ). In this ideal scenario, TensorE sees the below instruction sequence:

- LoadStationary (LS[0]) ( $128 \times 128$ )
- MultiplyMoving (MM[0]) ( $128 \times 512$ )
- LoadStationary (LS[1]) ( $128 \times 128$ )
- MultiplyMoving (MM[1]) ( $128 \times 512$ )
- ...

**Cost Model:** TensorE is a deeply pipelined engine; therefore, the engine can have several LS&MM instruction pairs in-flight at a given time. Due to this pipelining nature, it is often *not* useful to use end-to-end execution *latency* of a single instruction when estimating the instruction cost. Instead, we can focus on the **initiation interval** of such instructions, that is, the number of cycles between successive instruction launches. Therefore, we can estimate the cost of an instruction *I* by how soon TensorE can issue the next instruction after *I*.

For the sake of discussion, let's assume we have many back-to-back MM instructions with BF16/FP16/TF32/cFP8 input data type that reuse a single pre-loaded stationary inside TensorE. The initiation interval between subsequent MM instructions in this case is roughly  $\max(N, \text{MM\_INIT\_LATENCY})$ , where MM\_INIT\_LATENCY is 64 TensorE cycles on NeuronCore-v2, and  $N$  is the free axis size of moving of current MM (typically set to 512). For FP32 input data type, the instruction cost is roughly 4x higher than BF16/FP16/TF32/cFP8. Therefore, whenever possible, we recommend down-casting FP32 input matrix data type to one of BF16/FP16/TF32/cFP8 before performing matrix multiplications.

Figure below visualizes two pipelined MM instructions:

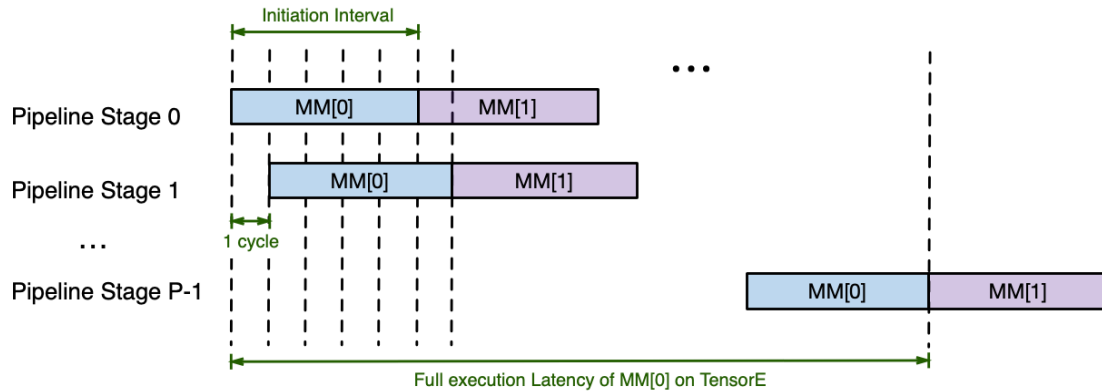
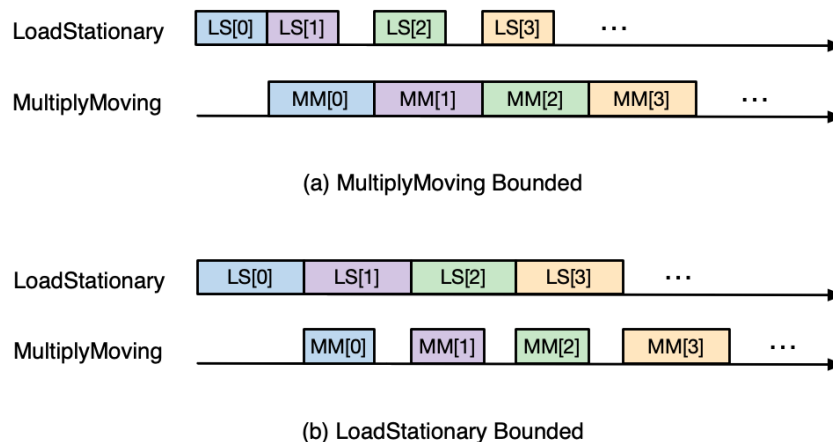


Fig. 7.16: Pipelined multiplyMoving instructions.

**Background LoadStationary:** In typical workloads, TensorE would be alternating between LS and MM instructions with different input matrices. In order to optimize TensorE’s utilization, we also enable a “background LoadStationary” capability, which allows loading of the next stationary tensor in parallel to the computation on the current stationary tensor.

As a result, depending on the relative sizes of the stationary and moving matrices, the overall TensorE performance can be bounded by either LS or MM instructions. Figure below visualizes these two cases. In the ideal scenario where stationary and moving use the largest tile sizes, TensorE should operate in case (a).



Possible execution timeline execution with background LoadStationary

**Fast LoadStationary:** Since LoadStationary is a pure data movement with no computation, TensorE can perform LoadStationary **up to 4x** faster than a MultiplyMoving with the same free axis size. Fast LoadStationary has an important performance implication on `nki.isa.nc_matmul`: When one of the input matrices has a small free axis size and the other has a large free axis size, we prefer to put the matrix with large free axis as the stationary matrix. For example, if we try to do a vector-matrix multiplication, it is recommended to put the matrix as stationary matrix and vector as moving matrix to get the best performance out of TensorE.

## Vector Engine

Vector Engine (VectorE) is specially designed to accelerate vector operations where every element in the output tensor typically depends on multiple elements from input tensor(s), such as vector reduction and element-wise operators between two tensors. VectorE consists of 128 parallel vector lanes, each of which can stream data from a SBUF/PSUM partition, perform mathematical operations, and write data back to each SBUF/PSUM partition in a deeply pipelined fashion.

**Data Types.** VectorE supports all NKI data types (details see [supported data types in NKI](#)) in both input and output tiles. *Arithmetic operations* are performed in FP32, with automatic zero-overhead input and output casting to and from FP32. Refer to `nki.isa` API reference manual for any instruction-specific data type requirements.

**Layout & Tile Size.** VectorE instructions expect the parallel axis of the input and output data to be mapped to the partition dimension. For example, the figure below shows reduction add of a  $N \times M$  matrix along the  $M$  dimension. Since each of  $N$  rows in the matrix can be reduced in parallel, the  $N$  dimension of the matrix should be mapped to the SBUF partition dimension. Refer to the [nki.isa API manual](#) for instruction-specific layout constraint of different VectorE instructions.

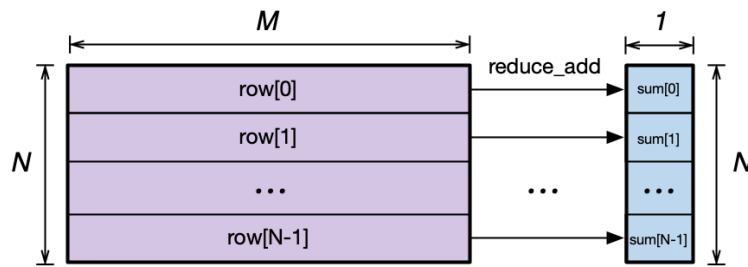


Fig. 7.17: Reduce add on Vector Engine.

In terms of tile size, the majority of VectorE instructions only have limitation on the input/output tile partition dimension size which must not exceed 128, while the free dimension size can be up to 64K elements for SBUF or 4K elements for PSUM. However, there are a few notable exceptions, such as [nki.isa.bn\\_stats](#) which further imposes free dimension size of input tile cannot exceed 512. Refer to the [nki.isa API manual](#) `<nki.language>` for instruction-specific tile size constraints.

## Cross-partition Data Movement

The VectorE also supports a limited set of cross-partition data movement within each group of 32 partitions. The figure below shows connectivity between SBUF and VectorE banks. VectorE consists of four Reshape and Compute banks: each Reshape Bank connects to 32 SBUF/PSUM partitions and outputs 32 parallel streams of data, while each Compute Bank can process 32 parallel data streams using 32 vector lanes. The Compute Bank can write back to 32 SBUF/PSUM partitions.

The Reshape Bank supports the following data movement:

1. *32x32 transpose*: Each Reshape Bank can read in 32 elements per SBUF/PSUM partitions and transpose the partition and free dimension of the incoming 32x32 matrix. This can be invoked by [nki.isa.nc\\_transpose](#) API by selecting VectorE as the execution engine.
2. *32 partition shuffle*: Each Reshape Bank can take an arbitrary *shuffle mask*  $SM^*$  of length 32. The integer value of  $SM[i]$  indicates the source partition ID (modulo 32) that the Reshape Bank output stream  $i$  will get. For example, we can broadcast partition[0] to partition[0-31] using a  $SM$  of 32 zeros. This can be invoked by [nki.isa.nc\\_stream\\_shuffle](#) API.

Refer [here](#) later in this doc for cross-bank data movement.

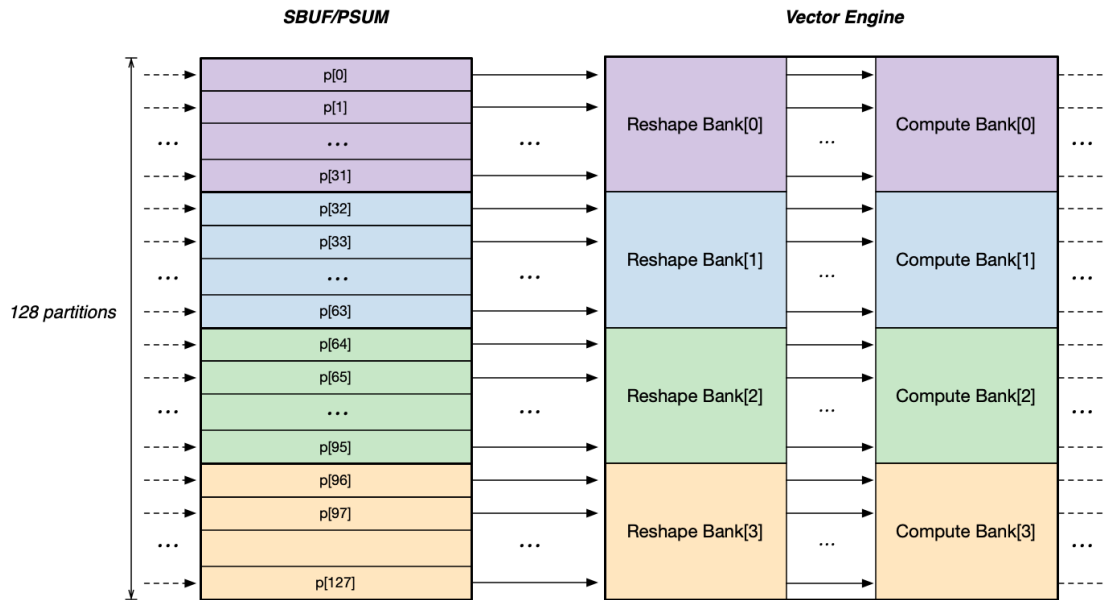


Fig. 7.18: Vector Engine reshape and compute banks.

## Performance Consideration

**128 Parallel Compute Lanes:** VectorE can perform computation with all 128 vector lanes in parallel, with each lane streaming data from/to one SBUF/PSUM partition. Therefore, the performance cost of a VectorE instruction using all 128 lanes is the same as an instruction that uses fewer than 128 lanes.

As a result, we recommend NKI developers to maximize the compute lanes used per VectorE instruction, that is, the partition axis size of input/output tiles of a single `nki.isa` or `nki.language` compute API call. When the partition axis size of input tiles is inevitably fewer than 128 partitions due to high-level operator definition, we could adopt an optimization called “partition vectorization” by packing multiple “small” VectorE instructions of the same operation into a single “large” Vector instruction. Refer to NKI Performance Guide for more detailed discussion of this optimization.

**Cost Model:** In the most common cases where the free axis size ( $N$ ) of the input tile(s) is sufficiently large ( $N > 128$ ), the execution cost of an instruction on VectorE is correlated to  $N$ :

- If there is only one input tile, most VectorE instructions can execute in roughly  $N$  cycles (example: `nki.isa.tensor_scalar`)
- If there are two input tiles, the instruction can execute in roughly  $2N$  cycles (example: `nki.isa.tensor_tensor`)

There are a few exceptions to the above rule, depending on the data types and instruction type. See [NKI ISA API doc](#) for instruction-specific instruction cost details.

In the rare cases where VectorE is running many back-to-back instructions either with  $N \ll 128$  or with every instruction depending on the output tile of the previous instruction, we need to add a static instruction overhead of 100 engine cycles to the above execution cost estimate.

The above rules are for general guidance only. To find out the exact instruction costs for your NKI kernel, you may capture a detailed instruction execution trace on device using [neuron-profiler](#).



## Scalar Engine

Scalar Engine (ScalarE) is specially designed to accelerate scalar operations where every element in the output tensor only depends on one element of the input tensor. In addition, ScalarE provides hardware acceleration to evaluate non-linear functions such as Gelu and Sqrt. The currently supported set of non-linear functions is listed in [here](#). It is worth noting that we can support any new non-linear functions on ScalarE as they come up in new ML model architectures through Neuron SDK software updates. Similar to VectorE, ScalarE consists of 128 parallel lanes, each of which can stream data from a SBUF/PSUM partition, perform mathematical operations, and write data back to each SBUF/PSUM partition in a deeply pipelined fashion.

**Data Types.** ScalarE supports all NKI data types (details see [supported data types in NKI](#)) in both input and output tiles. All internal computation is performed in FP32, with automatic zero-overhead input and output casting to and from FP32.

**Layout & Tile Size.** ScalarE typically evaluates scalar operations (such as, `nki.language.gelu`), which does not impose any input/output tile layout constraints. However, there are additional hardware features in ScalarE that will have layout constraints similar to VectorE (more discussion later).

In terms of tile size, ScalarE instructions only have limitation on the input/output tile partition dimension size which must not exceed 128, while the free dimension size can be up to 64K elements for SBUF or 4K elements for PSUM.

## Pipelined Multiply-Add

Each ScalarE compute lane also supports an additional multiply-add **before** the non-linear function (`func`) is applied in a pipeline fashion. Mathematically, ScalarE implements:

```
# Case 1: scale is SBUF/PSUM vector
# Input: 2D in_tile, 1D scale, 1D bias
# Output: 2D out_tile
for lane_id in range(in_tile.shape[0]):
    for k in range(in_tile.shape[1]):
        out_tile[lane_id][k] = func(in_tile[lane_id][k] * scale[lane_id]
                                   + bias[lane_id])

# Case 2: scale is a compile-time scalar constant in the instruction
for lane_id in range(in_tile.shape[0]):
    for k in range(in_tile.shape[1]):
        out_tile[lane_id][k] = func(in_tile[lane_id][k] * scale
                                   + bias[lane_id])
```

This functionality can be invoked using the `nki.isa.activation` API by specifying a scale for multiplication and bias for addition. The scale can either be a tile from SBUF/PSUM with one element/partition or a compile-time constant. On the other hand, the bias can only be a tile from SBUF/PSUM with one element/partition. A useful mental model for this capability is combining a `nki.isa.tensor_scalar` instruction with a non-linear function evaluation into a single instruction (2x speed-up than two separate instructions).

## Pipelined Reduction

Each ScalarE compute lane also supports reduction **after** the non-linear function (`func`) is applied in a pipeline fashion. On NeuronCore-v2, the reduction operator can only be addition.

Mathematically, ScalarE with accumulation enabled implements:

```
# Input: 2D in_tile, 1D scale (similarly for scalar scale), 1D bias
# Output: 2D out_tile, 1D reduce_res
for lane_id in range(in_tile.shape[0]):
    for k in range(in_tile.shape[1]):
        out_tile[lane_id][k] = func(in_tile[lane_id][k] * scale[lane_id]
                                   + bias[lane_id])
    reduce_res[lane_id] += out_tile[lane_id][k]
```

This functionality can be invoked using the `nki.isa.activation_reduce` API by specifying `reduce_op` as `nki.language.add` and `reduce_res` as the output reduction tile, passed by reference.

A useful mental model for this capability is combining a `nki.isa.activation` instruction with a `nki.isa.tensor_reduce` into a single API, which returns results from **both** APIs. Note, `nki.isa.activation_reduce` invokes two back-to-back ISA instructions on hardware, *Activate* and *ActReadAccumulator*. The *Activate* instruction performs the regular computation as specified in `nki.isa.activation` and also reduction at no additional cost. The reduction result is cached inside ScalarE after *Activate*. The *ActReadAccumulator* instruction is a low cost (roughly 64 ScalarE cycles on NeuronCore-v2) instruction to write the internal reduction result back to SBUF/PSUM, one element per partition.

## Performance Consideration

All the performance notes discussed for *Vector Engine* earlier are applicable to Scalar Engine, with one exception regarding instruction cost for two input tensors - ScalarE can only read up to one input tensor per instruction.

**Instruction Combination.** All `nki.isa.activation` instructions have the same execution cost, regardless of whether we enable the scale multiplication or bias add. Therefore, it is recommended to combine such multiply-add operations with non-linear function evaluation into a single ScalarE instruction if the computation allows it. This is highly useful for ML operators that are **not** TensorE heavy (not matmul-bound). Softmax is one such example, where we typically subtract the maximum value of the input elements before evaluating exponential function for numerical stability.

## GpSimd Engine

GpSimd Engine (GpSimdE) is intended to be a general-purpose engine that can run any ML operators that cannot be lowered onto the other highly specialized compute engines discussed above efficiently, such as applying a triangular mask to a tensor.

A GpSimdE consists of eight fully programmable processors that can execute arbitrary C/C++ programs. Therefore, this engine provides the hardware support for *Neuron Custom Operator*. In addition, each processor is a 512-bit vector machine that can run high-performance vectorized kernels. Every `nki.isa` API running on GpSimdE such as `nki.isa.iota` uses a vectorized kernel implementation that Neuron engineers hand-tune for the underlying processor ISA.

**Data Types.** Each processor in GpSimd supports vectorized computation for

- 16x FP32/INT32/UINT32, or
- 32x FP16/INT16/UINT16, or
- 64x INT8/UINT8

This is in contrast to ScalarE/VectorE which can only perform arithmetic operations in FP32. However, if the GpSimdE program chooses to, it can also access SBUF data of any *supported data types in NKI* and perform data casting to- and from-FP32 at no throughput cost similar to VectorE/ScalarE.

**Layout & Tile Size.** The layout and tile size requirements of GpSimdE highly depend on semantics of the exact instruction. Refer to the *nki.isa API reference guide* for these requirements.

**Memory Hierarchy.** In Trainium/Inferentia2, each GpSimdE processor has 64KB of local data RAM, also called tightly-coupled memory (TCM) as discussed in *Neuron Custom Operator*. The TCM is configured with a 3-cycle access latency and 512-bit data width. Therefore, TCM is often used to store intermediate computation results within a Neuron Custom Operator or GpSimdE instruction.

The eight processors in GpSimdE also have a high-bandwidth read/write interface connected to the SBUF. *Figure 7.19* below illustrates the GpSimdE connectivity to SBUF. Each processor connects to 16 SBUF partitions for both reading and writing: processor[0] connected to partition[0:15], processor[1] to partition[16:31] and so on. Each processor can programmatically send tensor read/write requests to SBUF to access data from the connected partitions. On the read side, once a read request is processed, the tensor read interface can deliver up to 512-bit of data from all 16 connected partitions collectively (up to 32-bit per partition) to the processor per cycle, which matches the 512-bit SIMD width. Similarly, on the write side, the tensor write interface can accept 512-bit of data for writing back to the connected SBUF partitions per cycle.

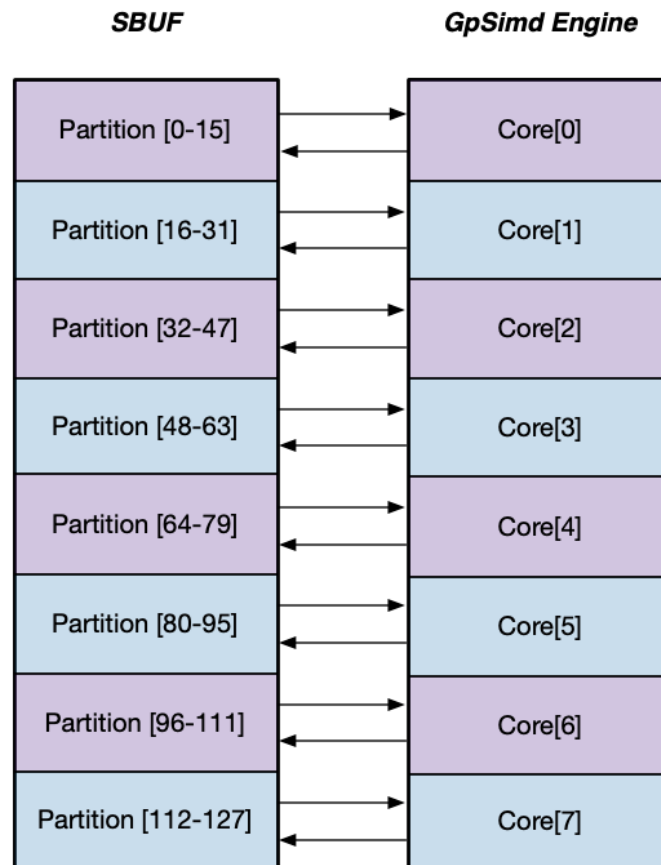


Fig. 7.19: Connectivity between GpSimdE and SBUF.

## Performance Consideration

**128 Parallel Compute Lanes:** Similar to VectorE and ScalarE, GpSimdE has 128 parallel compute lanes for 32-bit computation data types across SIMD lanes of all eight processors. Therefore, it is desirable to invoke GpSimdE instructions that will utilize all the parallel compute lanes, typically through accessing all 128 SBUF partitions for input and output. In addition, since each processor can also handle 32-wide 16-bit or 64-wide 8-bit data type computation, GpSimdE can effectively support 256 or 512 parallel compute lanes internally.

**Cost Model:** Unlike VectorE/ScalarE, there is no rule-of-thumb to estimate execution cost of a GpSimdE instruction. Refer to the [nki.isa](#) API reference manual to find out instruction-specific latency estimates.

## Data Movement

In this section, we will dive into the memory subsystem and discuss how to perform data movement between different memories and also how to do it efficiently. As a reminder, there are three main types of memory on a NeuronDevice: HBM, SBUF, and PSUM, from highest to lowest capacity. Figure below shows the specifications of these memories and their connectivity for one NeuronCore-v2:

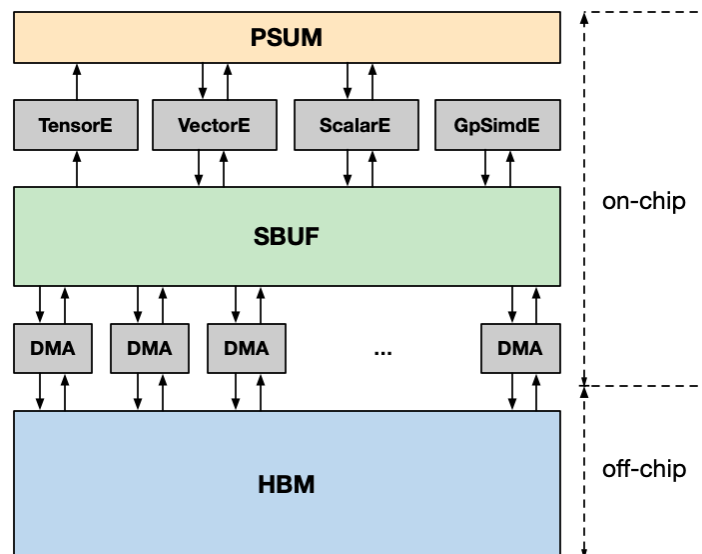


Fig. 7.20: Memory hierarchy.

As shown in the above figure, data movement between HBM and SBUF is performed using on-chip DMA (Direct Memory Access) engines, which can run in parallel to computation within the NeuronCore. Data movement between PSUM and SBUF is done through ISA instructions on the compute engines. However, different compute engines have different connectivity to SBUF/PSUM as indicated by the arrows in the figure. In addition, NeuronCore-v2 has the following restrictions:

1. VectorE and GpSimdE cannot access SBUF in parallel.
2. VectorE and ScalarE cannot access PSUM in parallel.

Therefore, VectorE and GpSimdE instructions that access SBUF must be serialized, similarly for VectorE and ScalarE instructions that access PSUM. This is enforced by Neuron Compiler during NKI kernel compilation, so NKI developers are not required to program such serializations.

The rest of this section will discuss the following topics in detail:

- Data movement between HBM and SBUF using DMAs.
- Accessing SBUF/PSUM tensors using compute engines.
- In-memory accumulation using TensorE and PSUM.

## Data movement between HBM and SBUF using DMAs

Each NeuronCore-v2 is equipped by 16 parallel DMA engines that can perform data movement between any addressable memories in the system. Here, we focus on using these DMA engines to move data between the local SBUF and HBM. Each DMA engine can process one **DMA transfer** at a time driving a peak bandwidth of 27 GiB/s, but all DMA engines can process different DMA transfers in parallel.

Each DMA transfer can gather a list of source **DMA buffers** and then scatter the data into another list of destination DMA buffers. Data within a DMA buffer must be continuous in the memory address map. There is some performance overhead at both DMA buffer and transfer levels, both of which can be amortized by moving a sufficiently large amount of data (more discussion below).

Next, let's examine how HBM and SBUF are laid out in the device memory address map. On one hand, HBM is logically a one-dimensional memory and hence occupies a flat chunk of continuous addresses in the address map. In the most common cases, an HBM tensor in NKI is also contiguous in the HBM address space.

On the other hand, SBUF is considered a two-dimensional memory with 128 partitions as discussed earlier [here](#). [Figure 7.21](#) shows how SBUF addresses fit in the device address map. `sbuf_base_addr` is a 64-bit address dependent on which NeuronCore-v2 on the device the SBUF is located in. The SBUF addresses start from the first byte of partition 0, increment along the free dimension first and then advance onto the next partition.

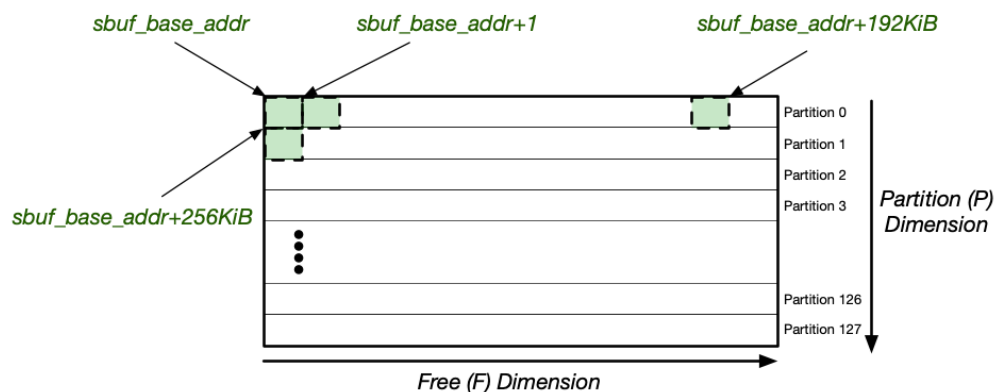


Fig. 7.21: SBUF memory address space.

As discussed in [NKI Programming Model](#), an SBUF tensor in NKI spans one or more partitions, with data starting at the same offset:

As a result, a data movement involving tensor in SBUF will require at least `tensor.shape[0]`, i.e., P dim size, different DMA buffers, since slices of tensor data from different SBUF partitions occupy non-contiguous memory in the address space. If the tensor data slice within each SBUF partition is not contiguous in the F dimension, more DMA buffers will need to be unrolled along the F dim. These DMA buffers are typically grouped into different DMA transfers so that multiple DMA engines can participate in the data movement to maximize memory bandwidth utilization.

In NKI, moving data from HBM to SBUF and from SBUF to HBM are done through [nki.language.load](#) and [nki.language.store](#) APIs, respectively. Neuron Compiler is responsible for converting each NKI API call to DMA transfers and assigning these transfers to different DMA engines. As an example, loading a 128x512 FP32 HBM tensor to SBUF is best done through 16 DMA transfers (one per DMA engine), each moving a scatter-gather list of 8 DMA buffers:

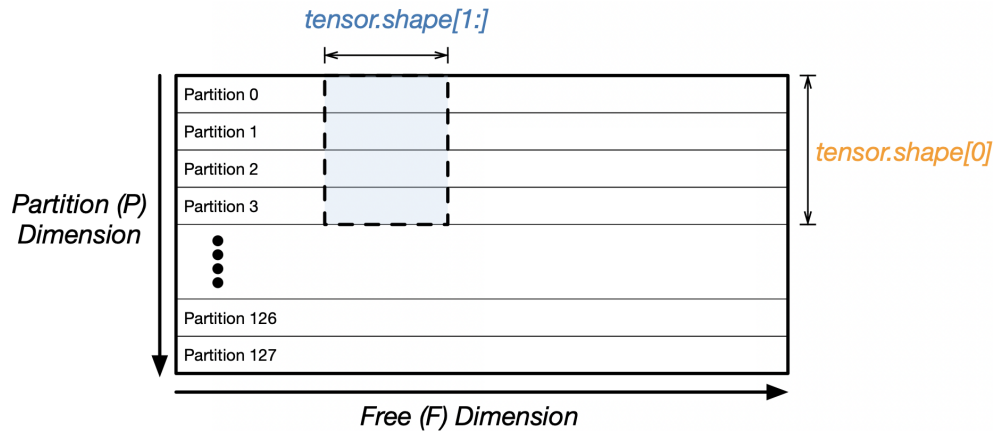


Fig. 7.22: SBUF tensor.

```
import neuronxcc.nki.language as nl
tile = nl.load(in_tensor[0:128, 0:512])
```

To achieve good performance out of the DMAs, we generally aim to:

1. Move a large amount of contiguous data in each DMA buffer to amortize DMA buffer overhead
2. Move a large amount of data in each DMA transfer to amortize DMA transfer overhead.
3. Invoke as many parallel DMA transfers on the available DMA engines as possible.

These goals ultimately boil down to a quick optimization rule: maximize **both free (4KiB or above) and partition (ideally 128) dimension sizes** when moving tensors between SBUF and HBM using `nki.language.load` and `nki.language.store`. Refer to the *NKI Performance Guide* for more information on optimizing performance of data movements between HBM and SBUF.

## Accessing SBUF/PSUM tensors using compute engines

Figure 7.23 shows a simplified timeline of how compute engines **stream** data in and out of on-chip SRAM (SBUF or PSUM). Refer to Figure 7.10 for the available connectivity between engines and SBUF/PSUM. At a high level, the compute engines are able to pipeline data reads, computation and writes along the F dimension of the src/dst tensors. In every cycle, each engine can read 128 elements across 128 SBUF/PSUM partitions, perform a computation on previously read 128 elements, and write 128 previously computed results to SBUF/PSUM. In other words, the P axis of a tensor is the *parallel* dimension for SBUF/PSUM data accessing, while the F axis of the tensor is the *time* dimension for data accessing.

When accessing SBUF/PSUM tensors in an instruction, we need to follow different rules in the P and F dimensions. First, hardware does not allow P dimension striding when accessing data from a single SBUF/PSUM tensor. Therefore, a valid src/dst tensor of an instruction must occupy a continuous number of partitions. In addition, the hardware further enforces which partition a tensor can start from (`start_partition`) based on the number of partitions the tensor occupies (`num_partition`). This is currently handled by the tensor allocator in Neuron Compiler during NKI kernel compilation process:

- If  $64 < \text{num\_partition} \leq 128$ , `start_partition` must be 0
- If  $32 < \text{num\_partition} \leq 64$ , `start_partition` must be 0 or 64
- If  $0 < \text{num\_partition} \leq 32$ , `start_partition` must be one of 0/32/64/96

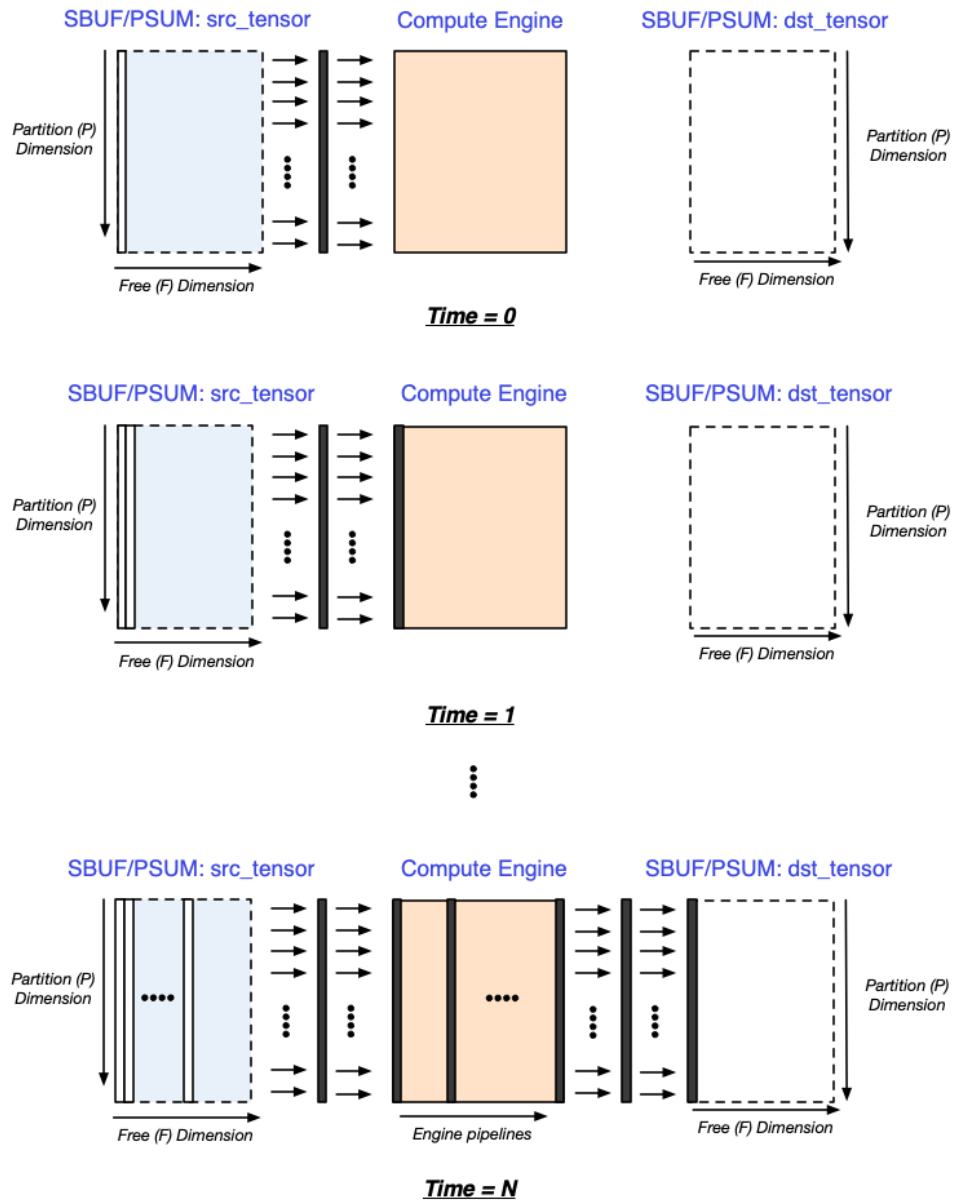


Fig. 7.23: Data streaming between SBUF and compute engine.

On the other hand, data accessing along the free dimension is a lot more flexible: the src/dst tensor of an engine instruction can support up to four-dimensional tensorized access pattern with a stride in each dimension within each partition. At the ISA level, each F axis in the tensor can have a size expressed in `uint16` and a stride expressed in `int16`, measured in data elements. As an example, if the tensor data type is BF16, and the stride of the most-minor F dimension is set to 10, then we will stride across 20B within a partition at a time. Refer to [Tile Indexing in NKI Programming Guide](#) to learn about how to index SBUF/PSUM tensors to achieve F dimension striding in NKI syntax.

Lastly, as implied in [Figure 7.23](#), when accessing a SBUF/PSUM tensor, all active partitions must follow the same F dimension access pattern. In other words, at every time step, the engine read/write interface will access data elements at the same *offset* within each active partition.

## Cross-Partition Connectivity

The majority of VectorE/ScalarE/GpSimdE instructions on NeuronCore-v2 require `src_tensor` and `dst_tensor` to occupy the same number of partitions. When the number of partitions involved exceeds 64, by the `start_partition` rule discussed above, the `src_tensor` and `dst_tensor` in such cases must both start from partition 0. Therefore, we effectively cannot perform any cross-partition data movement when `num_partition > 64`: each partition of `src_tensor` data will eventually flow into the corresponding partition in `dst_tensor`.

However, when `num_partition < 64`, VectorE/ScalarE/GpSimdE on NeuronCore-v2 supports two styles of cross-partition SBUF/PSUM data movement patterns: 1) cross-half movement for  $32 < \text{num\_partition} \leq 64$  and 2) cross-quadrant movement for  $0 < \text{num\_partition} \leq 32$ . Figure below illustrates these two patterns for `num_partition=64` and `num_partition=32`. The shaded portion of the Engine block indicates the active lanes for the given instruction. With these movement patterns, each partition in `src_tensor` still has a one-to-one mapping to each partition in `dst_tensor`.

## Performance Consideration

**Access pattern.** As discussed previously in the context of compute engine utilization, it is recommended to use as many partitions as possible when accessing SBUF/PSUM tensors to saturate the available data streaming bandwidth. In addition, accessing with a large stride in the most-minor (fastest) F dimension will incur performance penalty. When the most-minor F dimension stride is less than 16 bytes, SBUF/PSUM on NeuronCore-v2 can supply a peak bandwidth of 128 elements/cycle at 1.4 GHz for each tensor read/write interface. A 16-byte stride is equivalent to 4 elements for 32-bit data types, 8 elements for 16-bit data types or 16 elements for 8-bit data types. If the most-minor F dimension stride exceeds 16 bytes, the achievable bandwidth of each tensor read/write interface will be half of the peak bandwidth, which translates to roughly 50% performance hit on the instructions.

**Concurrent SBUF/PSUM accesses by engines.** As mentioned earlier, NeuronCore-v2 has the following on-chip RAM access restrictions:

1. Vector Engine and GpSimd Engine cannot access SBUF in parallel
2. Vector Engine and Scalar Engine cannot access PSUM in parallel

Despite these restrictions, SBUF is capable of driving peak bandwidth in each tensor read/write interface connected to VectorE/ScalarE/TensorE or GpSimdE/ScalarE/TensorE *simultaneously* without bandwidth interference. Similarly, PSUM can drive peak bandwidth for VectorE/TensorE or ScalarE/TensorE *simultaneously*.

**Tensor access overhead.** Initiating a tensor access request from an engine to its SBUF/PSUM read/write interface incurs a static overhead approximately 60 cycles on NeuronCore-v2. Compute engines can typically hide some of this latency through instruction level parallelism. However, it is still highly recommended to access tensors with large P and F dimension sizes whenever possible to amortize this overhead.



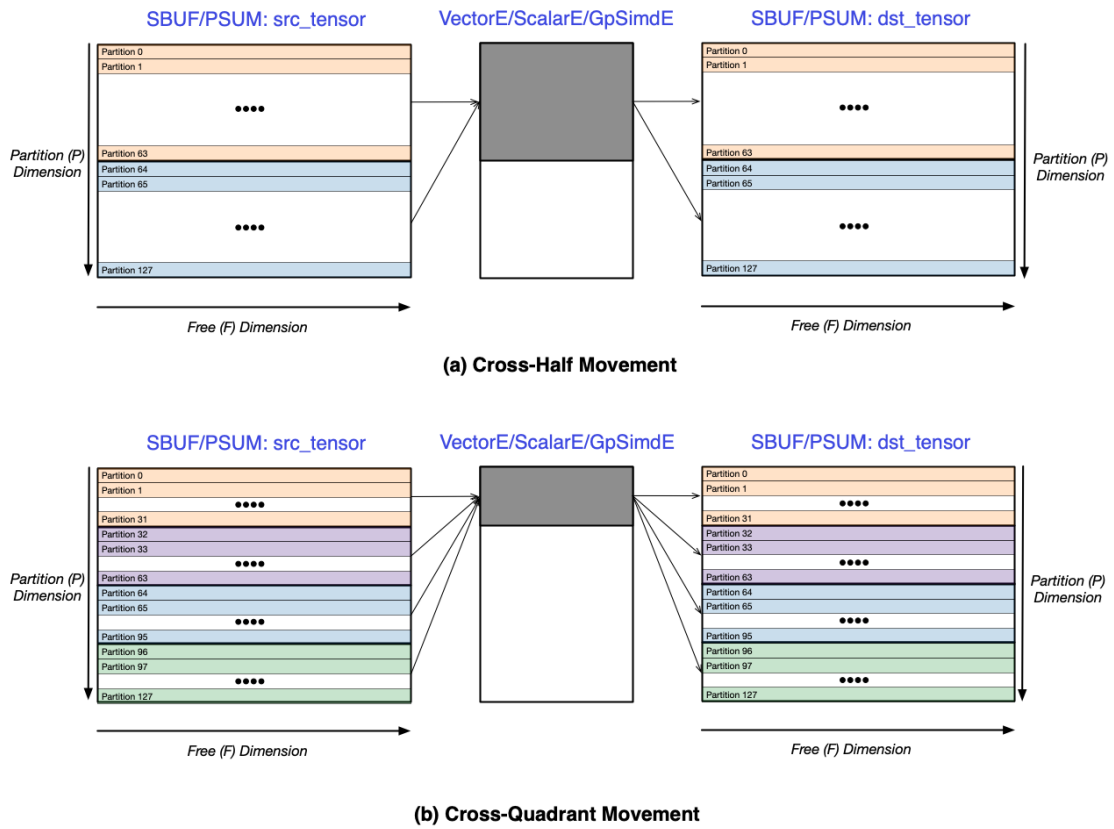


Fig. 7.24: Cross-partition connectivity.

## Near-memory accumulation in PSUM

As shown in Figure 7.10, both VectorE and ScalarE have read and write access to PSUM, while TensorE only has write access. In fact, PSUM is designed to be a landing buffer for TensorE with near-memory accumulation capabilities that allows read-accumulate-write to every 4B element in memory. Note, this accumulation mechanism can *only* be controlled by TensorE. VectorE and ScalarE can only access PSUM like a regular SRAM similar to SBUF.

Next, let's discuss how TensorE can write outputs to PSUM. As previously discussed, PSUM is organized into 128 *partitions*, each consisting of 16KB of memory. Each partition is further divided into 8 PSUM banks, with each bank holding up to 512 32-bit values. The output tile of a TensorE matrix multiplication instruction (`nki.isa.nc_matmul`) must **fit** into one PSUM bank per partition, which is the fundamental reason for the *free dimension size limitation* for the moving tensor. Every `nc_matmul` instruction can choose whether to *override* existing bank data with instruction output or *accumulate* instruction output into existing bank data element-wise.

The accumulation mode of PSUM is particularly useful when the high-level matmul operator has a contraction dimension (i.e., stationary/moving partition dimension of `nki.isa.nc_matmul`) greater than 128. As an example, let's assume the following matmul dimensions:

- `x.shape = [128, 256]`
- `y.shape = [256, 512]`

Figure below shows this matmul mathematically and also how we would tile the contraction dimension. With tiling, we slice both `x` and `y` in the contraction dimension to get `[x0, x1]` and `[y0, y1]` input tiles. To get the final output result, we need to perform:

- `output0 = matmul(x0, y0)`
- `output1 = matmul(x1, y1)`
- `output = output0 + output1`

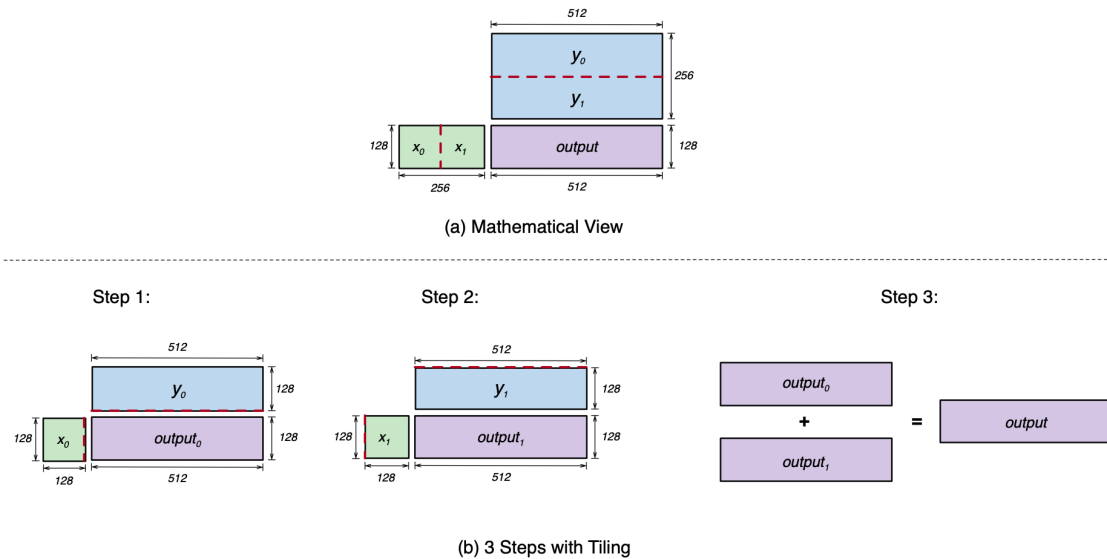


Fig. 7.25: Matmul tiling (mathematical view).

PSUM accumulation effectively combines Step 2 and 3 above into a single TensorE `nki.isa.nc_matmul` instruction. Assuming we have `x` in the transposed layout in SBUF, visually the above tiled matmul example will have two back-to-back `nki.isa.nc_matmul` instructions on TensorE:

Effectively, the first `nki.isa.nc_matmul` instruction overwrites the destination PSUM bank with the instruction output. The second instruction accumulates instruction output onto the previous instruction's result in the same PSUM.

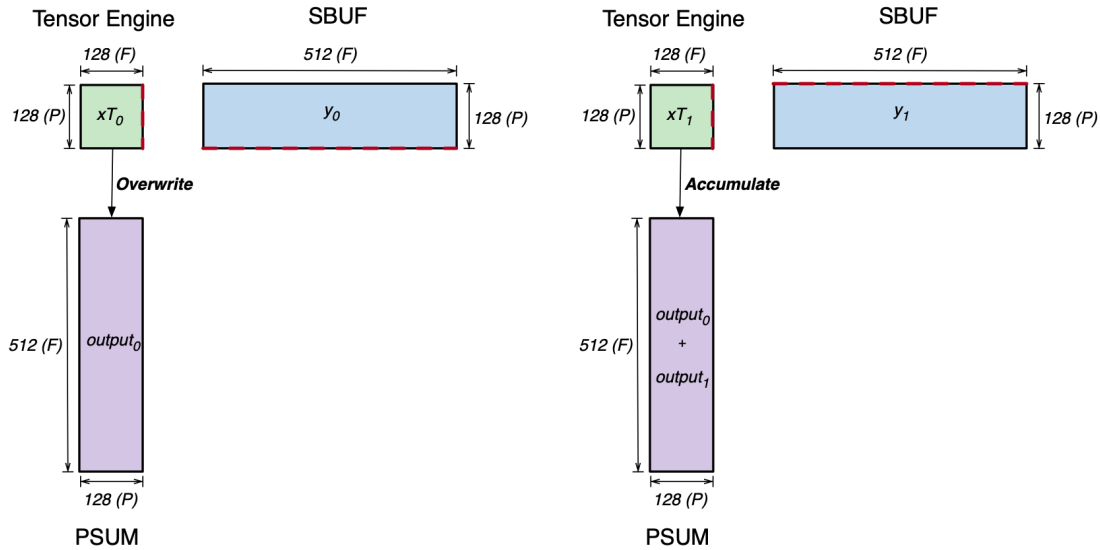


Fig. 7.26: Matmul tiling (hardware view).

The PSUM accumulation is always done in FP32. A series of TensorE matmul instructions with the first one writing to a PSUM bank and more subsequent instructions accumulating into the same PSUM bank data is called a *matmul accumulation group*.

In current release of NKI, the `nki.isa.nc_matmul` does not have an explicit control field to indicate *overwrite* or *accumulate* for the PSUM. Instead, NeuronCompiler relies on the following NKI code pattern to trigger PSUM accumulation:

```
# condition 1: a psum buffer with zeros
psum_buf = nl.zeros(..., buffer=nl.psum)

# condition 2: an affine range loop
for i in nl.affine_range(N):
    # condition 3: add matmul results from TensorEngine
    psum_buf += nl.matmul(stationary_tile, moving_tile) # or nisa.nc_matmul
```

Refer to the [Tiling Matrix Multiplications](#) tutorial for a detailed implementation.

**Note:** Due to current limitations in NKI, `psum_buf[...] = psum_buf + nisa.nc_matmul(stationary_tile, moving_tile)` will not reliably trigger the PSUM accumulation architecture feature. Therefore, even though this alternative syntax is functionally equivalent to the use of `+=`, it may get lowered to `nisa.tensor_tensor` on VectorEngine for accumulation instead, leading to much lower performance.

Finally, with 8 PSUM banks per partition, TensorE can have up to eight outstanding matmul accumulation groups, which allows flexible scheduling of matmul instructions on TensorE. Also, the extra buffering from multiple PSUM banks allows us to pipeline TensorE computation with other compute engines: TensorE can move onto the next accumulation group without waiting for VectorE/ScalarE to evict previous accumulation group results.

This document is relevant for: Inf2, Trn1, Trn2

This document is relevant for: Inf2, Trn1, Trn2

## Trainium2 Architecture Guide for NKI

This guide covers hardware architecture of third-generation NeuronDevices: Trainium2. We assume readers have gone through [Trainium/Inferentia2 Architecture Guide](#) in detail to understand the basics of NeuronDevice Architecture.

Fig. 7.27 shows a block diagram of a Trainium2 device, which consists of:

- 8 NeuronCores (v3).
- 4 HBM stacks with a total device memory capacity of 96GiB and bandwidth of 2.9TB/s.
- 128 DMA (Direct Memory Access) engines to move data within and across devices.
- 20 CC-Cores for collective communication.
- 4 NeuronLink-v3 for device-to-device collective communication.

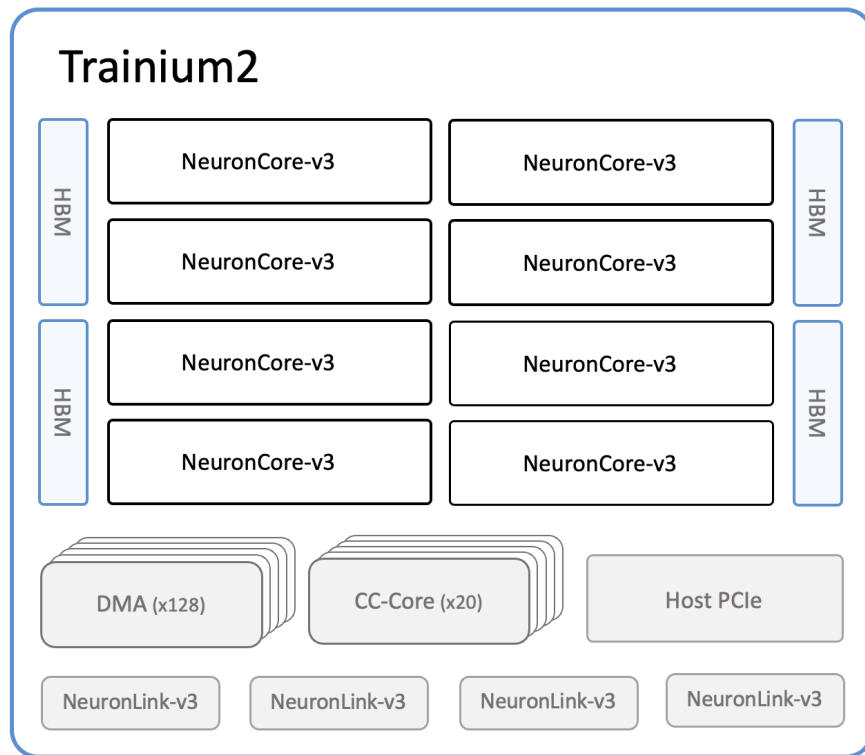


Fig. 7.27: Trainium2 Device Diagram.

For a high-level architecture specification comparison from Trainium1 to Trainium2, check out [Neuron architecture guide](#).

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Profiling NKI kernels with Neuron Profile

In this tutorial, we use Neuron Profile to view the execution trace of a NKI kernel captured on a NeuronCore. In doing so, we learn about:

- Installation and usage of Neuron Profile.
- Inspecting a detailed execution timeline of compute engine instructions and DMA engine activities generated from your NKI kernel.

As background, [Neuron Profile](#) is the tool you need to visualize where time is being spent during kernel execution on NeuronDevices, which is crucial for identifying performance bottlenecks and opportunities of your kernel. Neuron Profile produces runtime execution data for every instruction executed on each compute engine and also every data movement activity completed by DMA engines. Neuron Profile also reports key performance metrics such as compute engine and memory bandwidth utilization, which allows developers to quickly find out the achieved hardware efficiency of their kernel. Profiling typically has near zero overhead thanks to the dedicated on-chip profiling hardware in NeuronDevices.

### Install Neuron Profile

Make sure you have the latest version of the `aws-neuronx-tools`, which includes updated profiling support for NKI kernels. Neuron Profile is included within this package and is installed to `/opt/aws/neuron/bin`.

The `aws-neuronx-tools` package comes pre-installed on [Neuron DLAMIs](#). For detailed installation instructions see [Neuron Profile User Guide: Installation](#).

### Profile a NKI Kernel

#### Profile using `neuron-profile capture`

To profile a NKI kernel the required steps are (1) enable `NEURON_FRAMEWORK_DEBUG` to tell the compiler to save the NEFF file, (2) execute the NKI kernel to generate the NEFF, and (3) run `neuron-profile capture` to generate a NTFF profile. Each step is described in more detail below.

We will profile a NKI kernel which computes the element-wise exponential of an input tensor of any 2D shape. The rest of this tutorial will use a performance profile generated from this kernel as an example. Full code of `prof-kernel.py`:

```

1  """
2  Example kernel used to demmonstrate Neuron Profile.
3  """
4  import torch
5  from neuronxcc import nki
6  import neuronxcc.nki.language as nl
7  import math
8  import os
9  os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
10 os.environ["NEURON_CC_FLAGS"] = " --disable-dge "
11
12 @nki.jit
13 def tensor_exp_kernel(in_tensor):
14     """NKI kernel to compute elementwise exponential of an input tensor
15
16     Args:

```

(continues on next page)

(continued from previous page)

```

17     in_tensor: an input tensor of ANY 2D shape (up to SBUF size)
18     Returns:
19         out_tensor: an output tensor of ANY 2D shape (up to SBUF size)
20     """
21     out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
22                             buffer=nl.shared_hbm)
23
24     sz_p, sz_f = in_tensor.shape
25
26     i_f = nl.arange(sz_f)[None, :]
27
28     for p in nl.affine_range(math.ceil(sz_p / nl.tile_size.pmax)):
29         # Generate tensor indices for the input/output tensors
30         # pad index to pmax, for simplicity
31         i_p = p * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]
32
33         # Load input data from external memory to on-chip memory
34         # only read up to sz_p
35         in_tile = nl.load(in_tensor[i_p, i_f], mask=(i_p<sz_p))
36
37         # perform the computation
38         out_tile = nl.exp(in_tile)
39
40         # store the results back to external memory
41         # only write up to sz_p
42         nl.store(out_tensor[i_p, i_f], value=out_tile, mask=(i_p<sz_p))
43
44     return out_tensor
45
46 if __name__ == "__main__":
47     from torch_xla.core import xla_model as xm
48     device = xm.xla_device()
49
50     in_tensor = torch.rand((250, 512), dtype=torch.float32).to(device=device)
51
52     out_tensor = tensor_exp_kernel_(in_tensor)
53     print(f"output_nki={out_tensor}")

```

To profile this NKI kernel, follow these steps:

1. Enable Neuron debug output by setting the `NEURON_FRAMEWORK_DEBUG` environment variable. This will trigger the Neuron compiler to save the Neuron Executable File Format (NEFF) artifact to the current directory after compilation of your NKI kernel. The NEFF contains all hardware instructions required to execute your NKI kernel on a NeuronDevice, as well as metadata and debug info needed for profiling. For example, add the following lines to your NKI kernel source file:

```

import os
os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
os.environ["NEURON_CC_FLAGS"] = " --disable-dge "

```

**Note:** Use the flag `--disable-dge` to temporarily disable a new compiler feature which is interfering with DMA debugging information display in `neuron-profile`. This is highly recommended to improve NKI performance debugging

experience until we release a software fix for this issue.

2. Compile your NKI kernel to create a NEFF in your current directory:

```
$ python3 prof-kernel.py
```

**Note:** Find your NEFF named similarly to `MODULE_0_SyncTensorsGraph.13_12659246067793504316.neff`.

3. Profile the NEFF. This profiling step executes the NEFF on the NeuronDevice and records a raw execution trace into an Neuron Trace File Format (NTFF) artifact.

```
$ neuron-profile capture -n <path_to_neff> -s profile.ntff --profile-nth-exec=2
```

This will save your NTFF profile to `profile_exec_2.ntff`.

**Note:** The `--profile-nth-exec=2` option will profile your NEFF twice on the NeuronDevice and output a NTFF profile for the second iteration. This is recommended to avoid one-time warmup delays which can be seen in the first iteration of execution.

In *View Neuron Profile UI*, we will view the profile in a user-friendly format using the Neuron Profile UI.

## Profile using `nki.profile`

You may also use the `nki.profile` API to generate a NEFF and NTFF programmatically.

Below is an example NKI kernel decorated by `nki.profile`. Full code of `prof-kernel-profile.py`:

```
1  """
2  Example kernel used to demonstrate Neuron Profile with nki.profile.
3  """
4  from neuronxcc import nki
5  from neuronxcc.nki.typing import tensor
6  import neuronxcc.nki.language as nl
7  import math
8  from pathlib import Path
9
10 WORKING_DIRECTORY = Path.home() / 'reports'
11
12 @nki.profile(working_directory=WORKING_DIRECTORY, save_neff_name='file.neff', save_trace_
13 ↪ name='profile.ntff', profile_nth=2)
14 def tensor_exp_kernel(in_tensor):
15     """NKI kernel to compute elementwise exponential of an input tensor
16     Args:
17         in_tensor: an input tensor of ANY 2D shape (up to SBUF size)
18     Returns:
19         out_tensor: an output tensor of ANY 2D shape (up to SBUF size)
20     """
21     out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
22                             buffer=nl.shared_hbm)
```

(continues on next page)

(continued from previous page)

```

23  sz_p, sz_f = in_tensor.shape
24  i_f = nl.arange(sz_f)[None, :]
25  for p in nl.affine_range(math.ceil(sz_p / nl.tile_size.pmax)):
26      # Generate tensor indices for the input/output tensors
27      # pad index to pmax, for simplicity
28      i_p = p * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]
29      # Load input data from external memory to on-chip memory
30      # only read up to sz_p
31      in_tile = nl.load(in_tensor[i_p, i_f], mask=(i_p<sz_p))
32      # perform the computation
33      out_tile = nl.exp(in_tile)
34      # store the results back to external memory
35      # only write up to sz_p
36      nl.store(out_tensor[i_p, i_f], value=out_tile, mask=(i_p<sz_p))
37
38  return out_tensor
39
40  if __name__ == "__main__":
41      tensor_exp_kernel_(tensor[[250, 512], nl.float32])

```

To use *nki.profile* to create a NEFF file, NTFF profile, and dump reports in your specified directory, execute the example NKI kernel with:

```
$ python3 prof-kernel-profile.py
```

In *View Neuron Profile UI*, we will view the profile in a user-friendly format using the Neuron Profile UI.

## View Neuron Profile UI

Neuron Profile has an interactive web based UI used to view execution traces. In this section we will open Neuron Profile UI and view NKI specific profiling information. NKI specific information can be found in several places including instruction hover details, instruction click details, search results, and box select results. This section assumes that you followed the previous step to create a NEFF and NTFF.

To view the Neuron Profile web UI, execute the view command:

```
$ neuron-profile view -n <path_to_neff> -s <path_to_ntff> --db-bucket=my_kernel
```

The above command should print a URL that you can click to open the web UI:

```
View profile at http://localhost:3001/profile/my_kernel
```

---

**Note:** You must keep the view command running when viewing profiles.

---



---

**Note:** The `--db-bucket=my_kernel` argument is used to set a custom URL for the profile. Omitting this argument will generate a URL with a unique ID.

---

If `neuron-profile view` is run on a remote instance, you may need to use port forwarding to access the web UI. From your local machine, SSH to the remote instance and forward ports 3001 (the default `neuron-profile` HTTP server port) and 8086 (the default InfluxDB port). Then in the browser, go to `localhost:3001` to view the profiles.



```
$ ssh <user>@<ip> -L 3001:localhost:3001 -L 8086:localhost:8086
```

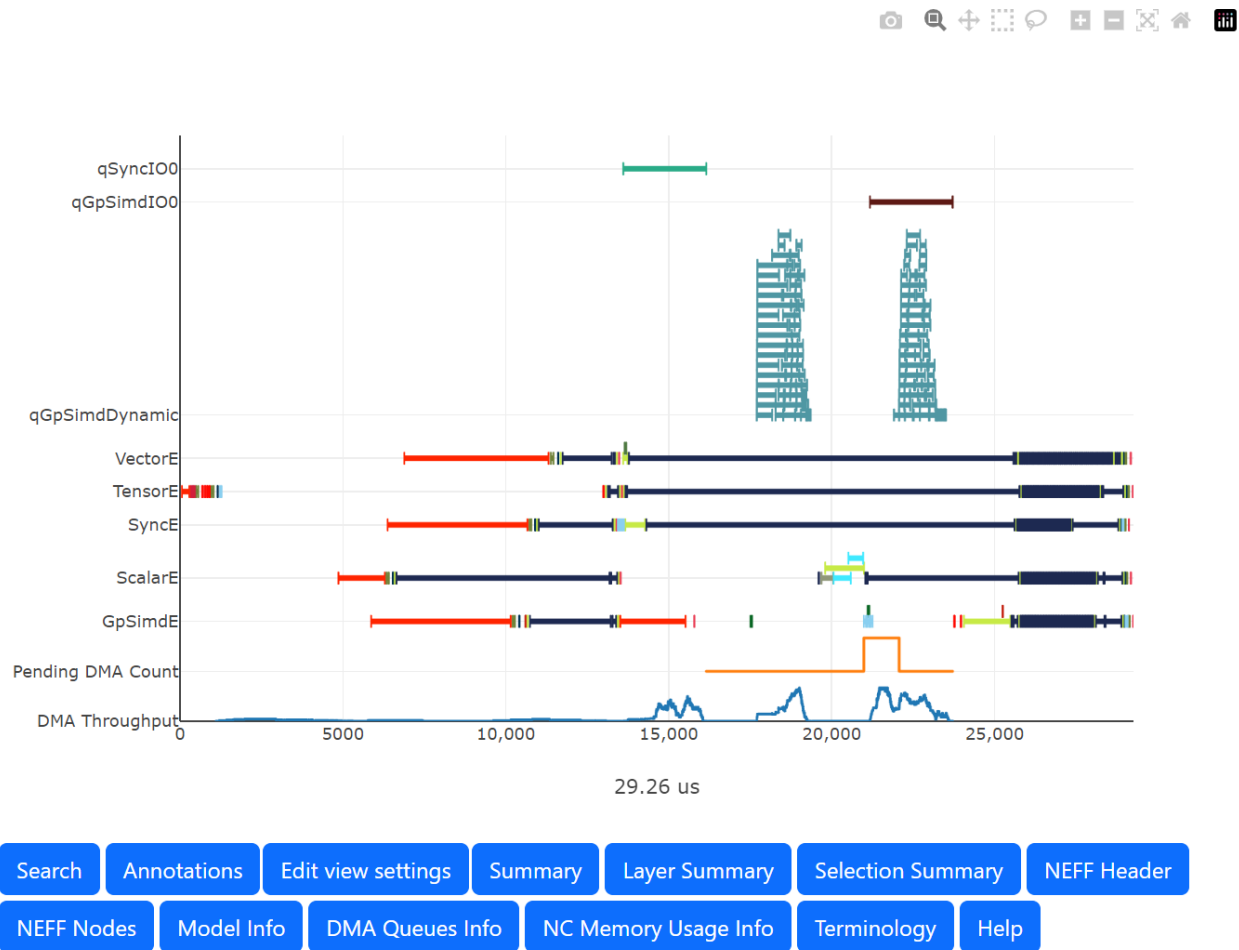


Fig. 7.28: Screenshot of the Neuron Profile UI.

If you hover over any engine instruction in the timeline with your mouse, you will see instruction details in a pop-up box.

If you click on any engine instruction in the timeline with your mouse, you will see instruction details in a panel below the timeline.

## Search

You can search for instructions associated with a specific line of NKI source code (for example “prof-kernel.py:33”). Or you can search for all instructions that have an associated line of NKI code by searching for “/./”, as seen below. This can be used to find the slowest NKI kernel instruction, for example line “prof-kernel.py:37” seen in Fig. 7.32 below. This helps with identifying performance bottlenecks while optimizing a NKI kernel. For help on what kinds of searches are possible, click on the help tooltip in the search panel. Search queries are saved into your browser’s URL so that you can share or revisit the same search using that URL.



Fig. 7.29: Instruction hover details including the line of NKI source code that generated this instruction.

## Event Details - ScalarE

**Timing:**

- **Start:** 19686 ns **i**
- **End:** 20054 ns **i**
- **Duration:** 368 ns **i**
- **Evt wait time:** 11 ns **i**

**Operation:**

- **Opcode:** ACT\_TABLE\_LOAD **i**
- **Operands:** table\_sel=0 **i**

**NKI:**


- **Nki source location:** [prof-kernel.py:33 \(Click to see source\)](#) **i**

**HLO:**

- **Layer:** \_custom-call.2 **i**
- **Hlo name:** %custom-call.2 = custom-call(%p1.2) **i**




Fig. 7.30: Instruction click details including the line of NKI source code that generated this instruction.



**Search** 

---

Category: instruction ▾

+ Add field  Search

Note: fields are joined by "AND"

Tip: enclose string with forward slashes to use regex, eg. /value/  

X Field: nki\_source\_location ▾

Fig. 7.31: Search panel and search help tooltip.

## Search Result Summary

Duration: 29.26 **us**Start Time: 0 **s**End Time: 29.26 **us**

Event Count: 9

Event Duration Sum: 1.21 **us**Event Duration Active: 1.16 **us** (3.95 %)

ScalarE (count = 3, active dur = 3.10%)

GpSimdE (count = 6, active dur = 0.85%)

Field: opcode

Field: nki\_source\_location

Count: 6

| Value                             | Value Count | Value Count % | Value Sum % | Dur Sum          | Dur %   | Dur Avg         |
|-----------------------------------|-------------|---------------|-------------|------------------|---------|-----------------|
| <a href="#">prof-kernel.py:30</a> | 1           | 16.67 %       |             | 42.00 <b>ns</b>  | 14.09 % | 42.00 <b>ns</b> |
| <a href="#">prof-kernel.py:37</a> | 5           | 83.33 %       |             | 256.00 <b>ns</b> | 85.91 % | 51.20 <b>ns</b> |

Fig. 7.32: Search results. The line of NKI source code that generated each instruction will appear in the search result summary.

## Box Select

You can click and drag on the timeline to select a range of instructions using the Box Select functionality. A summary will be produced that includes which lines of NKI source code produced these instructions. This helps with understanding a portion of the timeline. Selecting a large number of instructions may take some time to retrieve from the database.

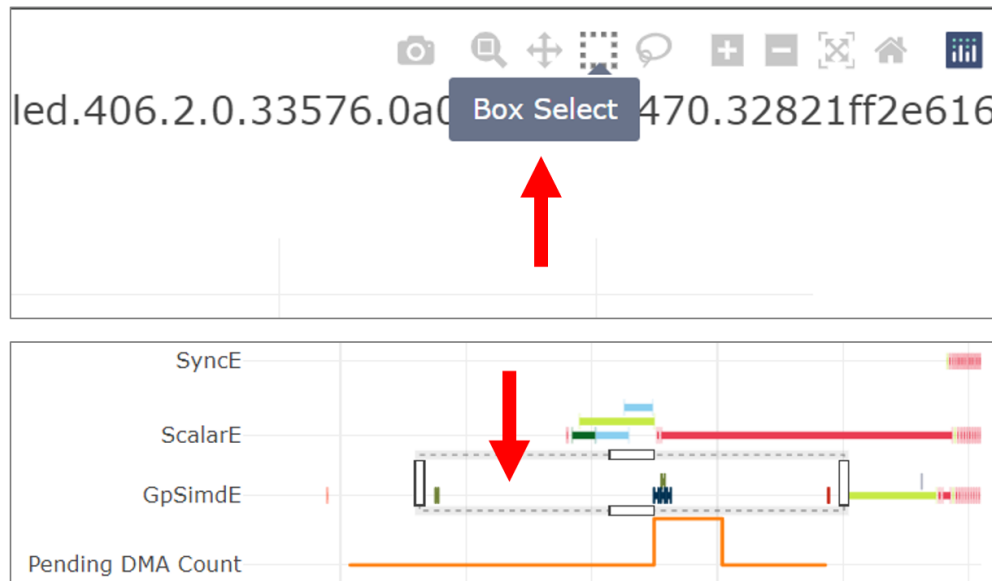


Fig. 7.33: Click on the “Box Select” button and then click and drag on a region of the timeline.

**Note:** An empty value for “nki\_source\_location” means that the instruction is not associated with a NKI source code line.

## View NKI Source Code in Neuron Profile

You may optionally include NKI source code file contents for display in Neuron Profile. This feature loads your NKI source code into an integrated code viewer, side-by-side with the execution timeline in the web UI. Including the source code makes it easier to navigate between instruction trace and NKI source code and also to track the version of code that produced the profile. Note, even without uploading the source code to Neuron Profile, the NKI source filename and line number are always available in instruction detail view as discussed in [View Neuron Profile UI](#).

To include NKI source code in the Neuron Profile UI you can use the `view` command with the `--nki-source-root` argument to pass in the folder of NKI source code:

```
$ neuron-profile view -n <path_to_neff> -s <path_to_ntff> --nki-source-root /home/ubuntu/
↪my_nki/ --db-bucket=my_kernel
```

To open the NKI source code viewer, click on an instruction that has a “Nki source location” field as shown in [Fig. 7.30](#). In the instruction’s details panel, the “nki\_source\_location” field should appear as a link. Clicking on the link will open the NKI source code viewer and highlight the associated line of NKI code. Inside the source code viewer, you can also click on any line of NKI source code to search for all instructions that were generated by that line of code.

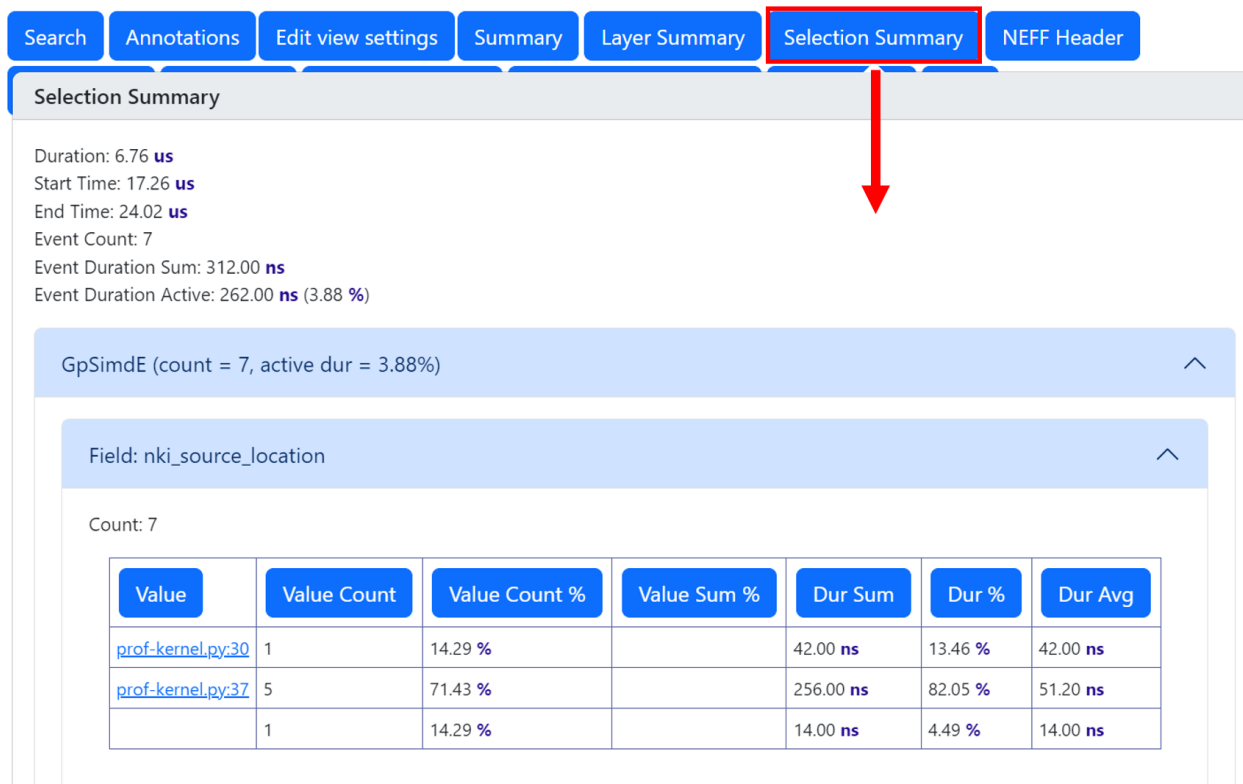


Fig. 7.34: Box select results. The line of NKI source code that generated each instruction will appear in the box select summary.

```

15 Args:
16     in_tensor: reference to an input tensor of ANY 2D shape (up to SBUF size)
17     out_tensor: reference to an output tensor of ANY 2D shape (up to SBUF size)
18     """
19     sz_p, sz_f = in_tensor.shape
20
21     i_f = nl.arange(sz_f)[None, :]
22
23     for p in nl.affine_range(math.ceil(sz_p / nl.tile_size.pmax)):
24         # Generate tensor indices for the input/output tensors
25         # pad index to pmax, for simplicity
26         i_p = p * nl.tile_size.pmax + nl.arange(nl.tile_size.pmax)[: , None]
27
28         # Load input data from external memory to on-chip memory
29         # only read up to sz_p
30         in_tile = nl.load(in_tensor[i_p, i_f], mask=(i_p<sz_p))
31
32         # perform the computation
33         out_tile = nl.exp(in_tile)
34
35         # store the results back to external memory
36         # only write up to sz_p
37         nl.store(out_tensor[i_p, i_f], value=out_tile, mask=(i_p<sz_p))
38
39 if __name__ == "__main__":
40     from torch_xla.core import xla_model as xm
41     device = xm.xla_device()
42
43     in_tensor = torch.rand((250, 512), dtype=torch.float32).to(device=device)
44     out_tensor = torch.zeros((250, 512), dtype=torch.float32).to(device=device)
45
46     tensor_exp_kernel(in_tensor, out_tensor)
47     print(f"output_nki={out_tensor}")
48

```

Fig. 7.35: NKI source code view.

## View Neuron Profile output as JSON

As an alternative to the Neuron Profile web UI, a JSON format output is available. The JSON output includes the profile summary and all events in the execution trace. To generate the JSON output, execute the following command:

```

$ neuron-profile view --output-format json --output-file profile.json -n <path_to_neff> -
↪s <path_to_ntff>
$ cat profile.json
{
  "summary": [
    {
      "total_time": 0.017,
      "event_count": 11215,
      [...]
    }
  ],
  "instruction": [
    {
      "timestamp": 10261883214,
      "duration": 148,
      "label": "TensorMatrix",
      "opcode": "MATMUL",
      "nki_source_location": "prof-kernel.py:33",
      [...]
    },
    [...]
  ]
}

```

## See also

- [Neuron Profile User Guide](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI Performance Guide

In this document, we describe a recipe to find performance bottlenecks of NKI kernels and apply common software optimizations to address such bottlenecks. During this process, we will showcase how to leverage [neuron-profile](#), a GUI-based performance profiler designed for NeuronDevices, to guide your performance optimization efforts. Before proceeding with this document, make sure to read through [NeuronDevice Architecture Guide](#) to familiarize yourself with Neuron hardware architecture.

Ideally, performance optimization efforts would end with one of two possible outcomes: the execution of a NKI kernel is either strictly **compute-bound** or **memory-bound**. In the context of NeuronDevices, compute-bound means at least one of the compute engines is active close to 100% of the kernel execution time (90%+ is considered good in practice), while memory-bound typically means the achieved device memory bandwidth utilization (MBU) is close to 100% (60%+ is considered good in practice). For compute-bound kernels that are matrix-multiplication dominated, we should also aim for close to 100% model flops utilization (MFU) in the execution. All of these metrics are available under the **Summary** tab in [neuron-profile](#) GUI:

|               |                               |             |
|---------------|-------------------------------|-------------|
| data_movement | dma_active_cycles             | 5350880     |
|               | dma_active_time               | 3.82 ms     |
|               | dma_active_time_percent       | 54.37 %     |
|               | dma_packet_time               | 50.07 ms    |
|               | dma_queue_count               | 33          |
|               | dma_transfer_average_bytes    | 402.286 KiB |
|               | dma_transfer_count            | 1792        |
|               | dma_transfer_time             | 7.99 ms     |
|               | hbm_read_bytes                | 640 MiB     |
|               | hbm_write_bytes               | 64 MiB      |
|               | input_queue_bytes             | 704 MiB     |
|               | inputs_and_weights_size_bytes | 192 MiB     |
|               | mbu_estimated_percent         | 25.61 %     |
|               | mbu_min_read_util_percent     | 6.00 %      |

Fig. 7.36: MBU metric in neuron-profile.

The rest of this document is divided into three sections, focusing on three categories of performance optimizations. The first section covers optimizations to maximize achieved arithmetic intensity, with the goal of minimizing compute engine idle periods due to unnecessary data movement. The second and third sections dive into optimizations to improve compute engine and data movement efficiency, respectively.



|               |                                      |           |
|---------------|--------------------------------------|-----------|
| tensor_engine | hfu_estimated_percent                | 30.24 %   |
|               | matmul_instruction_count             | 292964    |
|               | mfu_estimated_percent                | 25.55 %   |
|               | mfu_max_achievable_estimated_percent | 71.16 %   |
|               | tensor_engine_active_time            | 31.31 ms  |
|               | tensor_engine_active_time_percent    | 46.37 %   |
|               | tensor_engine_instruction_count      | 590227    |
|               | tensor_engine_instruction_time       | 132.55 ms |
| vector_engine | vector_engine_active_time            | 20.79 ms  |
|               | vector_engine_active_time_percent    | 30.78 %   |
|               | vector_engine_instruction_count      | 72328     |
|               | vector_engine_instruction_time       | 33.12 ms  |
| scalar_engine | activate_instruction_count           | 42212     |
|               | activate_instruction_time            | 30.88 ms  |
|               | scalar_engine_active_time            | 20.89 ms  |
|               | scalar_engine_active_time_percent    | 30.94 %   |
|               | scalar_engine_instruction_count      | 60895     |
|               | scalar_engine_instruction_time       | 31.76 ms  |
| sync_engine   | sync_engine_active_time              | 0.84 ms   |
|               | sync_engine_active_time_percent      | 1.25 %    |
|               | sync_engine_instruction_count        | 38204     |
|               | sync_engine_instruction_time         | 0.86 ms   |
| gpsimd_engine | gpsimd_engine_active_time            | 2.15 ms   |
|               | gpsimd_engine_active_time_percent    | 3.19 %    |
|               | gpsimd_engine_instruction_count      | 27243     |
|               | gpsimd_engine_instruction_time       | 2.58 ms   |

Fig. 7.37: Compute-related metrics in neuron-profile.

## Improving Arithmetic Intensity

Arithmetic intensity of a computation workload is commonly defined as the average number of computation operations performed per byte of data accessed from memory. In the context of NeuronDevices, the definition refers to data accessed from *device memory* (HBM), since the on-chip memory (SBUF) has sufficient bandwidth to keep all compute engines busy.

When arithmetic intensity is overly low, compute engines would be consuming data much faster than DMA engines fetching data from device memory into the on-chip memory SBUF. In this case, the execution is bounded by the available device memory bandwidth. Once arithmetic intensity is beyond certain threshold, that is, ratio of maximum compute throughput over memory bandwidth, the performance bottleneck shifts to how fast compute engines can perform computation, which leads to a compute-bound execution.

Figure below visualizes the [Roofline Model](#), which captures this idea by plotting the projected attainable compute throughput with respect to the arithmetic intensity of an algorithm.

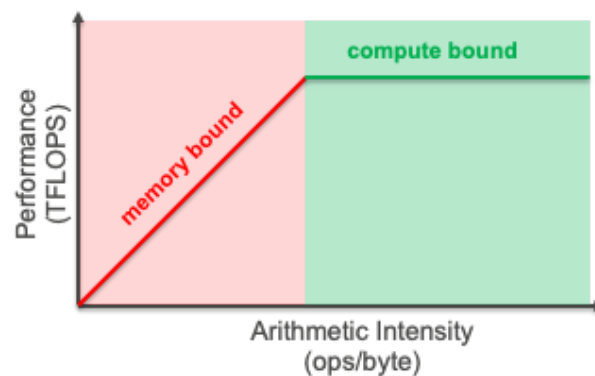


Fig. 7.38: The Roofline Model.

*Algorithmic* arithmetic intensity is an intrinsic characteristic of the particular workload and solely dependent on the compute algorithm. In reality, due to limited capacity in SBUF, the *achieved* arithmetic intensity of a NKI kernel implementation of such workload could be lower than the algorithmic arithmetic intensity. This could lead to excessive compute engine idle time blocked by completion of data movements. The two typical reasons behind this are *input data reloading* and *intermediate data spillage*. Let's discuss how to identify their symptoms in `neuron-profile` and how to mitigate these issues to improve arithmetic intensity next.

### Opt #1. Exploit temporal locality to minimize input data reloading

**Symptom:** In `neuron-profile`, if a NKI kernel triggers DMAs (`n1.load`) for the same input tensor multiple times, you would see the relevant DMA activities (on the timeline row with a label starting with `q` and ending with `I0`) being highlighted in an orange box. Hovering over the “+” sign of the box in top-left corner, a performance warning pop-up will show up, indicating which input tensor is being reloaded, the size of it and how many times it was reloaded. For example, figure below is a screenshot of such warning pop-up showing the `u` input tensor defined in my NKI kernel was reloaded ~7 times:

**Optimization:** Input tensor reloading could be avoided if the same data stay in SBUF across all the operations that consume it at different points of the execution. However, keeping too much data in SBUF across operations can increase the memory pressure in SBUF, leading to more spilling of intermediate data. Therefore, avoiding input reload should be a trade-off programmers need to make carefully. Figure below illustrates this trade-off conceptually.

A classic example of using this optimization technique is in a matrix multiplication kernel, where we need to exploit data reuse in the same rows of the left hand-side input matrix across different columns of the right hand-side matrix.

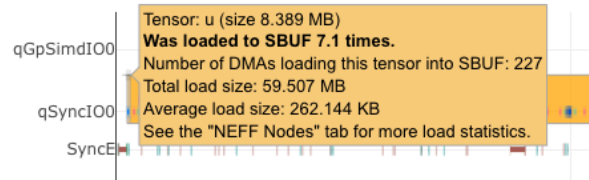
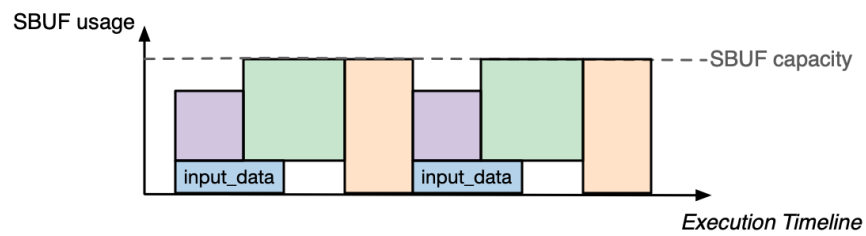
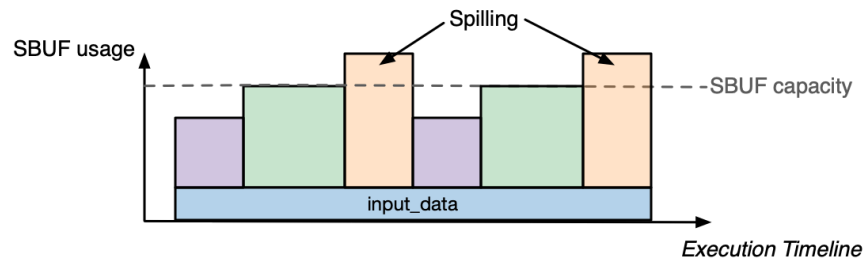


Fig. 7.39: Performance warning on input data reloading.



(a) Input Reloading



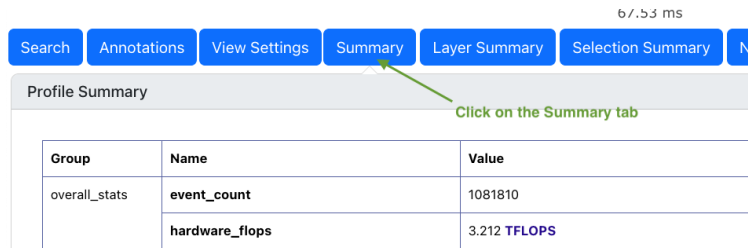
(a) No Input Reloading

Fig. 7.40: SBUF usage impact with and without input reloading.

See [Matmul NKI Tutorial Optimization 1-3](#) for more detailed discussion. Another great example is in the [Fused Mamba](#) kernel tutorial, where programmers can minimize reloading of largest input tensors through loop reordering.

## Opt #2. Fuse operations to minimize intermediate data spilling

**Symptom:** In neuron-profile , we can find many useful data movement related metrics in the Summary tab:



The screenshot shows the 'Summary' tab selected in the neuron-profile interface. A green arrow points to the 'Summary' tab with the text 'Click on the Summary tab'. Below the tabs, a table titled 'Profile Summary' displays overall statistics.

| Group         | Name           | Value        |
|---------------|----------------|--------------|
| overall_stats | event_count    | 1081810      |
|               | hardware_flops | 3.212 TFLOPS |

Fig. 7.41: neuron-profile Summary tab.

Below we highlight four relevant metrics to assess severity of data spilling under the `data_movement` section (tip: hovering over any metric name will show a detailed description of the metric):

Here, `spill_save_bytes` refers to the total size of intermediate data in bytes the workload spills from SBUF into device memory, while `spill_reload_bytes` indicates total size of spilled data in bytes the workload reloads back into SBUF. By comparing `spill_save_bytes` against `sb_read_bytes`, you can get a feel on how much of the data movement traffic from SBUF to device memory is related to spilling. Similarly, comparing `spill_reload_bytes` against `sb_write_bytes` indicates how much of traffic from device memory back to SBUF is related to spilling. If the spill related traffic takes up a significant portion (for example over 30%), it is likely worthwhile to take a close look at this optimization.

**Optimization:** To reduce spilling, the key is to find operator fusion opportunities in the kernel. To achieve fusion, we typically also need to slice up computation of each operator and perform computation for a portion of the input tensor at a time. As a simple example, assume a chain of operators `op0`  $\rightarrow$  `op1` on a large input tensor `kernel_in_hbm` that cannot fit in SBUF all at once. If we were to do the operators one at a time, we will effectively have the following sequence of events:

```
for tile in kernel_in_hbm:
    tile_sbuf = load(tile)
    op0_out_sbuf = op0(tile_sbuf)
    # compiler generated spilling, or NKI programmers explicitly perform a store
    spill_save(op0_out_sbuf, op0_out_hbm)

for tile in op1_out_device_memory:
    tile_sbuf = spill_reload(tile)
    op1_out_sbuf = op1(tile_sbuf)
    store(op1_out_sbuf, kernel_out_hbm)
```

However, if we fuse the operators from above:

```
for tile in kernel_in_hbm:
    tile_sbuf = load(tile)
    op0_out_sbuf = op0(tile_sbuf)
    op1_out_sbuf = op1(op0_out_sbuf)
    store(op1_out_sbuf, kernel_out_hbm)
```

Inside a NKI kernel, operator fusion is exactly done as the above through explicit loop fusion.

|               |   |             |
|---------------|---|-------------|
| data_movement | dma_active_cycles   | 70013528    |
|               | dma_active_time   | 58.34 ms    |
|               | dma_active_time_percent   | 86.39 %     |
|               | dma_packet_time   | 1.34 s      |
|               | dma_queue_count   | 1041        |
|               | dma_transfer_average_bytes  | 377.19 KiB  |
|               | dma_transfer_count  | 42186       |
|               | dma_transfer_time   | 250.62 ms   |
|               | hbm_read_bytes  | 11.387 GiB  |
|               | hbm_write_bytes   | 7.133 GiB   |
|               | input_queue_bytes   | 4.563 GiB   |
|               | inputs_and_weights_size_bytes   | 1.524 GiB   |
|               | mbu_estimated_percent   | 35.91 %     |
|               | mbu_min_read_util_percent   | 2.95 %      |
|               | output_queue_bytes  | 0 B         |
|               | psum_read_bytes   | 55.88 MiB   |
|               | Total size of all reads from the State Buffer. This includes DMAs reading from and instructions with input from the State Buffer. | 45.946 MiB  |
|               |   | 30458       |
|               |   | 205.569 MiB |
|               | sbuf_read_bytes   | 4.894 GiB   |
|               | sbuf_write_bytes  | 8.121 GiB   |
|               | spill_reload_bytes  | 4.636 GiB   |
|               | spill_save_bytes  | 2.572 GiB   |
|               | weight_queue_bytes  | 10.612 GiB  |
|               | weight_size_bytes   | 23.284 MiB  |

Fig. 7.42: Data movement metrics

One great use of this optimization is the self attention operator commonly found in Transformer models. Self attention performs a chain of operators: `matmul_0`  $\rightarrow$  `softmax`  $\rightarrow$  `matmul_1`, where `matmul_0` of a single attention head produces a large intermediate tensor shape that overflows SBUF in common Transformer models with a context length in the thousands. See [Fused Attention Tutorial](#) for more detailed discussion.

**Optimization Gotchas:** Certain code patterns in NKI might lead to unexpected spilling from programmers' perspectives. We are working on improving these in future releases. As an example, buffers sometimes need to be declared within the inner loop to avoid spilling. In other words, instead of:

```
buf = nl.ndarray((2, 4, nl.par_dim(128), 512), buffer=nl.sbuf)
for i0 in nl.affine_range(2):
    for i1 in nl.affine_range(4):
        buf[i0, i1, ....] = nl.load(...)
    ...
```

we need to implement:

```
for i0 in nl.affine_range(2):
    for i1 in nl.affine_range(4):
        buf = nl.ndarray((nl.par_dim(128), 512), buffer=nl.sbuf)
        buf[...] = nl.load(...)
```

With the above aforementioned optimizations, the kernel execution should achieve an arithmetic intensity that is somewhat close to the algorithmic arithmetic intensity. At this point, you should be able to observe from the execution timeline in `neuron-profile` whether the kernel spends more time in compute or DMA engines. The `engine/dma_active_time_percent` metrics reported in the Summary tab should also give you good hints. If your kernel execution is dominated by computation, we recommend going over [Optimizing Compute Efficiency](#) first to optimize compute efficiency. Otherwise, jump straight to [Optimizing Data Movement Efficiency](#) to understand how to optimize data movement efficiency.

## Optimizing Compute Efficiency

Compute efficiency optimizations typically fall into two categories:

1. “time” domain engine utilization: reduce engine idle time to keep the compute engine *on critical path* as busy as possible, such as enabling pipelining among engines.
2. “spatial” domain engine utilization: within the engine active periods, increase instruction efficiency to use as many hardware units within the engine as possible, such as combining multiple instructions into one.

Let’s dive into each category below.

### Reducing engine idle time

To improve the active time of a compute engine, we need to understand the exact reasons for the engine to enter an idle state. In `neuron-profile`, we can focus on the execution trace of the bottlenecked engine and zoom into the visually large engine idle gaps. For example, in the below profile, we expect VectorE to be the bottlenecked engine and therefore focus on the idle gaps on VectorE:

*Side note*, for faster GUI rendering, `neuron-profile` enables data sampling by default and “hides” certain instructions from the timeline with a large profile. To confirm whether an engine indeed has an idle gap, we recommend zooming into a smaller region of the profile and turn on “Show unsampled data” in `View Edit Settings` to make sure all instructions are rendered:



Fig. 7.43: Engine idle gaps.

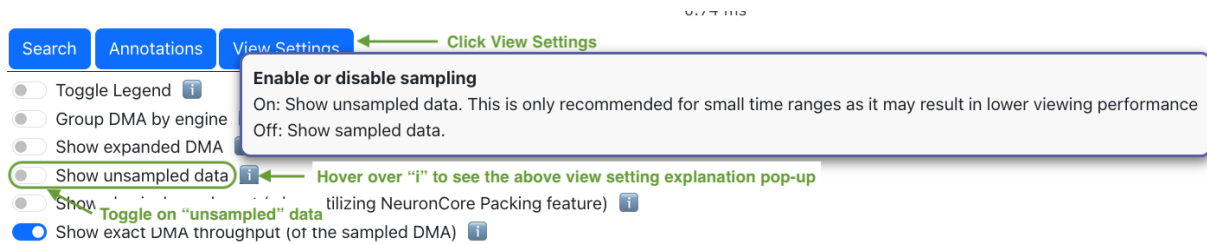


Fig. 7.44: Show unsampled data in neuron-profile.

For each engine idle gap, you can find out the reasons why the engine cannot execute instructions by inspecting the **semaphore wait condition** of the first instruction executed on the engine after the gap. Broadly speaking, these semaphore wait conditions are either waiting for 1) other compute engine instructions or 2) DMA activities to finish. We have different techniques to shrink the idle gaps caused by either of these wait conditions (that is, engine stall reasons).

### Opt #3. Overlap execution across compute engines through pipelining

**Symptom:** The semaphore wait condition of the first instruction after an idle gap is on a semaphore name that matches a compute engine name in NeuronCore: Vector, Scalar, GpSimd and Tensor. These semaphores are associated with instruction completion on the corresponding compute engine.

For example, the below TENSOR\_TENSOR instruction on VectorE is waiting for S[4] (Scalar) to reach a value of 36. This means VectorE was waiting for ScalarE to finish certain instructions.

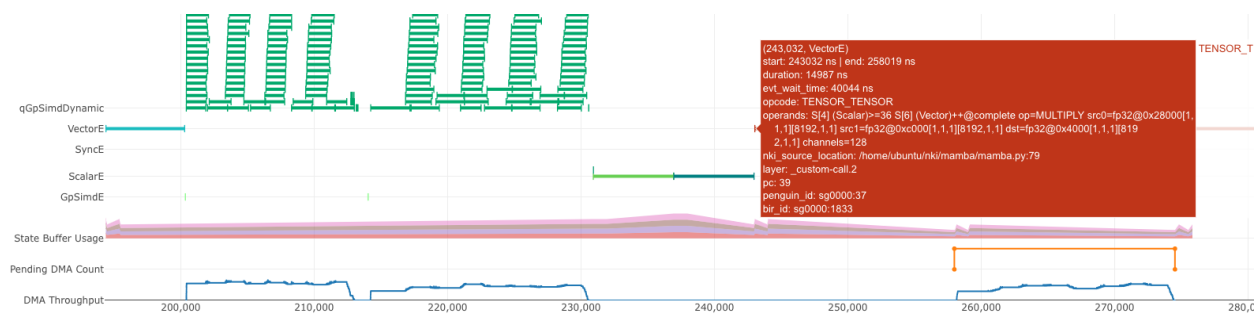


Fig. 7.45: Semaphore wait on another compute engine.

**Optimization:** When there is a sequence of operators on different compute engines, we can slice the computation in a way that the compute engines can process tiles of the original operator in a pipeline fashion. As an example, let's assume we have two operator back to back on a large (say, thousands of elements) tensor  $X$ :  $X \rightarrow \text{op0} \rightarrow Y \rightarrow \text{op1} \rightarrow Z$ .  $\text{op0}$  is performed on ScalarE while  $\text{op1}$  is on VectorE. For simplicity, let's assume tensor  $X/Y/Z$  have the same shape.

Figure below shows two possible execution timelines with and without engine pipelining. Without pipelining, VectorE is fully idle when ScalarE is executing  $\text{op0}$  on tensor  $X$  in the first half of the execution. Similarly, ScalarE is idle while VectorE is running  $\text{op1}$ . However, with pipelining, ScalarE is able to produce partial results in tiles and unblock VectorE as soon as the first tile is processed. Overall, engine pipelining shortens the end to end latency to complete  $\text{op0}$  and  $\text{op1}$ , through shrinking engine idle time and improving hardware utilization.

Choosing a proper tile size is crucial to the performance of such engine pipelining. It is up to NKI programmers to make this choice in kernel implementation and iterate on it using performance profiling data in neuron-profile. For complex kernels, we often need to schedule a pipeline among all engines: Tensor/Scalar/Vector/GpSimd Engine.

For example, in Transformer's self-attention layer, in addition to fusing  $\text{matmul\_0}(Q, K) \rightarrow \text{softmax} \rightarrow \text{matmul\_1}(\text{softmax\_out}, V)$  in a single kernel to minimize spilling as discussed in [Opt #2](#), we also need to form a complex engine pipeline for the operators to maximize utilization of the compute engines:

- $\text{matmul\_0/matmul\_1}$ : TensorE
- softmax:
  - exponential: ScalarE
  - summation: VectorE
  - scale by reciprocal of summation: ScalarE



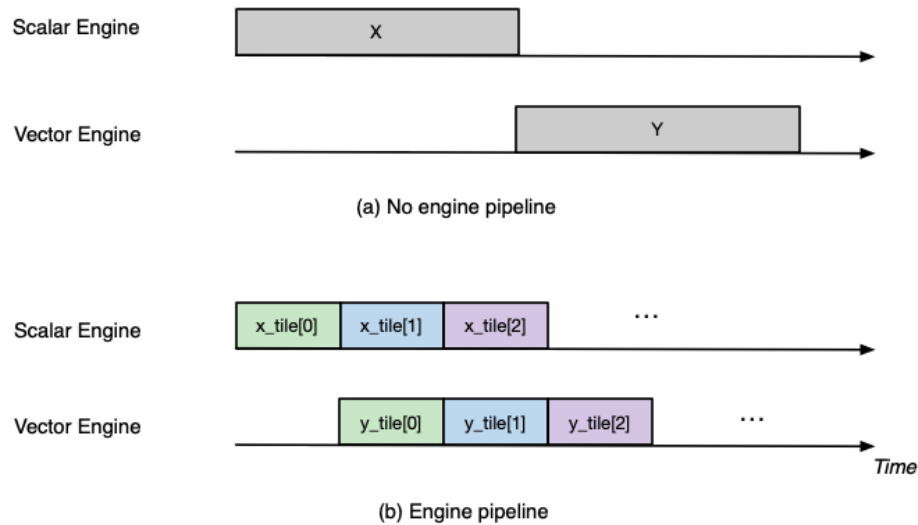


Fig. 7.46: Engine timeline with and without engine pipelining.

- for causal self attention, triangular masking: GpSimdE

See [Fused Self Attention](#) tutorial for more detailed discussion.

#### Opt #4. Overlap data loading with computation

**Symptom:** The semaphore wait condition of the first instruction after an idle gap is on a semaphore name that starts with letter q. These semaphores are associated with completion of DMA activities.

For example, hovering on an instruction will bring up the key instruction details as follows:

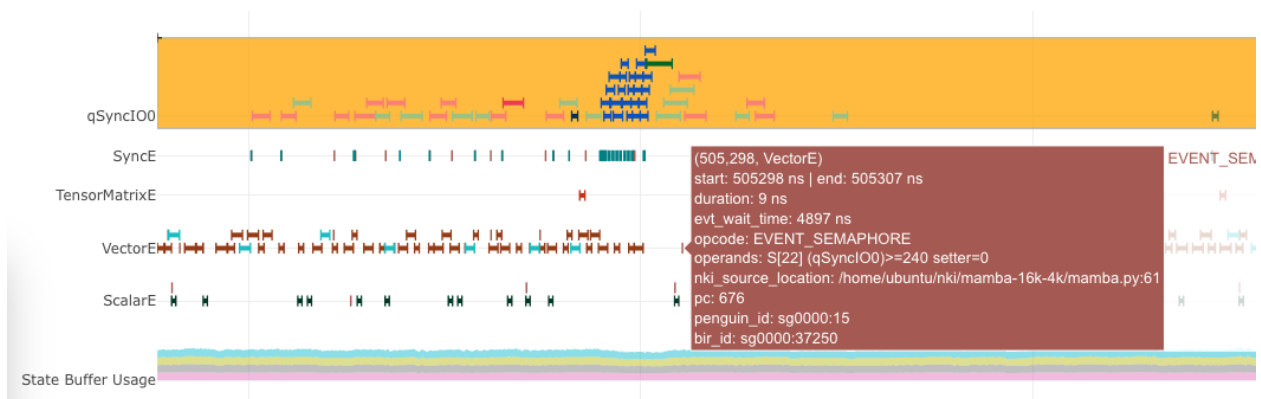


Fig. 7.47: Instruction waiting for input data loading.

In this particular screenshot, the `EVENT_SEMAPHORE` instruction could not start earlier even though `VectorE` was idle because it was waiting for semaphore `S[22]` (`qSyncIO0`) to reach a value of 240. The semaphore is only incremented whenever the corresponding DMA activities shown on the `qSyncIO0` execution trace are completed. Clicking on the DMA activities on `qSyncIO0` immediately before the `EVENT_SEMAPHORE` instruction, you may follow the `nki_source_location` to find out which line of code is related to this DMA activity (`n1.load()` call).

Similarly, if an instruction is blocked on S[47] (qSyncSpillReload0), that means it is blocked by DMA activities for spilling:

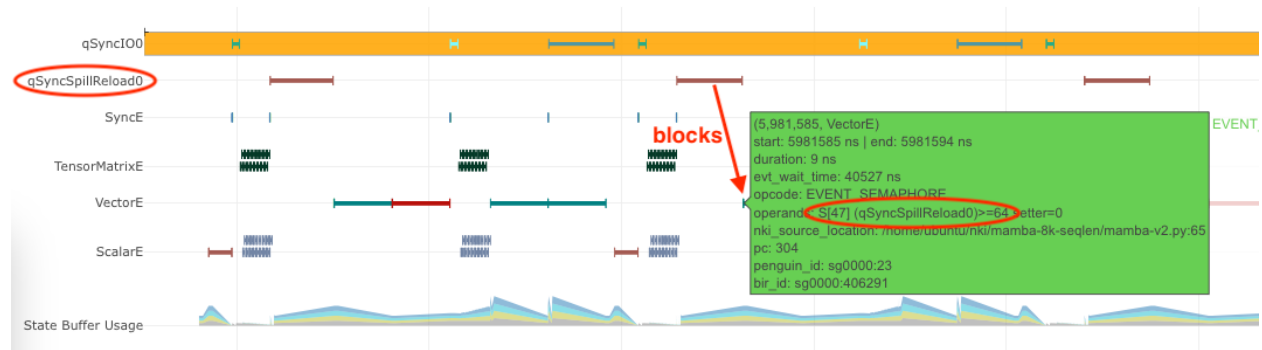


Fig. 7.48: Instruction waiting for spilled data reloading.

Clicking on the DMA activities on qSyncSpillReload0 immediately before the EVENT\_SEMAPHORE instruction, you may find out the name of the intermediate NKI tensor that was spilled/reloaded. For example, the below DMA transfer reloads the tensor named `deltaU` as defined in our NKI kernel. Note, spill/reload DMA transfers are generated by Neuron Compiler automatically by analyzing SBUF usage in NKI kernels. Therefore, these DMA transfers do not have an associated explicit NKI API call or `nki_source_location` information.

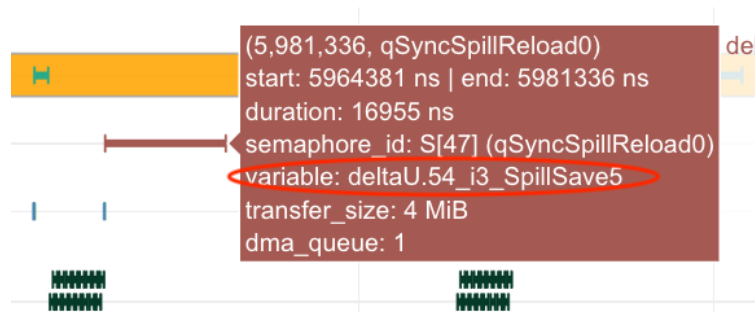


Fig. 7.49: Spilled tensor variable name.

**Optimization:** Overlapping data loading with compute is highly similar to enabling compute engine pipelining in Opt #3, since DMA engines can move data in parallel to compute engine execution, just like how compute engines can run different operators in parallel.

However, it is also possible that even after maximizing overlapping of compute and data movement the best you can, the data movement duration is still not hidden behind compute even though your kernel has a compute-bound arithmetic intensity. In these cases, the most common cause is the data movement in your kernel is not using the DMA engines *efficiently*. Refer to a [later section](#) to see relevant optimization techniques to improve DMA bandwidth utilization.

As a concrete example, we demonstrate how to properly overlap compute and data movement in a compute-bound (VectorE as the bottlenecked engine) kernel in [Mamba tutorial](#).

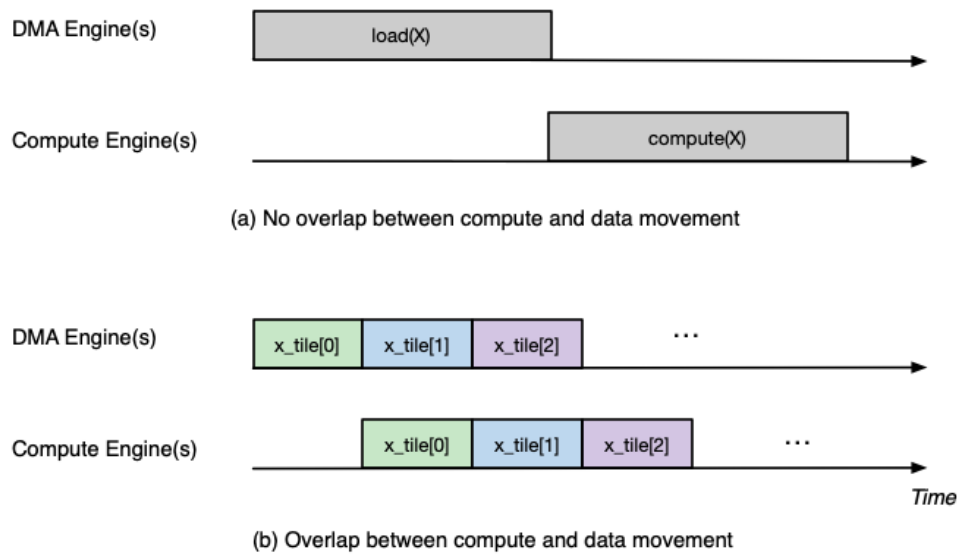


Fig. 7.50: DMA and engine timeline with and without overlapping.

## Improving engine efficiency

Once done with “avoiding engine idle gaps” as much as possible, we can focus on improving “engine efficiency” during the busy periods of the engine. We will start with two optimizations techniques that are generally applicable to all compute engines, followed by TensorE-specific optimization techniques.

### Opt #5a: Use sufficiently large input tiles in free dimension

**Symptom:** Certain operators might trigger many back-to-back instructions with small free dimension sizes in the input tensors. For example, in the below profile, ScalarE is busy with many repeated `activation` instructions with `IDENTITY` (scale/bias enabled) activation function, which is equivalent to calling `nki.isa.tensor_scalar(op0=n1.multiply, op1=add)` APIs. If you click on one of the instructions to pull up the instruction detailed view, you can see the source tensor access pattern is `fp32@20580[1, 1, 1][1, 1, 1]`, where the first set of bracket indicates 3D strides and the second set indicates 3D shape in FP32 elements. More detailed discussion of ISA access pattern can be found by clicking on the `i` button at the end of the `Operands` row.

In this example, each of the back-to-back instructions is reading **one** element per partition from SBUF, which would take about one engine cycle to perform useful computation within the instruction. Such instructions are extremely inefficient since the static instruction overhead in the order of ~100 cycles would be limiting the overall throughput.

To make things worse, these instructions also have data dependency (read after write) between consecutive instructions, which means the next instruction cannot start data read until the previous instruction has all of its output committed to the local SRAM. In neuron-profile, you can inspect data dependency between instructions by clicking on an instruction of interests (Inst1 in the below profile), which will highlight the clicked instruction and also the instruction that produces input for the clicked instruction (Inst0 in the below profile). The dependency information can also be viewed in the details “instruction dependency pcs”. In fact, all the neighboring instructions also have a similar dependency patterns in this profile.

With the above inefficiencies, the initiation interval (the time between the starting points of two consecutive instructions) for these instructions on ScalarE is around 189 ns (264 ScalarE cycles on NC-v2), which is much higher than the useful computation cost (one ScalarE cycle throughput-wise).

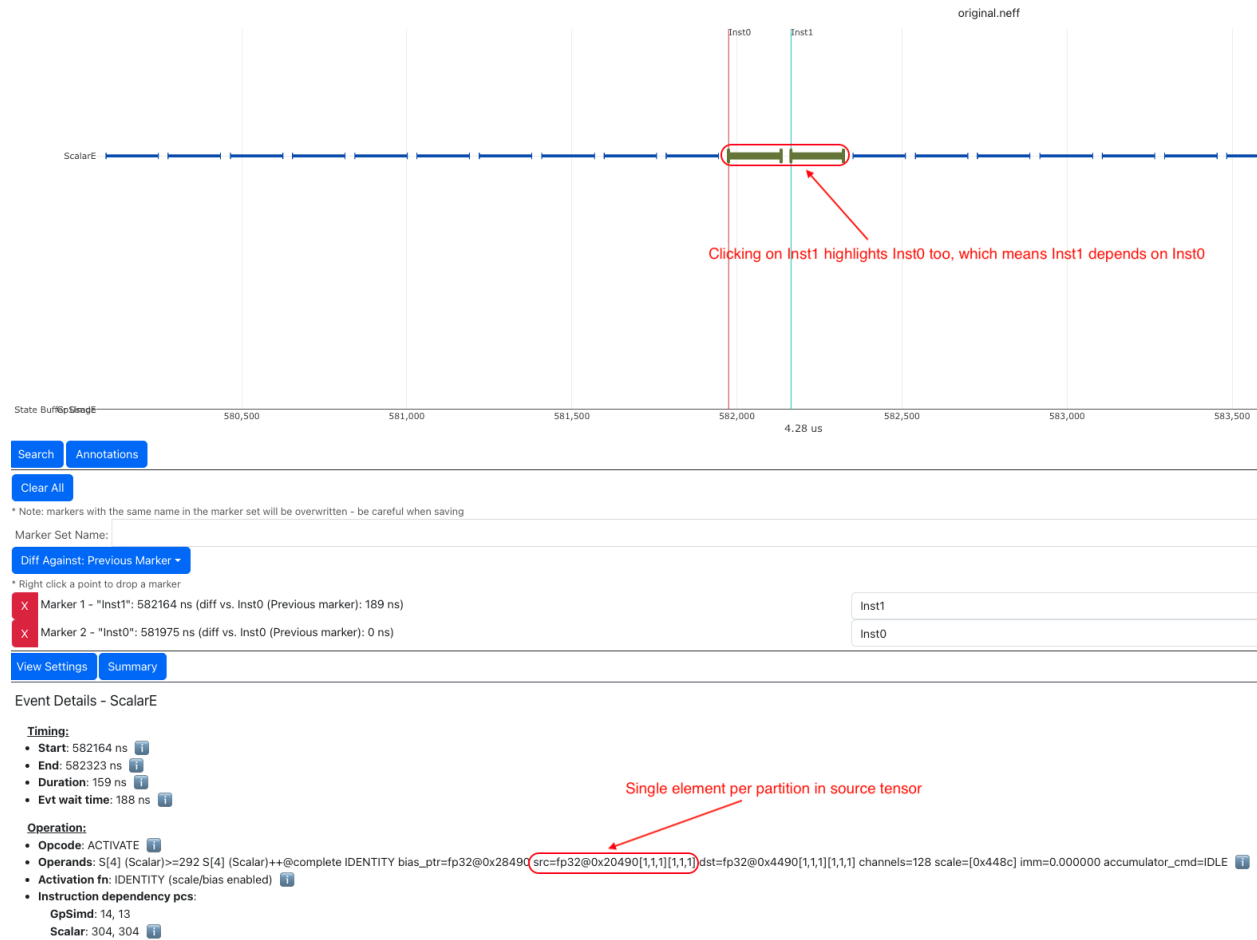


Fig. 7.51: Many back-to-back ScalarE instructions with small tensor shapes

**Optimization:** The trick of this optimization is to increase the free dimension size of instruction input tiles. As discussed in the [architecture guide](#), NeuronCore compute engines typically require at least 128 elements/partition in the source tensor to be efficient. However, it is worth mentioning that increasing free dimension sizes might not be trivial due to the high-level computation definition. We suggest developers walking through the [architecture guide](#) in detail to better understand capabilities of different compute engines, and mapping/reformulating the high-level operators onto the engines using the most suitable instructions. Such instructions could be invoked either through the high-level `[nki.language](api/nki.language)` or low-level `[nki.isa](api/nki.isa)` APIs.

In addition, keep in mind there is a trade-off in choosing the free dimension size in instruction input tiles: Too small of a tile size exposes significant instruction overhead leading to inefficient engine execution, while too large of a tile size often leads to inefficient pipelining between engines (working against [Opt #3](#)) and high memory pressure in SBUF (working against [Opt #2](#)).

As an example, a naive implementation of the prefix sum scan operation in Mamba v1 would trigger `seq_len` back-to-back single element `nki.isa.tensor_scalar` instructions as shown in the above profile example, where `seq_len` is the sequence length of the model typically in the range of thousands. A more efficient way to implement this operation is through a special VectorE instruction `nisa.tensor_tensor_scan`. See the [Mamba tutorial <tutorials/fused\\_mamba>](#) for more discussion.

### Opt #5b: Use sufficiently large input tiles in partition dimension

**Symptom:** When instructions use input/output tiles that span fewer than 128 partitions, they typically under-utilize the compute engine capabilities. This is because each SBUF/PSUM partition has a one-to-one mapping to parallel vector lanes in the compute engines. As an example, the `TENSOR_TENSOR` instruction (equivalent to `nki.tensor_tensor`) on VectorE takes a source tensor in SBUF that occupies 64 partitions only, as indicated by the `channels=64` instruction operand field. If we were to increase the `channels` field to 128, the instruction would have taken the same amount of time as `channels=64`.



Fig. 7.52: An instruction that read/write less than 128 partitions.

Similarly, for a `MultiplyMoving` instruction (Matmul opcode in neuron-profile) TensorE, if the instruction reads/writes tiles do not span the full SBUF/PSUM partitions, we would be underutilizing TensorE. As an example, the below `MultiplyMoving` instruction only writes to 96 partitions in PSUM, as indicated by the operand `128*96`, which means the instruction only uses 128 rows and 96 columns of the processing elements out of the available `128x128` systolic array.

**Optimization:** If we see **many back-to-back instructions** on the compute engine that have fewer than 128 partitions in the input/output tiles as discussed above, we should consider an optimization called “partition vectorization”.

As an example, say we have two `nki.isa.nc_matmul()` instructions with each generating a 64-partition PSUM tile of the same shape. Then VectorE needs to run `nki.isa.tensor_reduce()` on both tiles to generate a reduction result. Note, on `trn1/inf2`, VectorE cannot run the two independent `nki.isa.tensor_reduce()` instructions in parallel in this case, even though the total number of compute lanes required for these instructions does not exceed 128. To improve VectorE utilization in this case, we can:

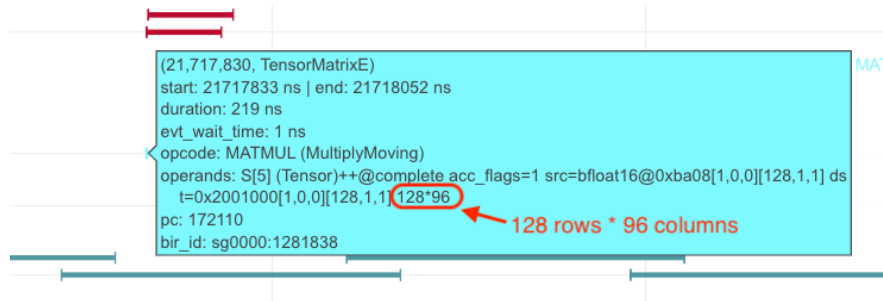


Fig. 7.53: MultiplyMoving instruction that uses <128 TensorE columns

1. The two `nc_matmul()` instructions write to disjoint PSUM partitions: partition 0-63 for the first `nc_matmul` and partition 64-127 for the second one.
2. Invoke a single `nki.isa.tensor_reduce()` instruction to process output of both `nki.isa.nc_matmul()` instructions.

The below pseudo-code illustrates the above computation without and with partition vectorization.

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl

#####
# option 1: No partition vectorization
# two 64-partition vector instructions running serially

# By default, NKI creates mm_tile0 and mm_tile1 in partition 0-63
mm_tile0 = nisa.nc_matmul(...)
mm_tile1 = nisa.nc_matmul(...)

# Both nki.isa.reduce instructions move data from psum partition 0-63
# in a serialized fashion
reduce0 = nisa.tensor_reduce(mm_tile0, ...)
reduce1 = nisa.tensor_reduce(mm_tile1, ...)

#####
# option 2: Partition vectorization
# vectorized into one 128-partition vector instructions

# Here, we explicitly declare a 128-partition tensor in PSUM
mm_tile = nl.zeros((128, ...), np.float32, buffer=nl.psum)

i_output0_p = nl.arange(64)[: , None]
i_output1_p = 64 + nl.arange(64)[: , None]
# Assign first part of mm_tile to partition 0-63
mm_tile[i_output0_p, ...] = nki.isa.nc_matmul(...)
# Assign second part of mm_tile to partition 64-127
mm_tile[i_output1_p, ...] = nki.isa.nc_matmul(...)

# A single nki.isa.reduce instruction, using all 128 partitions
reduce = nisa.tensor_reduce(mm_tile, ...)
```

Option #2 above is able to perform the reduction 2x faster, by vectorizing the partition dimension and performing a

single reduction instead of two.

## Opt #6: Combine instructions

**Symptom:** Even though the majority of popular ML models are matrix multiplication heavy, certain operators can be vector/scalar operation heavy instead, such as self-attention in Transformer models. These operators typically have a performance bottleneck in VectorE or ScalarE or both. As an example, the below profile shows the inner loop of self attention, where either VectorE or ScalarE is busy at any moment in time, while TensorE has clear engine idle gaps.

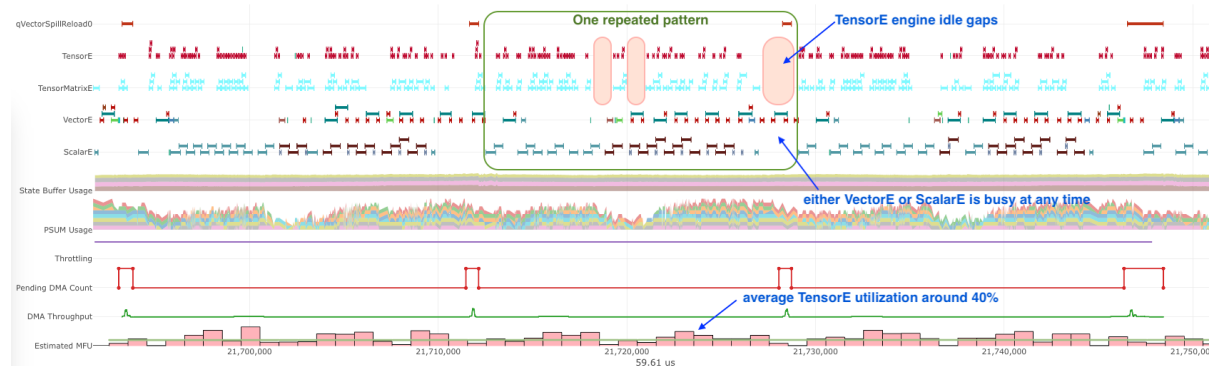


Fig. 7.54: A VectorE/ScalarE-bound profile.

**Optimization:** A common optimization to tackle vector/scalar-operation-heavy operators is **combining instructions** using low-level `nki.isa` APIs. Combining instructions can leverage the deep pipelined stages within VectorE and ScalarE engine data path to increase hardware utilization per instruction and reduce the instruction count. Check out the [architecture guide](#) to learn what operations can be done in a pipeline fashion in a single VectorE/ScalarE instruction.

For example, below pseudo-code showcase combining three instructions into a single one on ScalarE. `impl 1` and `impl 2` are functionally equivalent, but `impl 2` is 3x faster in terms of latency by touching the input data only once and running all three operations (multiply, add, exp) in a pipeline.

```
import neuronxcc.nki.isa as nisa
import neuronxcc.nki.language as nl

# input: data (tile[128, 512]), scale (tile[128, 1]) , bias (tile[128, 1])

# impl 1:
scaled = nl.multiply(data, scale)
shifted = nl.add(scaled, bias)
exp = nl.exp(shifted)

# impl 2:
exp = nisa.activation(nl.exp, data,
                     bias, scale)
```

Check out [nki.isa APIs](#) to understand low-level ISA API semantics, limitations, engine mapping, and rough estimates of performance cost.

See [Fused Mamba](#) tutorial for a concrete example to combine matrix-vector multiplication and exponential evaluation in a single `nisa.activation` instruction. Similarly, in [Fused Self Attention](#) tutorial, we combine the subtraction of the maximum with exponential in a single `nisa.activation` instruction in the Softmax operator.

## Opt #7: TensorE only: Leverage fast weight load

**Symptom:** Let's consider a matrix multiplication between two matrices of shape  $[M, K]$  and  $[K, N]$ , with one of the following conditions:

1.  $M$  is significantly smaller than 128, while  $N$  is much larger than 128, or
2. the other way around:  $N$  is significantly smaller than 128, while  $M$  is much larger than 128

In NKI, if the matrix with  $\min(M, N)$  dimension is mapped to the **stationary tensor** (x input tensor in `n1.matmul` and `nisa.nc_matmul`) for the TensorE LoadStationary instruction (details see [architecture guide](#)), we will typically end up under-utilizing TensorE more severely compared to mapping such matrix to the **moving tensor**.

In neuron-profile, programmers can identify also this inefficient case by inspecting the `src` access patterns for LoadStationary and MultiplyMoving instructions on TensorE. For example, the below screenshot indicates a stationary tensor with 1 element per partition and a moving tensor with 128 elements per partition:

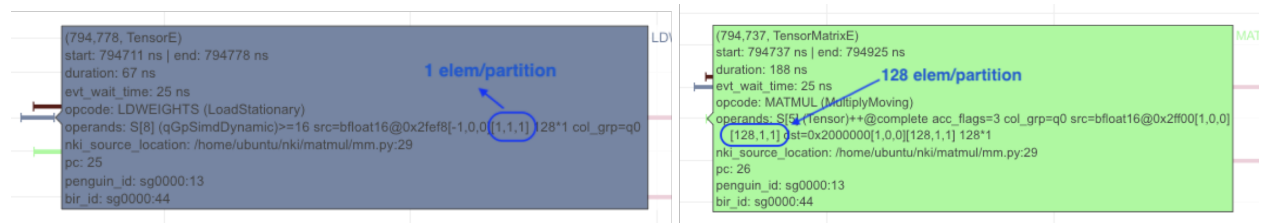


Fig. 7.55: Example instructions for matrix-vector multiplication.

If you have many back-to-back TensorE instructions with the above pattern, we recommend applying the below optimization.

**Optimization:** The key idea of this optimization is to simply swap the stationary and moving tensor positions for the given matmul in NKI, in order to leverage the “Fast LoadStationary” support in TensorE (more discussion in [architecture guide](#)). To better understand the intuition behind this, let's walk through a concrete example.

Consider a  $[1, 128] \times [128, 128]$  matrix multiplication as below:

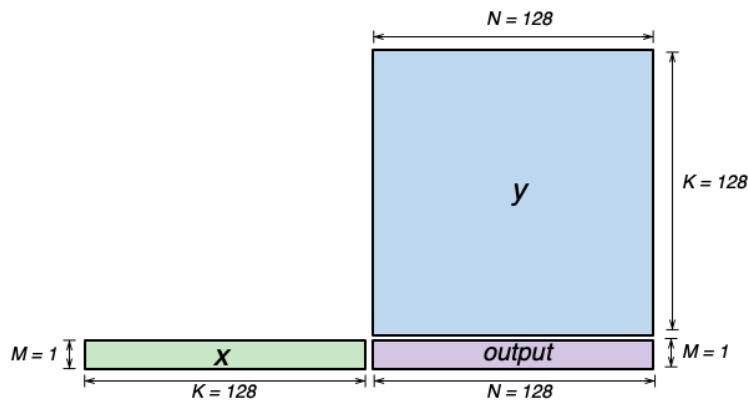


Fig. 7.56: Illustration of matrix-vector multiplication.

Since  $K=128$  is the contraction dimension, it will get mapped to the partition dimension of the SBUF for both the  $x$  and  $y$  matrices.  $M$  and  $N$  will therefore get mapped to the free dimension of the SBUF, and we will refer to  $x$  as the “short” tensor, and  $y$  as the “long” tensor (short and long in the free dimension, respectively). We have two possible ways of performing this computation on the TensorE, which we'll refer to as “Short Moving” and “Short Stationary”, depending on which tensor has the short free dimension.



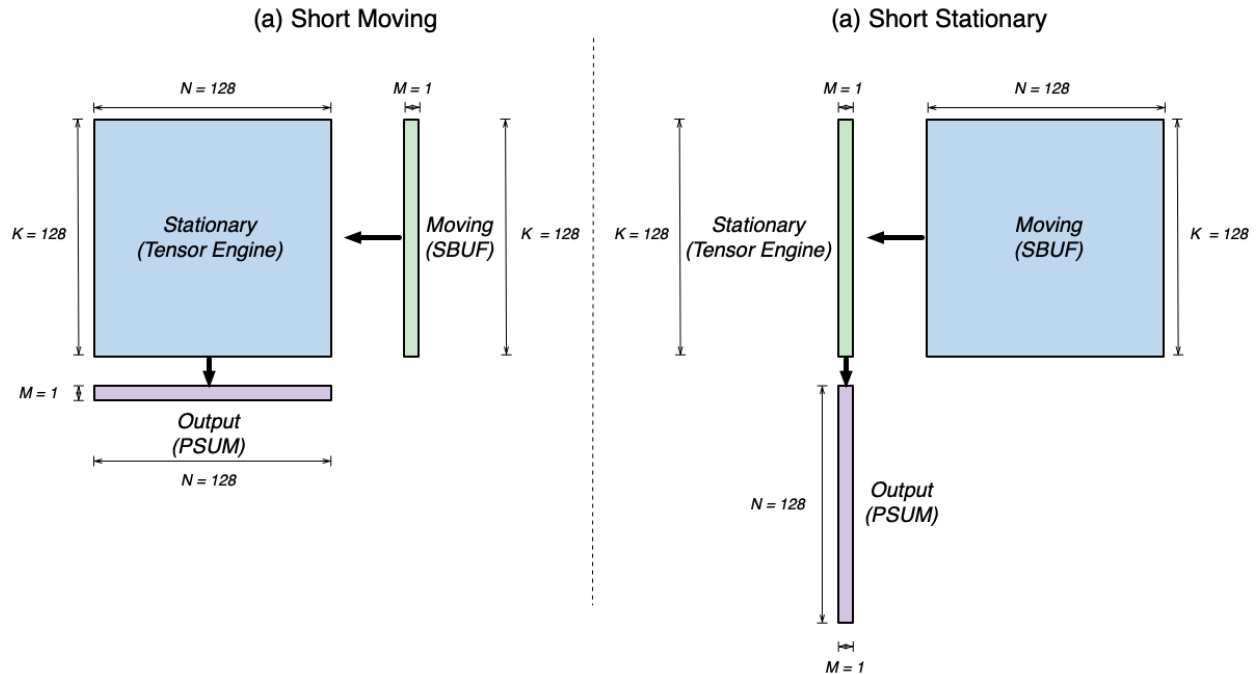


Fig. 7.57: Two possible TensorE instruction mapping for matrix-vector multiplication.

Based on the multiplication property of transpose, we have  $A \times B = (B.T \times A.T).T$ . Meanwhile, based on the semantics of TensorE, when we want to compute  $A \times B$ , we need to call `nc_matmul(A.T, B)`, and for  $B.T \times A.T$ , we need to call `nc_matmul(B.T.T, A.T) -> nc_matmul(B, A.T)`. Notice how the parameters to `nc_matmul` are swapped! Thus, when we swap stationary and moving tensors and perform the matrix multiplication, the output tensor will be transposed from the original output.

Recall, if there is a difference in initiation interval between `LoadStationary` and `MultiplyMoving`, one of them can end up limiting the throughput of TensorE:

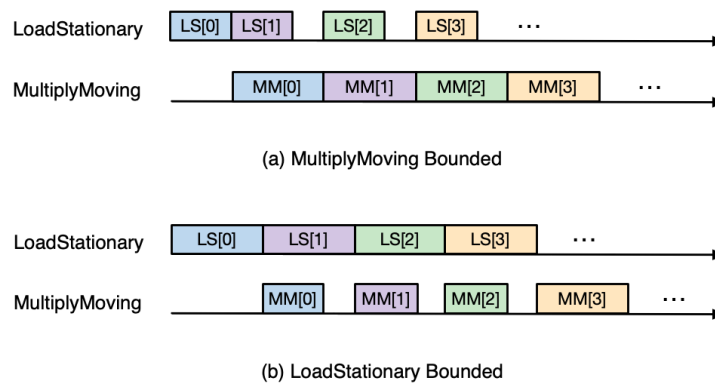


Fig. 7.58: Two possible TensorE performance characteristics.

In the above scenarios, we expect TensorE performance to be bound by whichever instruction reads the longer tensor - `LoadStationary` in “Short Moving”, and `MultiplyMoving` in “Short Stationary”. However, with TensorE Fast LoadStationary, TensorE can perform `LoadStationary` **up to 4x** faster than a `MultiplyMoving` with the same free axis size.

So in the two above scenarios:

1. Short Moving - LoadStationary initiation interval is roughly equal to the number of elements divided by 4 (because of fast LoadStationary), and MultiplyMoving initiation interval is dominated TensorE instruction turnaround time `MM_INIT_LATENCY` (64 cycles on `trn1`). Therefore, we have  $LS\_II \approx 128/4 = 32$  cycles, and  $MM\_II \approx \max(1, MM\_INIT\_LATENCY=64 \text{ cycles})$  which leads to issuing a MM roughly every 64 cycles.
2. Short Stationary - MultiplyMoving initiation interval will dominate, which leads to issuing a MM roughly every 128 cycles.

Because of the above, we will prefer to map short tensors to the moving tensor in `MultiplyMoving` instruction in `TensorE`.

A classic example is a matrix-vector product. This is commonly seen in auto-regressive token generation in LLMs, where most of the matmuls occur only on a single token (vector) as the feature map, while the weight tensor remains large and hence must be broken into tiles to meet `TensorE` tile size constraints.

## Opt #8: TensorE only: Mitigating overhead from tensor transposes

**Symptom:** Since `TensorE` accounts for over 90% of the hardware FLOPS on a `NeuronCore`, we would like the engine to perform useful computations as much as possible, especially in matmul-heavy kernels. The most common “not useful” computation that could occupy precious `TensorE` cycles is tensor PF-transposes, which swap the partition and free dimensions of a NKI tile. When you have a profile with `TensorE` visually extremely busy, we recommend doing a sanity check on how much of the `TensorE` activities are performing transposes. One easy way to check is by selecting `Instruction Type` as the `Instruction Grouping in View Settings`:

With this instruction coloring, `TensorE` instructions will be highlighted in two different colors: one for Transpose and one for Regular (useful matmuls). As an example, the below profile has an execution trace with `TensorE` being the performance bottleneck. Visually, we can see the bulk of the `TensorE` execution is for regular matmuls, but there is a noticeable chunk of engine time spent on transpose-induced instructions in red. Note, the colors for transpose versus regular instructions are chosen randomly by the profiler each time. You should hover over the instructions to check the `Instruction Type` field on the pop-up to confirm the color mapping.

**Optimization:** The key goal of this optimization is to reduce the number of transpose-induced instructions on `TensorE`, when such instructions are taking up a large portion of the execution. Before diving into techniques to reduce transposes, it is important to understand the root cause of these transposes.

At a high level, tensor transposes are needed to adjust the data layout of tensors to match the partition dimension requirements of different ISA instructions. Refer to the [architecture guide](#) for layout requirements of each compute engine. Transposes are inserted explicitly into NKI kernels through `nl.transpose` or `nisa.sb.transpose` APIs, or calling `nl.matmul` with `transpose_x=False`. These transposes are most commonly lowered down to `Tensor Engine`.

Broadly speaking, there are 2 different types of tensor transposes, with different root causes:

1. IO tensor transpose (abbreviated as IO transpose)
2. intermediate tensor transpose (abbreviated as intermediate transpose)

**IO transpose.** These transposes are \*\* done on NKI kernel IO (input/output) tensors, which must reside in device memory in current NKI releases. The transposes are needed when the NKI compute API consuming input tensors or producing the output tensors expect a different layout than their IO layout in device memory. To simplify discussion, we dive into input tensor layout discussion below, but the same reasoning also applies to output tensors.

For example, say we have an input tensor in device memory with layout `[out_channel=128, in_channel=128]` (major-to-minor ordering), but the `nisa.nc.matmul` call in our NKI kernel expects `[in_channel, out_channel]` as input tile layout. In this case, we can perform a `nl.load` to load the input into SBUF, with `out_channel` being the partition dimension because `out_channel` is the most major dimension in device memory. Then, a PF-transpose on `TensorE` is required before the loaded data can be consumed by `nisa.nc.matmul`. Alternatively, we can invoke `nl.load_transpose2d` to transpose the input tensor on the fly in the DMA engine, with a major caveat of much lower

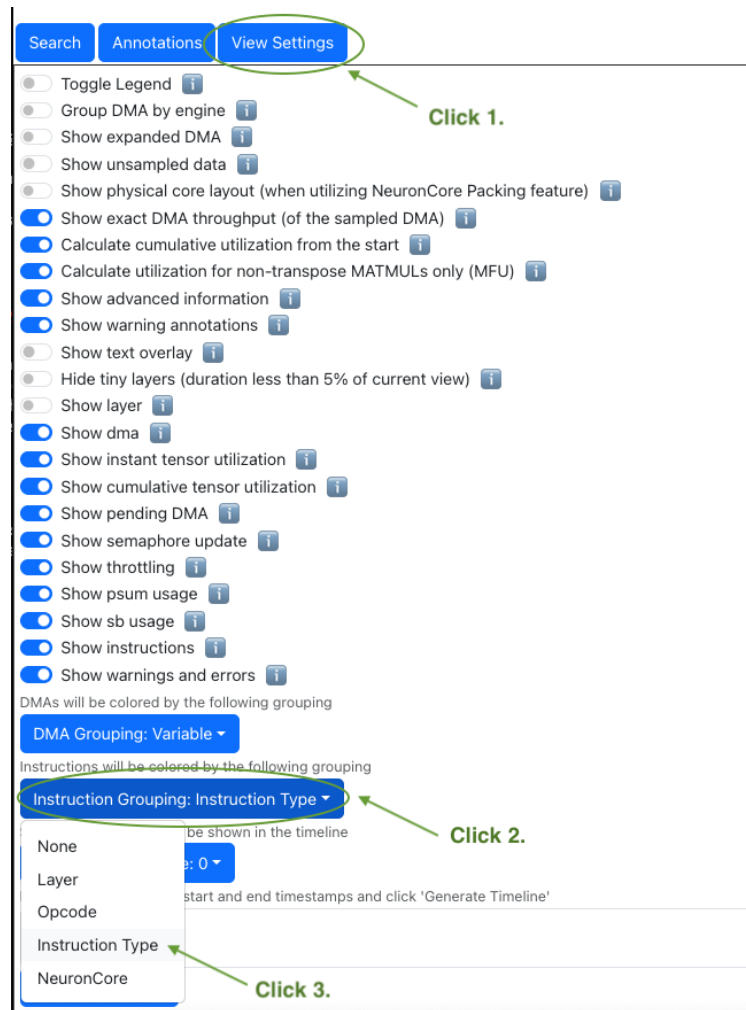


Fig. 7.59: Change view settings to visualize transposes.



Fig. 7.60: Example timeline with a transpose instruction type.

DMA bandwidth compared to `nl.load`. `nl.load_transpose2d` could make sense in a compute-bound kernel, but should certainly be avoided in memory-bound kernels.

Either way, an IO transpose is inevitable here *due to* the IO tensor layout choice we made as NKI programmers. In the naive case scenario where we only care about reaching the best performance for a single kernel, we can carefully decide on the IO tensor layout to make sure it is compatible with the NKI compute API layout requirements. When the input tensor is consumed by multiple compute APIs with conflicting layout requirements, IO-transposes cannot be avoided but should still be minimized as much as possible with a careful trade-off.

However, NKI kernels are often injected into a larger model defined at the framework such as PyTorch and JAX, in which case the kernel IO tensors are also input/output of the surrounding framework operators. These cases will require more complex reasoning on the optimal IO tensor layout for the NKI kernel, but the optimization goal of minimizing IO transposes remains the same.

One last complexity in deciding IO tensor layout is the layout choice also has a potential impact on DMA efficiency. See more discussion in a [later section](#) discussion optimizing data movement efficiency.

**Intermediate Transpose.** These transposes are done on intermediate tensors produced within a NKI kernel. These transposes arise due to layout requirement mismatches between producer and consumer NKI compute APIs.

There are two common techniques to reduce intermediate transposes: 1) swapping moving/stationary tensors in `nisa.nc_matmul` (or equivalently, `nl.matmul`) and 2) mapping a computation to an alternative engine with different layout requirements.

One example for technique 1) is in an operator chain commonly seen in Transformer models: `linear_layer` → `layernorm`. Normally, we tend to map the weight `[hidden_size, 4xhidden_size]` tensor in `linear_layer` to the stationary tensor and the input feature map `[hidden_size, seq_len]` to the moving tensor when performing `nisa.nc_matmul` on TensorE. The output feature map of this `matmul` will be in a layout of `[4xhidden_size, seq_len]`. However, the first step in `layernorm` to calculate mean and variance, `nisa.bn_stats`, requires `4xhidden_size` to be the free dimension because we need to calculate mean/variance within a single token. Therefore, a naive implementation of this operator chain will trigger a PF-transpose between the `nisa.nc_matmul` and `nisa.bn_stats` instructions. However, if we were to instead map the weight tensor to the moving tensor and input feature map to stationary tensor, we can skip this PF-transpose entirely because the `nisa.nc_matmul` output will be in the expected layout by `nisa.bn_stats`.

An example for technique 2) is in a similar operator chain: `linear_layer` → `RMSnorm` with the same intermediate tensor dimensions as the above example. `RMSnorm` is considered a cheaper normalization operator compared compared `LayerNorm`, because it replaces the mean/variance calculation with squared and summation. Unlike `nisa.bn_stats` for mean/variance calculations which must be done along the free dimension, for `RMSnorm` the scalar squared operator has no layout requirement and the summation can be done along either dimensions: use `Vec-torE nisa.tensor_reduce` for free dimension summation or use TensorE `nisa.nc_matmul` for partition dimension summation (see [TensorE alternative use case](#) in the architecture guide). Since `RMSnorm` can be done with either `[4xhidden_size, seq_len]` or `[seq_len, 4xhidden_size]`, we should make the layout choice based on more surrounding operator: `RMSnorm` in Transformer models is typically followed by yet another `linear_layer`, which requires the `[4xhidden_size, seq_len]` layout. Therefore, to minimize intermediate transposes in an operator chain like `linear_layer` → `RMSnorm` → `linear_layer`, we should map the weight tensor of the first `linear_layer` to the stationary tensor and leverage TensorE to perform cross-partition summation for `RMSnorm`.

## Optimizing Data Movement Efficiency

The key goal of optimizing memory-bound kernels is to keep the DMA engines running at high bandwidth utilization as much as possible. If you are seeing major DMA engine idle gaps in neuron-profile, you should first find ways to hide compute behind DMA activities using techniques discussed in [Opt #4](#). The rest of this section is going to focus on optimizations to improve DMA bandwidth utilization. All the optimizations below are applicable to a common symptom: computation blocked by DMA activities, which are keeping the DMA engines “busy” but at low bandwidth utilization (< 60%):



Fig. 7.61: Busy DMA engines with relatively idle compute engines.

Note, the current NKI release only supports running a kernel on a single NeuronCore (subject to changes in future releases). Therefore, the optimizations below will focus solely on movement between device memory and on-chip memory SBUF for now.

### Opt #9: Perform sufficiently large DMA transfers

**Symptoms:** A quick way to determine whether the DMA transfers are moving large enough amount of data per transfer is to visualize the DMA activities per engine in neuron-profile:

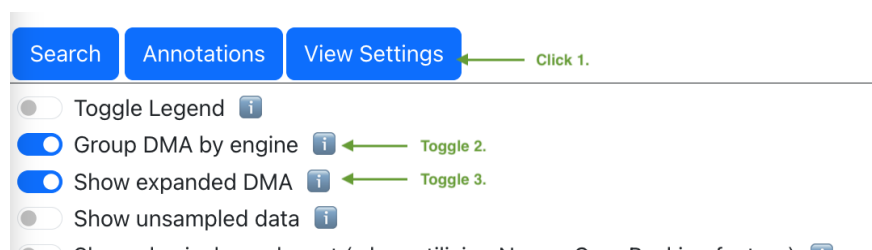


Fig. 7.62: Change view settings to visualize DMA transfer per DMA engine.

With the above view settings, each DMA transfer will be shown with a continuous bar on the execution trace, grouped by DMA engines. Below is a profile example with small DMA transfers going on all 16 DMA engines. Visually, we can see DMA engine empty gaps (due to DMA overhead) are taking up more time than active DMA transfers. Hovering over some of DMA transfers, we can also see a transfer size of 4B, which is extremely tiny. For reference, the transfer size on Trainium/Inferential2 should be larger than 32KiB to achieve ideal bandwidth.

For comparison, here’s another profile with sufficiently large DMA transfers, achieving close 70% DMA throughput utilization:



Fig. 7.63: Example timeline with tiny DMA transfers.

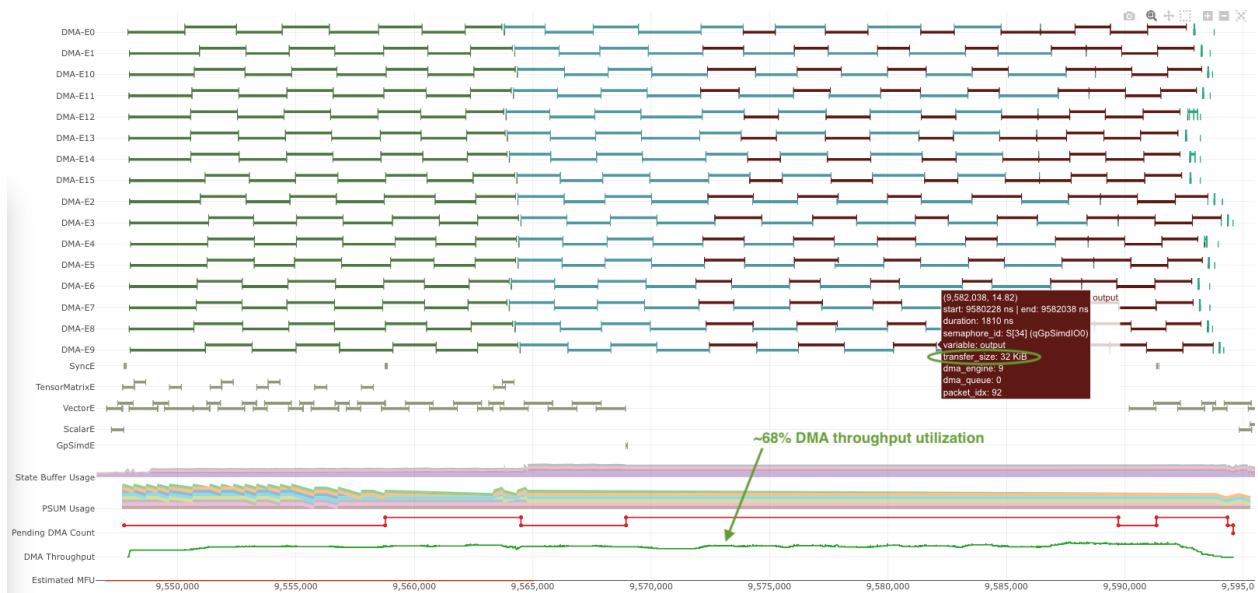


Fig. 7.64: Example timeline with large DMA transfers.

**Optimizations:** Refer to the architecture guide for more detailed discussion on DMA engines and intuitions behind the need for large DMA transfer sizes to achieve good DMA efficiency. Here, we will discuss simple rule of thumbs in NKI to trigger large DMA transfers: maximize the partition and free dimension sizes in both `nl.load` and `nl.store`. For example, the below data loading will trigger 16 DMA transfers that can be run on all 16 DMA engines, which each transfer loading 8 SBUF partitions' worth of data with a transfer size of 32KiB:

```
import neuronxcc.nki.language as nl

def load_store_32kib_contiguous(in_tensor, out_tensor):
    # both in_tensor and out_tensor have FP32 data type, 4B/element
    assert in_tensor.dtype == out_tensor.dtype == nl.float32
    # both have shape 128x1024 in device memory
    assert in_tensor.shape == out_tensor.shape == [128, 1024]

    # partition dim size is at maximum supported by the architecture: 128
    # free dim size is at the ideal size to achieve good bandwidth usage: 1024
    # Beyond 1024 has diminished return on bandwidth and
    # runs the risk of degrading compute/data movement pipelining efficiency
    i_p, i_f = nl.mgrid[0:128, 0:1024]

    # This access pattern should map to 16 DMA transfers (1 transfer/DMA engine),
    # with each DMA transfer moving 8 partitions worth of data:
    # 8 partitions * 1024 elements * 4B/element = 32 KiB
    data_tile = nl.load(in_tensor[i_p, i_f])

    # Do some useful computation
    ...

    # Store, similar size as the load
    nl.store(out_tensor[i_p, i_f], data_tile)
```

### Opt #10: Minimize use of DMA transposes.

**Symptom:** Excessive use of DMA transposes, invoked through `nl.load_transpose2d`, can degrade DMA bandwidth significantly. In `neuron-profile`, you can find out whether `nl.load_transpose2d` is taking up substantial amount of execution time by using the search functionality, which will highlight all the DMA activities that perform transposes on the fly:

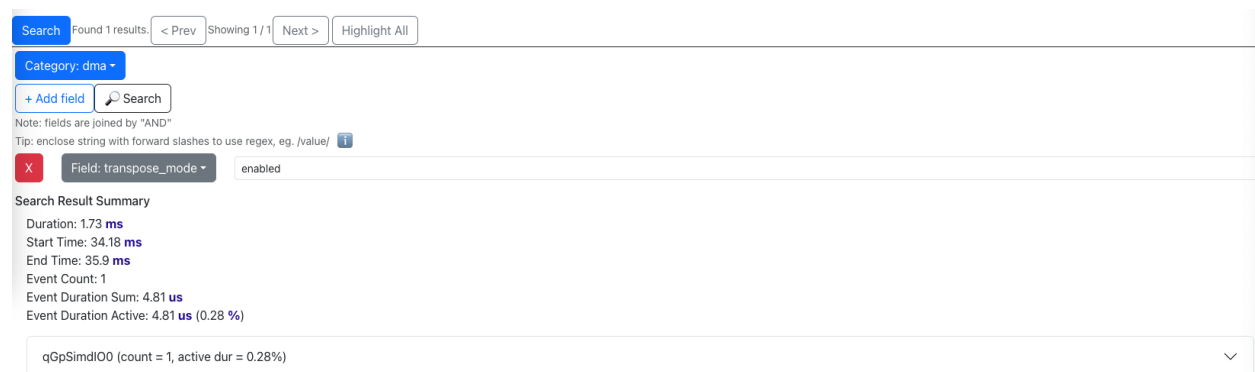


Fig. 7.65: Search for DMA activities that perform transposes.



**Optimizations:** Refer to *Opt #8* for a detailed discussion on how to eliminate the need of transposes on device memory input data. When the transposes are inevitable and the kernel is memory bound, we recommend replacing `nl.load_transpose2d` with `nl.load()` and `nisa.nc_transpose()`. For example, if you have an `in_tensor` of shape [8192, 128] in device memory but you would like an SBUF tile of shape [128, 8192] spread across 128 partitions for computation, the following two code snippets can achieve the same functionality:

```
# Option 1, low DMA bandwidth usage:
sbuf_opt1 = nl.load_transpose2d(in_tensor[0:8192, 0:128])

# Option 2, better DMA bandwidth usage, fastest transpose:
sbuf_opt2 = nl.ndarray((128, 8192), dtype=in_tensor.dtype)
for i_in_tile in nl.affine_range(8192 // 128):
    i_start = i_in_tile*128
    current_tile = nl.load(in_tensor[i_start:i_start+128, 0:128])
    sbuf_opt2[0:128, i_start:i_start+128] = nisa.nc_transpose(current_tile)
```

Option 2 above is especially great for cases where `nl.load_transpose2d` is slowing down data movement in the critical path and TensorE is otherwise idle. Occasionally Option 1 can still be the right call, when the amount of data to be transposed is small and the overhead of `nl.load_transpose2d` can be well hidden behind other useful computation.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## NKI Direct Allocation Developer Guide

In this document, we will discuss how to perform direct allocation in NeuronCore on-chip memory (SBUF and PSUM) correctly and efficiently to improve NKI kernel performance. This document is organized into five sections:

- *Background: NKI Tensors Semantics*
- *Introduction to NKI Direct Allocation API*
- *Development Best Practices*
- *Common Errors*
- *Known Limitations*

### Background: NKI Tensors Semantics

As discussed in *NKI programming model*, a multi-dimensional NKI Tensor in SBUF/PSUM must have a dimension mapped to the partition (P) dimension of the physical memory, labeled using `nl.par_dim`. We also define any NKI Tensor with the first dimension as the partition dimension is considered a NKI Tile, which is the data type that NKI compute APIs operate on. The remaining dimensions after the partition dimension in a NKI Tile are considered free (F) dimensions. The free dimensions describe how data is organized within each SBUF/PSUM partition.

To introduce NKI allocation API, let's define the block (B) dimension as any dimension before the partition dimension in a NKI Tensor. Therefore, a NKI Tensor has three types of dimensions: (B, P, F). Note, a NKI Tensor can have one or more dimensions in both B and F, but there can only be one dimension in P due to Neuron ISA requirements. The block dimension effectively describes how many (P, F) NKI Tiles the tensor has, which commonly corresponds to how many NKI API compute API invocations we need to process the entire tensor.

As an example, we can declare a NKI Tensor in SBUF as follows:



```
nki_tensor = nl.ndarray((16, nl.par_dim(128), 512), dtype=nl.bfloat16, buffer=nl.sbuf)

for i_block in nl.affine_range(16):
    nki_tensor[i_block, :, :] = nl.load(...)
    ...                          = nl.exp(nki_tensor[i_block, :, :])
```

Here, `nki_tensor` has a block dimension of 16, and we use it as an intermediate SBUF tensor for loading data from device memory and feeding inputs to the exponential operator. The different tiles in `nki_tensor` are processed by different iterations of a loop constructed with `nl.affine_range(...)`, indicating no loop-carried dependency across iterations.

In fact, the block dimension of `nki_tensor` as defined above is only considered **logical** in NKI. Neuron Compiler inspects the above code and uses a heuristic-driven allocator to decide how many physical tiles are allocated to each tensor. The key performance goal of the allocator is to achieve instruction parallelism across different engines in NeuronCore while minimizing memory usage in the on-chip memory.

Let's first consider the case where `nki_tensor` has only one physical tile, `T`, allocated in SBUF. The different loop iterations will end up completely serialized:

```
i_block = 0
1. nl.load(nki_tensor[0, :, :]) => write ``T``
2. nl.exp(nki_tensor[0, :, :]) => read ``T``

i_block = 1
3. nl.load(nki_tensor[1, :, :]) => write ``T``
4. nl.exp(nki_tensor[1, :, :]) => read ``T``

...
```

Here, we say only one logical tile in `nki_tensor` is alive at a time, because there is only one physical tile as the backing storage for `nki_tensor`. However, in NeuronCore, `nl.load` and `nl.exp` are executed using two independent resources: DMA Engine and Scalar Engine. In this serialized execution, there is no instruction parallelism achieved between these engines.

An obvious improvement is to adopt “double buffering”, by allocating two physical tiles for `nki_tensor`, `T0` and `T1`. Now, the execution enables much better computation and data movement overlapping:

```
i_block = 0
1. nl.load(nki_tensor[0, :, :]) => write ``T0``

i_block = 0 & 1
2. nl.load(nki_tensor[1, :, :]) => write ``T1`` | nl.exp(nki_tensor[0, :, :]) => read
↳ ``T0``

i_block = 1 & 2
3. nl.load(nki_tensor[1, :, :]) => write ``T0`` | nl.exp(nki_tensor[1, :, :]) => read
↳ ``T1``

...
```

Here, we reuse, or rotate, the same physical tiles across the loop iterations. No physical tile is being read from and written to in the same time step, while the DMA and Scalar Engines can operate in parallel. Besides DMA and Scalar Engines, NeuronCore also consists of Tensor, Vector, Gpsimd Engines that can execute instructions in parallel.

Given the amount of parallelism available in hardware and the complex parallel programs seen in common machine learning workloads, the heuristic-based memory allocator in Neuron Compiler may not yield the optimal allocation

decisions. Bad allocation decisions typically lead to sub-optimal engine parallelism and/or on-chip memory over-subscription causing excessive spills of intermediate data to device memory. With NKI direct allocation API, programmers can now bypass the compiler allocator and take full control of memory allocation in SBUF/PSUM for NKI Tensors.

## Direct Allocation API

This section will go over the SBUF allocation in detail, including `ncc.sbuf.alloc()` API that provides the most flexibility for tensor allocation and `ncc.sbuf.mod_alloc()` API that provides ease-of-use using a modulo allocation strategy. Both of these APIs can be used to replace the automatic allocated buffer type `buffer=nl.sbuf` when declaring a NKI tensor:

```
# Automatic allocation
nki_tensor = nl.ndarray((16, nl.par_dim(128), 512), ..., buffer=ncc.sbuf.auto_alloc())
nki_tensor = nl.ndarray((16, nl.par_dim(128), 512), ..., buffer=nl.sbuf) # alias of auto_
↪ alloc

# Direct allocation, full flexibility
nki_tensor = nl.ndarray((16, nl.par_dim(128), 512), ..., buffer=ncc.sbuf.alloc(...))

# Direct allocation, modulo allocation
nki_tensor = nl.ndarray((16, nl.par_dim(128), 512), ..., buffer=ncc.sbuf.mod_alloc(...))
```

The PSUM allocation APIs, `ncc.psum.alloc()` and `ncc.psum.mod_alloc()` follow a highly similar design. For more information on the semantics of these APIs, check out [API reference page for allocation control](#).

### `ncc.sbuf.alloc()`

This SBUF allocation API enables user to control:

- the number of physical tiles to allocate for a given NKI Tensor, and
- the exact mapping between logical tile and physical tile in SBUF

`ncc.sbuf.alloc()` accepts a single input parameter, `func`, which is a user-defined callable object that takes in:

1. a tuple of integers `idx` representing a logical block index,
2. an integer `pdim_size` for the number of partitions the logical tile has, and
3. an integer `fdim_size` for the number of bytes the logical tile has per partition.

The `func` returns a tuple of two integers, (`start_partition`, `byte_addr`), representing the memory location of the mapped physical tile for the given logical block. `start_partition` indicates the starting partition of physical tile and must follow these ISA rules:

- If  $64 < \text{pdim\_size} \leq 128$ , `start_partition` must be 0
- If  $32 < \text{pdim\_size} \leq 64$ , `start_partition` must be 0 or 64
- If  $0 < \text{pdim\_size} \leq 32$ , `start_partition` must be one of 0/32/64/96

The `byte_addr` indicates the byte offset into each partition the physical tile allocation starts from. For example, on NeuronCore-v2, a valid `byte_addr` can be any integer values from 0 (inclusive) to  $192\text{KiB} - 16\text{KiB} = (192 - 16) * 1024$  (exclusive). 192KiB is the physical size of a SBUF partition and 16KiB is allocated for compiler internal usage. Refer to [NeuronDevice Architecture Guide](#) for the physical SBUF partition size on each NeuronCore version. In addition, `byte_addr` must be aligned to `nki.language.constants.sbuf_min_align`.

At compile time, the compiler will statically evaluate `func` over indices of all the logical tiles defined in the NKI tensor to calculate physical addresses for each tile. As an example, consider the following simple allocation that allocates four physical tiles back to back along the free dimension of SBUF, with every logical tile mapped to a different physical tile sequentially.

```
def simple_1d_alloc_func(idx, pdim_size, fdim_size):
    idx, = idx # unpack the tuple
    return (0, idx * fdim_size)

t = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
               buffer=ncc.sbuf.alloc(simple_1d_alloc_func))
```

In this example, the compiler will query `simple_1d_alloc_func` with `idx` ranging from `(0, )` to `(3, )`, `pdim_size=128`, and `fdim_size=512*sizeof(nl.bfloat16)=1024`. We can visualize the final allocation in Fig. 7.66.

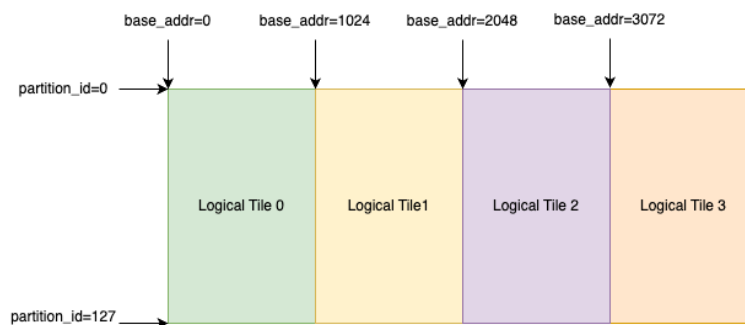


Fig. 7.66: Visualization of `simple_1d_alloc_func` in SBUF.

This `ncc.sbuf.alloc` API provides great flexibility through the customizable function to perform logical to physical tile mapping. With Python closures, the function can carry arbitrary metadata, which enables programmers to define their own memory allocator. As another example, here's a simple allocator that queries a global variable `next_addr` to keep track of the next available byte address in the free dimension.

```
next_addr = 0
def simple_1d_alloc_factory(total_fdim_size):
    base_addr = next_addr
    next_addr += total_fdim_size

    def simple_1d_alloc_func(idx, pdim_size, fdim_size):
        # unpack the tuple
        idx, = idx

        # hard-code to partition 0, since each tile takes up 128 partitions
        start_partition = 0

        return (start_partition, base_addr + idx * fdim_size)

    return simple_1d_alloc_func

# Using simple_1d_alloc_factory, next_addr is automatically incremented.
# Physical tiles of t0 and t1 start at 0 and 4096, respectively
t0 = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
```

(continues on next page)

(continued from previous page)

```

        buffer=ncc.sbuf.alloc(simple_1d_alloc_factory(512*2*4)))
t1 = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
        buffer=ncc.sbuf.alloc(simple_1d_alloc_factory(512*2*4)))

```

### ncc.sbuf.mod\_alloc()

Alternative to the `ncc.sbuf.alloc()` API which requires programmers to define an allocation algorithm from scratch, NKI also provides the `ncc.sbuf.mod_alloc()` API which invokes a pre-defined modulo allocation scheme in Neuron Compiler.

Modulo allocation works as follows. Suppose that we allocate **two** physical tiles for a tensor with a logical shape of (8, `par_dim(128)`, 512). The eight logical tiles are assigned to the two physical tiles by taking a modulo of two on the logical tile index (that is, block index). Therefore, logical tiles with index (0, ), (2, ), (4, ), (6, ) share the same physical tile, while logical tiles (1, ), (3, ), (5, ), (7, ) share the other physical tile.

The `ncc.sbuf.mod_alloc` API takes four input parameters:

1. `base_addr` indicates the starting byte offset within each SBUF partition of the physical tiles.
2. `base_partition` indicates the starting SBUF partition of the physical tiles.
3. `num_par_tiles` indicates the number of physical tiles to be allocated along the partition dimension of SBUF. This is only applicable for tiles that use fewer than 64 partitions per ISA constraints.
4. `num_free_tiles` indicates the number of physical tiles to be allocated along the free dimension of SBUF.

Given the above input parameters and the modulo allocation scheme, Neuron Compiler is then able to calculate the physical tile memory location, (`start_partition`, `byte_addr`) for each logical tile in the tensor. Note, this is the same information that the callable allocation function passed into `ncc.sbuf.mod_alloc()` would return. See [API reference manual for ncc.sbuf.mod\\_alloc](#) for the exact formula to calculate (`start_partition`, `base_addr`).

Next, we discuss a common use case of `ncc.sbuf.mod_alloc`, which specifies only the `base_addr` and `num_free_tiles` fields while leaving the remaining parameters to default (`base_partition=0` and `num_par_tiles=(1,)`).

```

nki_tensor = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
        buffer=ncc.sbuf.mod_alloc(base_addr=0, num_free_tiles=(2, )))

```

This produces the following allocation:

Table 7.8: Modulo Allocation Example

| Logical Tile Index | Physical Tile start_partition | Physical Tile byte_addr   |
|--------------------|-------------------------------|---|
| (0, )              | 0                             | $0 + (0 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 0$    |
| (1, )              | 0                             | $0 + (1 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 1024$ |
| (2, )              | 0                             | $0 + (2 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 0$    |
| (3, )              | 0                             | $0 + (3 \% 2) * 512 * \text{sizeof}(\text{nl.bfloat16}) = 1024$ |

The above example is an easy way to implement double buffering without having to define a callable function manually like how we did for `ncc.sbuf.alloc()`. We can also implement multi-buffering using `ncc.sbuf.mod_alloc()` by changing the value of `num_free_tiles` (or `num_par_tiles` when each tile occupies less than 64 partitions).

## Development Best Practices

First and foremost, direct allocation APIs are considered advanced NKI features for performance optimizations. We highly recommend using direct allocation API only after your kernel is functionally correct with automatic allocation. Automatic allocation is invoked when NKI tensors are declared with `buffer=nl.sbuf` (alias of `ncc.sbuf.auto_alloc`) or `buffer=nl.psum` (alias of `buffer=ncc.psum.auto_alloc`).

The rest of this section goes over best practices of direct allocation APIs to optimize kernel performance.

### #1. Hoist allocations outside of the loop-nests you want to block across

To parallelize a loop, every tensor used in the loop must have multiple live tiles so that the different hardware engines can read/write from/to different memory locations in parallel. To achieve this, make sure to allocate tensors with logical block dimensions above any loop you want to run in parallel. For example, the following loop will be serialized because `t` has only one tile alive,

```
for i in affine_range(8):
    t = nl.ndarray((128, 512), dtype=..., buffer=ncc.sbuf.mod_alloc(base_addr=0))
    t[i] = ...
    # do something with t
```

To improve parallelism, programmers should hoist the tensor declaration and allocation above the loop, like this:

```
t = nl.ndarray((8, 128, 512), dtype=...,
              buffer=ncc.sbuf.mod_alloc(base_addr=0, num_free_tiles=(8,)))
for i in affine_range(8):
    t[i] = ...
    # do something with t
```

### #2. Avoid PSUM Bank Collisions

As discussed in *NeuronDevice Architecture Guide*, PSUM in each NeuronCore has eight banks that can accumulate TensorE matrix multiplication results independently. Especially in complex loops where PSUM tensors have multiple logical block dimensions, programmers should pay close attention to PSUM bank allocations so that they do not collide.

*More examples coming soon.*

## Common Errors

This section goes over the common compilation error programmers may encounter while using direction allocation APIs.

### #1. Mixing direct allocation with automatic allocation.

Automatically allocated tensors from default arguments or lowering need to be explicitly passed to those NKI APIs or their allocations will collide. When direct allocation is used, all tensors, including the tensor returned from a instruction, in that kernel must also use direct allocation. For example,

```
t = nl.load(input) # t is a new tensor, this will fail

# the correct way
t = nl.ndarray(shape=..., buffer=ncc.sbuf.alloc(...))
t[...] = nl.load(input)
```

### #2. Calling compute APIs that introduce implicit tensors while direction allocation is used.

Certain NKI compute APIs implicitly create constant or intermediate tensors that are not in programmers' control. For example, invoking `nisa.nc_transpose` with `engine=nisa.tensor_engine` creates an identity matrix under the hood, which cannot be allocated explicitly using direct allocation APIs. Similarly, many high-level `nki.language` APIs, such as `nl.softmax <api/generated/nki.language.softmax>`, are lowered down to several `nisa` APIs in the compiler. The intermediate tensors between these lowered `nisa` APIs also cannot be explicitly allocated by NKI programmers.

Therefore, due to restrictions discussed in common error #1, such APIs are not allowed when direction allocation is used in the kernel. A compiler error would occur when this is violated.

### #3. Lifetime Conflicts

Each NKI kernel has its own address space, and any physical tiles that must be alive simultaneously due to compute definitions of the kernel should be assigned unique addresses in the kernel address space. For example, tensors below have partially overlapping physical memory addresses, which would cause errors if the two tensors need to be alive at the same time.

```
# t0 physical tiles occupy:
# partition [0:128],
# byte_addr [0:512*num_free_tiles*sizeof(nl.bfloat16)] = [0:2048]
t0 = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
                buffer=ncc.sbuf.mod_alloc(base_addr=0, num_free_tiles=(2, )))

# t1 physical tiles occupy:
# partition dim - [0:128],
# free dim (byte_addr) - [1024:1024+512*num_free_tiles*sizeof(nl.bfloat16)] = [1024:3072]
t1 = nl.ndarray((4, par_dim(128), 512), dtype=nl.bfloat16,
                buffer=ncc.sbuf.mod_alloc(base_addr=1024, num_free_tiles=(2, )))
```

Another common lifetime conflict error is when the number of physical tiles is insufficient to hold all the logical tiles that need to be alive at the same time. For example,

```
# Lifetime conflict #
t1 = nl.ndarray((8, par_dim(128), 512),
                buffer=ncc.sbuf.mod_alloc(byte_addr=0, num_free_tiles=(2, )))

for i in nl.affine_range(8):
    t1[i] = nl.load(...)
```

(continues on next page)

(continued from previous page)

```

# End of loop: we need all eight logical tiles in t1 to be
# alive in SBUF so that we can start the next loop.
# Only two tiles can be alive according to our allocation above -> ERROR

for i in nl.affine_range(8):
    result[i] = nl.exp(t1[i])

#####

# Correct way #
for i in nl.affine_range(4):
    t1 = nl.ndarray((2, par_dim(128), 512),
                    buffer=ncc.sbuf.mod_alloc(byte_addr=0, num_free_tiles=(2, )))
    for i in nl.affine_range(2):
        t1[i] = nl.load(...).
        result[i] = nl.exp(t1[i]) # t[i] are dead after iteration, thus no error

```

Neuron Compiler has built-in checks for such lifetime conflicts. For example, when there is a PSUM tensor lifetime conflict, an error like “[SCH713] Violation of accumulation group interleaving” will be thrown. However, **the lifetime checks in the current release are in-complete**, which may not catch all the lifetime violations in the kernel.

If the kernel using direction allocation API generates incorrect results numerically, one plausible cause is the kernel has tensor lifetime conflicts that are not caught during compilation. One way to verify this is to re-compile the same kernel with automatic allocation forced on using the `force_auto_alloc` decorator.

## Known Limitations

- When direct allocation API is used, HBM tensors cannot be declared unless they are used as kernel outputs.
  - All tensors declared with `buffer=nl.shared_hbm` must be returned as the result of the kernel.
  - Tensors declared with `buffer=nl.hbm` or `buffer=nl.private_hbm` are not allowed.
  - A compilation error `Non IO HBM tensor is not supported in allocated NKI kernel: <list of tensor names>` will be thrown when such a tensor is encountered.
- For `ncc.psum.mod_alloc`, the `base_addr` and `start_partition` input fields must be 0. This implies that only one physical tile can live in a PSUM bank at a time and the PSUM tile must start from partition 0.
- A PSUM tile cannot cross bank boundaries. Therefore, the size of the free dimension of each tile has a maximum of 2KiB, or 512 FP32 elements.
- The compiler’s ability to check for race condition and lifetime conflicts is limited. It is not guaranteed to catch all race conditions.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI Block Dimension Migration Guide

The SBUF/PSUM tensors in NKI used to allow block dimensions in front of the partition dimension. The block dimension support has been removed due the following reasons.

- Removing block dimensions does not hurt the expressivity of NKI.
- Block dimension is a pure software concept and does not have direct hardware mapping.
- The block dimension is unintuitive and causes confusion.
- Using block dimension has no inherit performance benefit, particularly using block dimension has no relationship with memory throughput whatsoever.
- Multi-buffering is implicit with block dimension. Removing block dimension will make multi-buffering more natural.

This document will first explain the semantics of block dimensions in detail, then it will provide information on how to migrate existing code that uses block dimensions while maintain the functional correctness and performance.

### What are block dimensions?

Consider the following NKI tensor.

```

1 a = nl.ndarray((4, 8, nl.par_dim(128), 2, 512), buffer=nl.sbuf)
2
3 # - (4, 8): (B) block dimensions
4 # - 128: (P) partition dimension
5 # - (2, 512): (F) free dimension

```

As explained in the *Direct Allocation Guide*, a NKI tensor has three types of dimensions:  $(B, P, F)$ . The partition dimension maps to the partition dimension of the physical memory, and the free dimensions describe how data is organized in each SBUF/PSUM partition. The block dimensions described how many physical  $(P, F)$  tiles the tensor has.

The block dimension of tensors is a **logical** dimension and is a pure software concept. The compiler analyzes the memory dependency and allocates physical address to each tiles. **This means that the physical tiles may not be alive in the memory simultaneously**, and in most of the cases they don not. Consider the following code snippet that access the tensor *a*.

```

1 @nki.jit
2 def exp_func(inp):
3     output = nl.ndarray((4, 8, 128, 2, 512), dtype=float32,
4                         buffer=nl.shared_hbm)
5     a = nl.ndarray((4, 8, nl.par_dim(128), 2, 512), dtype=float32, buffer=nl.sbuf)
6     for i in range(4):
7         for j in range(8):
8             a[i, j] = nl.load(inp[i, j])
9             a[i, j] = nl.exp(a[i, j])
10            nl.store(output[i, j], value=result)

```

At the very minimum, only 1 physical tile of *a* needs to be alive. Then the execution is completely serialized. Essentially, all physical tiles would have the exact same memory address.

```

1 Physical Address Map
2

```

(continues on next page)



(continued from previous page)

```

3 output[0, 0] --> Partition 0 - 128, Free 0 - 2048B
4 output[0, 1] --> Partition 0 - 128, Free 0 - 2048B
5 ...

```

Instead, compiler could choose to allocate 2 physical tiles to  $a$ , then the dma copy from HBM to SBUF can overlap with the exponential operation. In other word, **the block dimension allows compiler to perform space-time tradeoff at liberty.**

```

1 Physical Address Map
2
3 output[0, 0] --> Partition 0 - 128, Free 0 - 2048B
4 output[0, 1] --> Partition 0 - 128, Free 2048 - 4096B
5 output[0, 2] --> Partition 0 - 128, Free 0 - 2048B
6 output[0, 3] --> Partition 0 - 128, Free 2048 - 4096B
7 ...

```

When performing the migration, it is important to understand the dependency relationship between blocks and choose the correct migration method accordingly.

## Migration for SBUF tensors

**If blocks need to be alive at the same time, move the block dimension into free dimension**

```

1 a = nl.ndarray((8, par_dim(128), 512), buffer=nl.sbuf, dtype=bfloat16)
2
3 # ----> Migrate to
4 a = nl.ndarray((128, 8, 512), buffer=nl.sbuf, dtype=bfloat16)

```

As an example, all 8 blocks of `add_buf` needs to be alive at the same time when the first for loop finishes. Therefore, the block dimension need to be fold into the free dimension.

```

1 @nki.jit
2 def sb_blocks(inp):
3     res = nl.ndarray(shape=(8, 128, 512), dtype=inp.dtype, buffer=nl.shared_hbm)
4     add_buf = nl.ndarray(shape=(8, nl.par_dim(128), 512), dtype=inp.dtype, buffer=nl.
5     ↪sbuf)
6     for i in range(8):
7         add_buf[i] = nl.load(inp[i])
8     for i in range(8):
9         nl.store(res[i], add_buf[i])
10    return res
11
12 # should migrate to
13 @nki.jit
14 def sb_blocks_migrated(inp):
15     res = nl.ndarray(shape=(8, 128, 512), dtype=inp.dtype, buffer=nl.shared_hbm)
16     add_buf = nl.ndarray(shape=(128, 8, 512), dtype=inp.dtype, buffer=nl.sbuf)
17     for i in range(8):
18         add_buf[0:128, i, 0:512] = nl.load(inp[i])
19     for i in range(8):

```

(continues on next page)

(continued from previous page)

```

19     nl.store(res[i], add_buf[0:128, i, 0:512])
20     return res

```

If blocks does not need to be alive at the same time, remove the block dimension and hoist it down

```

1  a = nl.ndarray((8, par_dim(128), 256))
2  for i in nl.affine_range(8):
3      <do something with a[i]>
4
5  # should be transformed to ....
6  for i in nl.affine_range(8):
7      a = nl.ndarray((128, 256))
8      <do something with a>

```

As an example, all 8 blocks of add\_buf does not need to be alive at the same time. We can remove the block dimension and hoist down the tensor inside the loop.

```

1  @nki.jit
2  def sb_blocks(inp):
3      res = nl.ndarray(shape=(8, 128, 512), dtype=inp.dtype, buffer=nl.shared_hbm)
4      add_buf = nl.ndarray(shape=(8, nl.par_dim(128), 512), dtype=inp.dtype, buffer=nl.
    ↪ sbuf)
5      for i in range(8):
6          add_buf[i] = nl.load(inp[i])
7          nl.store(res[i], add_buf[i])
8      return res
9
10 # should migrate to
11 @nki.jit
12 def sb_blocks_migrated(inp):
13     res = nl.ndarray(shape=(8, 128, 512), dtype=inp.dtype, buffer=nl.shared_hbm)
14     for i in range(8):
15         add_buf = nl.ndarray(shape=(128, 512), dtype=inp.dtype, buffer=nl.sbuf)
16         add_buf[0:128, 0:512] = nl.load(inp[i])
17         nl.store(res[i], add_buf[0:128, 0:512])
18     return res

```

**Warning:** To preserve performance, it is important to hoist down the tensor inside the loop.

It is important to note that the dependency relationship between loop iterations is different in `sb_blocks_migrated` and the following `sb_blocks_migrated_incorrect`.

```

1  @nki.jit
2  def sb_blocks_migrated_incorrect(inp):
3      res = nl.ndarray(shape=(8, 128, 512), dtype=inp.dtype, buffer=nl.shared_hbm)
4      add_buf = nl.ndarray(shape=(128, 512), dtype=inp.dtype, buffer=nl.sbuf)
5      for i in range(8):
6          add_buf[0:128, 0:512] = nl.load(inp[i])

```

(continues on next page)

(continued from previous page)

```

7     nl.store(res[i], add_buf[0:128, 0:512])
8     return res

```

In `sb_blocks_migrated`, compiler could unroll the loop and materialize multiple copies of the tensor `add_buf`. However, in the `sb_blocks_migrated_incorrect`, the execution will be serialized because the loop carries dependency on `add_buf`.

## Migration for PSUM tensors

---

**Note:** To be filled, the backend support for removing blocks in PSUM tensor is still in progress.

---

## Migration of direct allocation & multi-buffering

---

**Note:** For more information on direct allocation API, please refer to *Direct Allocation Guide*

---

When we have block dimensions, we allocate interleaved address for blocks to achieve multi-buffering.

```

1  def interleave_alloc_func(idx, pdim_size, fdim_size):
2      """
3      This function assumes 1d block dimension, and will allocate unique
4      address by modulo of 2.
5
6      For a tensor of 4 blocks, block 0 and 2 will have the same address, while
7      block 1 and 3 will have the same address that is different to that of 0 and 2.
8      """
9      # unpack the tuple
10     idx, = idx
11
12     # hard-code to partition 0, since each tile takes up 128 partitions
13     start_partition = 0
14
15     return (start_partition, (idx % 2) * fdim_size)
16
17 @nki.jit
18 def copy_func(inp):
19     output = nl.ndarray((4, 128, 512), dtype=float32, buffer=nl.shared_hbm)
20     a = nl.ndarray((4, nl.par_dim(128), 512), dtype=float32, buffer=ncc.sbuf.
21     ↪ alloc(interleave_alloc_func))
22     for i in range(4):
23         a[i] = nl.load(inp[i])
24         nl.store(output[i], value=a[i])

```

After removing the block dimension, we could write the following to implement the same multi-buffering, which is actually more natural and closer to that on CPU.

```

1  def interleave_alloc_func(idx, pdim_size, fdim_size):
2      """

```

(continues on next page)

(continued from previous page)

```

3  This function assumes 1d block dimension, and will allocate unique
4  address by modulo of 2.
5
6  For a tensor of 4 blocks, block 0 and 2 will have the same address, while
7  block 1 and 3 will have the same address that is different to that of 0 and 2.
8  """
9  # unpack the tuple
10 assert idx == () # We don't have any block dimension
11
12 # hard-code to partition 0, since each tile takes up 128 partitions
13 start_partition = 0
14
15 return (start_partition, (idx % 2) * fdim_size)
16
17 @nki.compiler.skip_middle_end_transformations
18 @nki.jit
19 def exp_func(inp):
20     output = nl.ndarray((4, 128, 512), dtype=nl.float32, buffer=nl.shared_hbm)
21     a = nl.ndarray((128, 2, 512), dtype=nl.float32, buffer=ncc.sbuf.alloc(interleave_alloc_
22     ↪ func))
23     for i in range(4):
24         a[0:128, i % 2, 0:512] = nl.load(inp[i])
25         nl.store(output[i], value=a[0:128, i % 2, 0:512])

```

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## NKI Tutorials

The full source code of the following tutorials can be also viewed on the [nki-samples](#) repository on GitHub.

*This document is relevant for:* Inf2, Trn1, Trn2

### Single program, multiple data tensor addition

In this tutorial we write a simple tensor addition kernel using NKI in PyTorch and JAX. In doing so, we learn about:

- The NKI syntax and the *SPMD programming model*.
- Best practices for validating and benchmarking your custom kernel against a reference native PyTorch or JAX implementation.

---

**Note:** This tutorial is written using the SPMD programming model in NKI. However, as discussed in *NKI programming guide*, adopting the SPMD programming model has **no** impact on performance of NKI kernel, and therefore is considered **optional** in current NKI release.

---

## PyTorch

### Compute kernel

We start by defining the compute kernel that has large tensor inputs, but operates on a subset of the tensor at a tile size of [128, 512]. The partition dimension tile size is chosen according to the tile size restrictions (*nki.language.tile\_size.pmax*), while the free dimension tile size is chosen arbitrarily (512).

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3
4
5 @nki.jit
6 def nki_tensor_add_kernel_(a_input, b_input):
7     """NKI kernel to compute element-wise addition of two input tensors
8
9     This kernel assumes strict input/output sizes can be uniformly tiled to [128,512]
10
11     Args:
12         a_input: a first input tensor
13         b_input: a second input tensor
14
15     Returns:
16         c_output: an output tensor
17     """
18     # Create output tensor shared between all SPMD instances as result tensor
19     c_output = nl.ndarray(a_input.shape, dtype=a_input.dtype, buffer=nl.shared_hbm)
20
21     # Calculate tile offsets based on current 'program'
22     offset_i_x = nl.program_id(0) * 128
23     offset_i_y = nl.program_id(1) * 512
24
25     # Generate tensor indices to index tensors a and b
26     ix = offset_i_x + nl.arange(128)[: , None]
27     iy = offset_i_y + nl.arange(512)[None, :]
28
29     # Load input data from device memory (HBM) to on-chip memory (SBUF)
30     # We refer to an indexed portion of a tensor as an intermediate tensor
31     a_tile = nl.load(a_input[ix, iy])
32     b_tile = nl.load(b_input[ix, iy])
33
34     # compute a + b
35     c_tile = a_tile + b_tile
36
37     # store the addition results back to device memory (c_output)
38     nl.store(c_output[ix, iy], value=c_tile)
39
40     # Transfer the ownership of `c_output` to the caller
41     return c_output

```

In this example:

1. We define the NKI kernel in `nki_tensor_add_kernel_`, decorate it with the *nki.jit* decorator to call the nki compiler to compile the kernel.

2. Inside, we first allocate tensor `c_output` as the result of the kernel
3. Next, we define offsets into the tensors, based on the ID of the worker executing the code (`nl.program_id`), and generate tile indices using these offsets with `nl.arange`. We use advanced indexing here to showcase how it works. Basic indexing with slicing can also work. See *NKI Programming Model* for more information on different tensor indexing modes.
4. We use `nl.program_id` to enable SPMD execution (single-program, multiple-data, see *SPMD: Launching Multiple Instances of a Kernel*), where each worker only operates on a (sub-tensor) tile of the input/output tensors. By accessing its own `program_id`, each worker can calculate the offsets it needs to access the correct tiles.
5. The first axis of the tensor (mapped to the partition-dimension) is tiled into blocks of 128, based on hardware restrictions (see *Tile Size Considerations*). The second axis (mapped to the free-dimension) is tiled into blocks of 512 (no tile-size constraint, since the addition operation is performed on the Vector engine, the only restriction is on-chip memory capacity).
6. We then load sub-tensors data from tensors `a_input` and `b_input` using `nl.load`, to place the tiles `a_tile` and `b_tile` in the on-chip memory (SBUF)
7. We sum them to compute `c_tile`, and store it back to DRAM in the relevant portion of the `c_output` tensor, using `nl.store`. Since both inputs and output are the same shape, we can use the same set of indices to access all three tensors.
8. At the end, we use `return` statement to transfer the ownership of tensor `c_output` to the caller of the kernel.

## SPMD execution

We declare a helper function, to launch the compute-kernel with appropriate grid/block sizes, to perform the computation over the whole input tensors.

```

1 def nki_tensor_add(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the
5     ↪ inputs/outputs
6
7     Args:
8         a_input: a first input tensor, of shape [N*128, M*512]
9         b_input: a second input tensor, of shape [N*128, M*512]
10
11     Returns:
12         a tensor of shape [N*128, M*512], the result of a_input + b_input
13     """
14     # The SPMD launch grid denotes the number of kernel instances.
15     # In this case, we use a 2D grid where the size of each invocation is 128x512
16     grid_x = a_input.shape[0] // 128
17     grid_y = a_input.shape[1] // 512
18
19     return nki_tensor_add_kernel_[grid_x, grid_y](a_input, b_input)

```

We are using a two-dimensional grid, where the first dimension of the tensor is tiled in the X dimension of the grid, while the second dimension is tiled in the Y dimension of the grid. In this scenario we assume that tensor sizes are a multiple of maximum tile sizes allowed, so we do not need to handle partial tiles.

## Launching kernel and testing correctness

To execute the kernel, we prepare tensors a and b, and call the `nki_tensor_add` helper function. We also verify the correctness of the NKI kernel against, torch by comparing the outputs of both, using `torch.allclose`:

```

1 import torch
2 from torch_xla.core import xla_model as xm
3
4 if __name__ == "__main__":
5     device = xm.xla_device()
6
7     a = torch.rand((256, 1024), dtype=torch.bfloat16).to(device=device)
8     b = torch.rand((256, 1024), dtype=torch.bfloat16).to(device=device)
9
10    output_nki = nki_tensor_add(a, b)
11    print(f"output_nki={output_nki}")
12
13    output_torch = a + b
14    print(f"output_torch={output_torch}")
15
16    allclose = torch.allclose(output_torch, output_nki, atol=1e-4, rtol=1e-2)
17    if allclose:
18        print("NKI and Torch match")
19    else:
20        print("NKI and Torch differ")
21
22    assert allclose

```

Output:

```

2023-12-29 15:18:00.000558: 14283 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2023-12-29 15:18:00.000559: 14283 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↳'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↳workdir/49f554a2-2c55-4a88-8054-cc9f20824a46/model.MODULE_5007921933048625946+d41d8cd9.
↳hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/49f554a2-2c55-4a88-8054-
↳cc9f20824a46/model.MODULE_5007921933048625946+d41d8cd9.neff', '--verbose=35']
.
Compiler status PASS
output_nki=tensor([[0.9297, 0.8359, 1.1719, ..., 0.4648, 0.2188, 0.9336],
[0.3906, 1.3125, 0.8789, ..., 1.6562, 1.7734, 0.9531],
[0.6445, 1.1406, 1.3281, ..., 0.9531, 0.8711, 0.9336],
...,
[0.4023, 0.6406, 1.5312, ..., 0.7617, 0.7734, 0.3359],
[0.8125, 0.7422, 1.2109, ..., 0.8516, 1.2031, 0.5430],
[1.3281, 1.2812, 1.3984, ..., 1.2344, 0.8711, 0.5664]],
device='xla:1', dtype=torch.bfloat16)
2023-12-29 15:18:02.000219: 14463 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2023-12-29 15:18:02.000220: 14463 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↳'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↳workdir/2e135b73-1c3b-45e4-a6f0-2c4b105c20e5/model.MODULE_
↳10032327759287407517+d41d8cd9.hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/

```

(continues on next page)

(continued from previous page)

```

↪2e135b73-1c3b-45e4-a6f0-2c4b105c20e5/model.MODULE_10032327759287407517+d41d8cd9.neff',
↪'--verbose=35']
*
Compiler status PASS
output_torch=tensor([[0.9297, 0.8359, 1.1719, ..., 0.4648, 0.2188, 0.9336],
                    [0.3906, 1.3125, 0.8789, ..., 1.6562, 1.7734, 0.9531],
                    [0.6445, 1.1406, 1.3281, ..., 0.9531, 0.8711, 0.9336],
                    ...,
                    [0.4023, 0.6406, 1.5312, ..., 0.7617, 0.7734, 0.3359],
                    [0.8125, 0.7422, 1.2109, ..., 0.8516, 1.2031, 0.5430],
                    [1.3281, 1.2812, 1.3984, ..., 1.2344, 0.8711, 0.5664]],
                    device='xla:1', dtype=torch.bfloat16)
2023-12-29 15:18:03.000797: 14647 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↪neuron-compile-cache
2023-12-29 15:18:03.000798: 14647 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↪'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↪workdir/74f8b6ae-76d9-4dd8-af7f-e5e1c40a27a3/model.MODULE_5906037506311912405+d41d8cd9.
↪hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/74f8b6ae-76d9-4dd8-af7f-
↪e5e1c40a27a3/model.MODULE_5906037506311912405+d41d8cd9.neff', '--verbose=35']
*
Compiler status PASS
NKI and Torch match

```

Note that the tensor values you see will differ from what's printed above, since this example uses `torch.rand` to initialize the inputs.

## JAX

### Compute kernel

We can reuse the same NKI compute kernel defined for PyTorch above.

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3
4
5 @nki.jit
6 def nki_tensor_add_kernel(a_input, b_input):
7     """NKI kernel to compute element-wise addition of two input tensors
8
9     This kernel assumes strict input/output sizes can be uniformly tiled to [128,512]
10
11     Args:
12         a_input: a first input tensor
13         b_input: a second input tensor
14
15     Returns:
16         c_output: an output tensor
17     """
18     # Create output tensor shared between all SPMD instances as result tensor
19     c_output = nl.ndarray(a_input.shape, dtype=a_input.dtype, buffer=nl.shared_hbm)

```

(continues on next page)



(continued from previous page)

```

20
21 # Calculate tile offsets based on current 'program'
22 offset_i_x = nl.program_id(0) * 128
23 offset_i_y = nl.program_id(1) * 512
24
25 # Generate tensor indices to index tensors a and b
26 ix = offset_i_x + nl.arange(128)[: , None]
27 iy = offset_i_y + nl.arange(512)[None, :]
28
29 # Load input data from device memory (HBM) to on-chip memory (SBUF)
30 # We refer to an indexed portion of a tensor as an intermediate tensor
31 a_tile = nl.load(a_input[ix, iy])
32 b_tile = nl.load(b_input[ix, iy])
33
34 # compute a + b
35 c_tile = a_tile + b_tile
36
37 # store the addition results back to device memory (c_output)
38 nl.store(c_output[ix, iy], value=c_tile)
39
40 # Transfer the ownership of `c_output` to the caller
41 return c_output

```

## SPMD execution

Now we can also declare a helper function, to launch the compute-kernel with appropriate grid/block sizes, to perform the computation:

```

1 def nki_tensor_add(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the
5     ↪ inputs/outputs
6
7     Args:
8         a_input: a first input tensor, of shape [N*128, M*512]
9         b_input: a second input tensor, of shape [N*128, M*512]
10
11     Returns:
12         a tensor of shape [N*128, M*512], the result of a_input + b_input
13     """
14
15     # The SPMD launch grid denotes the number of kernel instances.
16     # In this case, we use a 2D grid where the size of each invocation is 128x512
17     grid_x = a_input.shape[0] // 128
18     grid_y = a_input.shape[1] // 512
19
20     return nki_tensor_add_kernel_[grid_x, grid_y](a_input, b_input)

```

We are using a two-dimensional grid, where the first dimension of the tensor is tiled in the X dimension of the grid, while the second dimension is tiled in the Y dimension of the grid. In this scenario we assume that tensor sizes are a

multiple of maximum tile sizes allowed, so we do not need to handle partial tiles.

### Launching kernel and testing correctness

To execute the kernel, we prepare arrays `a` and `b`, and call the `nki_tensor_add` helper function. We also verify the correctness of the NKI kernel against, JAX by comparing the outputs of both, using `jax.numpy.allclose`:

```

1 import jax
2 import jax.numpy as jnp
3
4 if __name__ == "__main__":
5
6     seed_a, seed_b = jax.random.split(jax.random.PRNGKey(42))
7     a = jax.random.uniform(seed_a, (256, 1024), dtype=jnp.bfloat16)
8     b = jax.random.uniform(seed_b, (256, 1024), dtype=jnp.bfloat16)
9
10    output_nki = nki_tensor_add(a, b)
11    print(f"output_nki={output_nki}")
12
13    output_jax = a + b
14    print(f"output_jax={output_jax}")
15
16    allclose = jnp.allclose(output_jax, output_nki, atol=1e-4, rtol=1e-2)
17    if allclose:
18        print("NKI and JAX match")
19    else:
20        print("NKI and JAX differ")
21
22    assert allclose

```

Output:

```

.
Compiler status PASS
.
Compiler status PASS
.
Compiler status PASS
output_nki=[[0.992188 1.27344 1.65625 ... 0.90625 1.34375 1.77344]
 [0 0.90625 1.34375 ... 0.390625 0.703125 0.914062]
 [0.5 0.390625 0.703125 ... 1.22656 1.15625 1.01562]
 ...
 [1.98438 1.98438 1.98438 ... 1.33594 1.64062 1.35938]
 [0.992188 1.33594 1.64062 ... 1.16406 1.67188 1.20312]
 [1.49219 1.16406 1.67188 ... 1.375 1 1.6875]]
.
Compiler status PASS
output_jax=[[0.992188 1.27344 1.65625 ... 0.90625 1.34375 1.77344]
 [0 0.90625 1.34375 ... 0.390625 0.703125 0.914062]
 [0.5 0.390625 0.703125 ... 1.22656 1.15625 1.01562]
 ...
 [1.98438 1.98438 1.98438 ... 1.33594 1.64062 1.35938]
 [0.992188 1.33594 1.64062 ... 1.16406 1.67188 1.20312]

```

(continues on next page)

(continued from previous page)

```
[1.49219 1.16406 1.67188 ... 1.375 1 1.6875]]
.
Compiler status PASS
NKI and JAX match
```

Note that the array values you see will differ from what's printed above, since this example uses `jax.random.uniform` to initialize the inputs.

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- NKI baremetal implementation: `spmd_tensor_addition_nki_kernels.py`
- **PyTorch implementation: `spmd_tensor_addition_torch.py`**
  - You must also download `spmd_tensor_addition_nki_kernels.py` into the same folder to run this PyTorch script.
- **JAX implementation: `spmd_tensor_addition_jax.py`**
  - You must also download `spmd_tensor_addition_nki_kernels.py` into the same folder to run this PyTorch script.

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

Run NKI baremetal implementation:

```
python3 spmd_tensor_addition_nki_kernels.py
```

Run PyTorch implementation:

```
python3 spmd_tensor_addition_torch.py
```

Run JAX implementation:

```
python3 spmd_tensor_addition_jax.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Single program, multiple data tensor addition using multiple Neuron Cores

In this tutorial we reuse the *simple tensor addition kernel*, but directly control how our kernels and tensors are distributed across multiple neuron cores.

Doing so, we expand our knowledge about:

- The NKI syntax and the *SPMD programming model*.
- `nki.language.spmd_dim()` and `nki.language.nc()`

## PyTorch

## Reusing existing compute kernel in helper function

We start by reusing the `nki_tensor_add_kernel_` compute kernel that has large tensor inputs, but operates on a subset of the tensor at a tile size of `[128, 512]`. The partition dimension tile size is chosen according to the tile size restrictions (`nki.language.tile_size.pmax`), while the free dimension tile size is chosen arbitrarily (512).

```

1 def nki_tensor_add_nc2(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors using
   ↪ multiple Neuron cores.
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the
   ↪ inputs/outputs.
5     a_input and b_input are sharded across Neuron cores, directly utilizing Trn2
   ↪ architecture capabilities
6
7     Args:
8         a_input: a first input tensor, of shape [N*128, M*512]
9         b_input: a second input tensor, of shape [N*128, M*512]
10
11     Returns:
12         a tensor of shape [N*128, M*512], the result of a_input + b_input
13     """
14
15     # The SPMD launch grid denotes the number of kernel instances.
16     # In this case, we use a 2D grid where the size of each invocation is 128x512
17     # Since we're sharding across neuron cores on the 1st dimension we want to do our
   ↪ slicing at
18     # 128 per core * 2 cores = 256
19     grid_x = a_input.shape[0] // (128 * 2)
20     grid_y = a_input.shape[1] // 512
21
22     # In addition, we distribute the kernel to physical neuron cores around the first
   ↪ dimension
23     # of the spmd grid.
24     # This means:
25     # Physical NC [0]: kernel[n, m] where n is even
26     # Physical NC [1]: kernel[n, m] where n is odd
27     # notice, by specifying this information in the SPMD grid, we can use multiple neuron
   ↪ cores
28     # without updating the original `nki_tensor_add_kernel_` kernel.
29     return nki_tensor_add_kernel_[nl.spmd_dim(grid_x, nl.nc(2)), grid_y](a_input, b_input)

```

In this example:

1. We reuse the NKI kernel in `nki_tensor_add_kernel_` which is decorated with the `nki.jit` decorator to call the nki compiler to compile the kernel.
2. Recall this kernel defines offsets into the tensors based on the ID of the worker executing the code (`nl.program_id`), and generates tile indices using these offsets with `nl.arange`.
3. Using SPMD execution as discussed in *SPMD: Launching Multiple Instances of a Kernel*, note that each worker only operates on a (sub-tensor) tile of the input/output tensors. By accessing its own `program_id`, each worker can calculate the offsets it needs to access the correct tiles.

4. When multiple Neuron Cores are specified in the SPMD launch grid, these tensors are further sharded across available cores. On Trainium 2, we have 2 local cores that have shared HBM.
5. As before, the first axis of the tensor (mapped to the partition-dimension) is tiled into blocks of 128, based on hardware restrictions (see *Tile Size Considerations*). The second axis (mapped to the free-dimension) is tiled into blocks of 512 (no tile-size constraint, since the addition operation is performed on the Vector engine, the only restriction is on-chip memory capacity).
6. `nl.store` for kernels running on both cores will write to an `c_output` in shared HBM, dramatically increasing the throughput of the computation.

## SPMD execution

1. We want to shard the workload across 2 cores, so for every `nl.nc(2)` we determine our initial `axis=0` to be 128 from the expected slice size in the kernel \* the number of cores = 256.
2. Thus we alter our previous sample and change `grid_x` to `a_input.shape[0] // (128 * 2)` to account for this.
3. Launch the kernel with launch grid `[nl.spmd_dim(grid_x, nl.nc(2)), grid_y]`

As before, we are using a two-dimensional grid where the first dimension of the tensor is tiled in the X dimension of the grid while the second dimension is tiled in the Y dimension of the grid. We similarly assume that tensor sizes are a multiple of maximum tile sizes allowed, so we do not need to handle partial tiles.

However, this time we also directly specify how each instance of our kernel will be distributed across multiple local Neuron Cores such that:

```
# Physical NC [0]: kernel[n, m] where n is 0 or even
# Physical NC [1]: kernel[n, m] where n is odd
```

## Launching kernel and testing correctness

To execute the kernel, we prepare tensors `a` and `b`, and call the `nki_tensor_add_nc2` helper function. We also verify the correctness of the NKI kernel against, torch by comparing the outputs of both, using `torch.allclose`:

```
1 import torch
2 from torch_xla.core import xla_model as xm
3
4 if __name__ == "__main__":
5     device = xm.xla_device()
6
7     a = torch.rand((512, 2048), dtype=torch.bfloat16).to(device=device)
8     b = torch.rand((512, 2048), dtype=torch.bfloat16).to(device=device)
9
10    output_nki = nki_tensor_add_nc2(a, b)
11    print(f"output_nki={output_nki}")
12
13    output_torch = a + b
14    print(f"output_torch={output_torch}")
15
16    allclose = torch.allclose(output_torch, output_nki, atol=1e-4, rtol=1e-2)
17    if allclose:
18        print("NKI and Torch match")
```

(continues on next page)

(continued from previous page)

```

19 else:
20     print("NKI and Torch differ")
21
22 assert allclose

```

Output:

```

2023-12-29 15:18:00.000558: 14283 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2023-12-29 15:18:00.000559: 14283 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↳'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↳workdir/49f554a2-2c55-4a88-8054-cc9f20824a46/model.MODULE_5007921933048625946+d41d8cd9.
↳hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/49f554a2-2c55-4a88-8054-
↳cc9f20824a46/model.MODULE_5007921933048625946+d41d8cd9.neff', '--verbose=35']
.
Compiler status PASS
output_nki=tensor([[1.459  1.488  1.607  ... 1.217  0.7354 1.457 ]
[1.793  0.7373 0.8877 ... 1.813  0.8936 1.39  ]
[0.7285 0.9473 1.531  ... 1.04   1.302  0.8413]
...
[0.7705 1.195  1.047  ... 1.307  0.588  0.7725]
[1.21   1.719  1.209  ... 1.171  0.583  0.5034]
[1.307  1.521  0.9526 ... 0.5825 1.518  0.673 ]],
device='xla:1', dtype=torch.bfloat16)
2023-12-29 15:18:02.000219: 14463 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2023-12-29 15:18:02.000220: 14463 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↳'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↳workdir/2e135b73-1c3b-45e4-a6f0-2c4b105c20e5/model.MODULE_
↳10032327759287407517+d41d8cd9.hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/
↳2e135b73-1c3b-45e4-a6f0-2c4b105c20e5/model.MODULE_10032327759287407517+d41d8cd9.neff',
↳'--verbose=35']
.
Compiler status PASS
output_torch=tensor([[1.459  1.488  1.607  ... 1.217  0.7354 1.457 ]
[1.793  0.7373 0.8877 ... 1.813  0.8936 1.39  ]
[0.7285 0.9473 1.531  ... 1.04   1.302  0.8413]
...
[0.7705 1.195  1.047  ... 1.307  0.588  0.7725]
[1.21   1.719  1.209  ... 1.171  0.583  0.5034]
[1.307  1.521  0.9526 ... 0.5825 1.518  0.673 ]],
device='xla:1', dtype=torch.bfloat16)
2023-12-29 15:18:03.000797: 14647 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2023-12-29 15:18:03.000798: 14647 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd: [
↳'neuronx-cc', '--target=trn1', 'compile', '--framework', 'XLA', '/tmp/neuroncc_compile_
↳workdir/74f8b6ae-76d9-4dd8-af7f-e5e1c40a27a3/model.MODULE_5906037506311912405+d41d8cd9.
↳hlo.pb', '--output', '/tmp/neuroncc_compile_workdir/74f8b6ae-76d9-4dd8-af7f-
↳e5e1c40a27a3/model.MODULE_5906037506311912405+d41d8cd9.neff', '--verbose=35']
.
Compiler status PASS
NKI and Torch match

```

Note that the tensor values you see will differ from what's printed above, since this example uses `torch.rand` to initialize the inputs.

## JAX

### Helper function and SPMD execution

We can reuse the same NKI compute kernel defined for PyTorch above and declare a helper function to launch the compute-kernel with appropriate grid/block sizes, to perform the computation:

```

1 def nki_tensor_add_nc2(a_input, b_input):
2     """NKI kernel caller to compute element-wise addition of two input tensors using
   ↪ multiple Neuron cores.
3
4     This kernel caller lifts tile-size restriction, by applying the kernel on tiles of the
   ↪ inputs/outputs.
5     a_input and b_input are sharded across Neuron cores, directly utilizing Trn2
   ↪ architecture capabilities
6
7     Args:
8         a_input: a first input tensor, of shape [N*128, M*512]
9         b_input: a second input tensor, of shape [N*128, M*512]
10
11     Returns:
12         a tensor of shape [N*128, M*512], the result of a_input + b_input
13     """
14
15     # The SPMD launch grid denotes the number of kernel instances.
16     # In this case, we use a 2D grid where the size of each invocation is 128x512
17     # Since we're sharding across neuron cores on the 1st dimension we want to do our
   ↪ slicing at
18     # 128 per core * 2 cores = 256
19     grid_x = a_input.shape[0] // (128 * 2)
20     grid_y = a_input.shape[1] // 512
21
22     # In addition, we distribute the kernel to physical neuron cores around the first
   ↪ dimension
23     # of the spmd grid.
24     # This means:
25     # Physical NC [0]: kernel[n, m] where n is even
26     # Physical NC [1]: kernel[n, m] where n is odd
27     # notice, by specifying this information in the SPMD grid, we can use multiple neuron
   ↪ cores
28     # without updating the original `nki_tensor_add_kernel` kernel.
29     return nki_tensor_add_kernel_[nl.spmd_dim(grid_x, nl.nc(2)), grid_y](a_input, b_input)

```

As before, we are using a two-dimensional grid where the first dimension of the tensor is tiled in the X dimension of the grid, while the second dimension is tiled in the Y dimension of the grid. We similarly assume that tensor sizes are a multiple of maximum tile sizes allowed, so we do not need to handle partial tiles.

However, this time we also directly specify how each instance of our kernel will be distributed across multiple local Neuron Cores such that:

```
# Physical NC [0]: kernel[n, m] where n is 0 or even
# Physical NC [1]: kernel[n, m] where n is odd
```

## Launching kernel and testing correctness

To execute the kernel, we prepare arrays `a` and `b`, and call the `nki_tensor_add_nc2` helper function. We also verify the correctness of the NKI kernel against, JAX by comparing the outputs of both, using `jax.numpy.allclose`:

```
1 import jax
2 import jax.numpy as jnp
3
4 if __name__ == "__main__":
5
6     seed_a, seed_b = jax.random.split(jax.random.PRNGKey(42))
7     a = jax.random.uniform(seed_a, (512, 2048), dtype=jnp.bfloat16)
8     b = jax.random.uniform(seed_b, (512, 2048), dtype=jnp.bfloat16)
9
10    output_nki = nki_tensor_add_nc2(a, b)
11    print(f"output_nki={output_nki}")
12
13    output_jax = a + b
14    print(f"output_jax={output_jax}")
15
16    allclose = jnp.allclose(output_jax, output_nki, atol=1e-4, rtol=1e-2)
17    if allclose:
18        print("NKI and JAX match")
19    else:
20        print("NKI and JAX differ")
21
22    assert allclose
```

Output:

```
.
Compiler status PASS
.
Compiler status PASS
.
Compiler status PASS
output_nki=[[0.992188 1.27344 1.65625 ... 0.90625 1.34375 1.77344]
 [0 0.90625 1.34375 ... 0.390625 0.703125 0.914062]
 [0.5 0.390625 0.703125 ... 1.22656 1.15625 1.01562]
 ...
 [1.98438 1.98438 1.98438 ... 1.33594 1.64062 1.35938]
 [0.992188 1.33594 1.64062 ... 1.16406 1.67188 1.20312]
 [1.49219 1.16406 1.67188 ... 1.375 1 1.6875]]
.
Compiler status PASS
output_jax=[[0.992188 1.27344 1.65625 ... 0.90625 1.34375 1.77344]
 [0 0.90625 1.34375 ... 0.390625 0.703125 0.914062]
 [0.5 0.390625 0.703125 ... 1.22656 1.15625 1.01562]
 ...
 .
```

(continues on next page)



(continued from previous page)

```
[1.98438 1.98438 1.98438 ... 1.33594 1.64062 1.35938]
[0.992188 1.33594 1.64062 ... 1.16406 1.67188 1.20312]
[1.49219 1.16406 1.67188 ... 1.375 1 1.6875]]
```

Compiler status PASS

NKI **and** JAX match

Note that the array values you see will differ from what's printed above, since this example uses `jax.random.uniform` to initialize the inputs.

## Download all source code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- **NKI baremetal implementation: `spmd_multiple_nc_tensor_addition_nki_kernels.py`**
  - You must also download `spmd_tensor_addition_nki_kernels.py` into the same folder to run this script.
- **PyTorch implementation: `spmd_multiple_nc_tensor_addition_torch.py`**
  - You must also download `spmd_multiple_nc_tensor_addition_nki_kernels.py` and `spmd_tensor_addition_nki_kernels.py` into the same folder to run this PyTorch script.
- **JAX implementation: `spmd_multiple_nc_tensor_addition_jax.py`**
  - You must also download `spmd_multiple_nc_tensor_addition_nki_kernels.py` and `spmd_tensor_addition_nki_kernels.py` into the same folder to run this PyTorch script.

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

Run NKI baremetal implementation:

```
python3 spmd_multiple_nc_tensor_addition_nki_kernels.py
```

Run PyTorch implementation:

```
python3 spmd_multiple_nc_tensor_addition_torch.py
```

Run JAX implementation:

```
python3 spmd_multiple_nc_tensor_addition_jax.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Transpose2D

In this tutorial, we transpose a tensor along two of its axes using NKI. In doing so, we learn about:

- The NKI syntax and programming model.
- Multi-dimensional memory address patterns in NKI.

As background, there are two main types of transposition in NKI:

1. Transposition between the partition-dimension axis and one of the free-dimension axes, which is achieved via the `nki.isa.nc_transpose` instruction.
2. Transposition between two axes on the free-dimension, which is achieved via a `nki.language.copy` instruction, with indexing manipulation in the free axis to re-arrange the data.

In this example, we'll focus on the second case: consider a three-dimensional input tensor  $[P, F1, F2]$ , where the  $P$  axis is mapped to the different SBUF partitions and the  $F1$  and  $F2$  axes are flattened and placed in each partition, with  $F1$  being the major dimension. Our goal in this example is to transpose the  $F1$  and  $F2$  axes with a parallel dimension  $P$ , to re-arrange the data within each partition. [Figure](#) below illustrates the input and output tensor layouts.

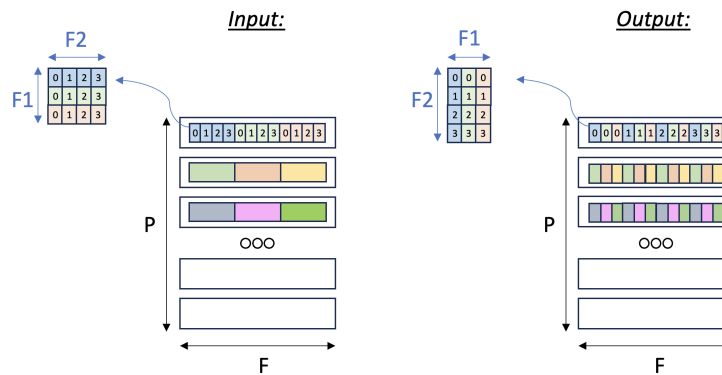


Fig. 7.67: Tensor F1:F2 Transpose

## PyTorch

### Compute kernel

```

1  import neuronxcc.nki as nki
2  import neuronxcc.nki.language as nl
3
4
5  @nki.jit
6  def tensor_transpose2D_kernel_(in_tensor, shape2D):
7      """
8      NKI kernel to reorder the elements on axis[1] of the input tensor.
9
10     Every row of the input tensor is a flattened row-major 2D matrix.
11     The shape2D argument defines the dimensions of the flattened matrices (#rows,#cols).
12     Our goal in this kernel is to transpose these flattened 2D matrices, i.e. make them (
13     ↪ #cols,#rows).

```

(continues on next page)

(continued from previous page)

*Example:*

```
in_tensor = [a0,a1,a2,a3,b0,b1,b2,b3,c0,c1,c2,c3]
shape2D = (3,4)
```

*this means that in\_tensor has 3 rows and 4 columns, i.e. can be represented as:*

```
[a0,a1,a2,a3]
[b0,b1,b2,b3]
[c0,c1,c2,c3]
```

*after transpose, we expect to get:*

```
[a0,b0,c0]
[a1,b1,c1]
[a2,b2,c2]
[a3,b3,c3]
```

*Thus, out\_tensor is expected to be [a0,b0,c0,a1,b1,c1,a2,b2,c2,a3,b3,c3]**Args:*

```
in_tensor: an input tensor
shape2D: tuple representing the dimensions to be transposed: (#rows, #cols)
out_tensor: an output (transposed) tensor
"""
```

```
out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
                        buffer=nl.shared_hbm)
```

*# Gather input shapes*

```
sz_p, _ = in_tensor.shape
```

*# Load input data from external memory to on-chip memory*

```
in_tile = nl.load(in_tensor)
```

*# Performing f1/f2 transpose**# =====**# The desired transpose pattern is provided as an input:*

```
sz_f1, sz_f2 = shape2D
```

*# We're going to need 3 indices to perform f1:f2 transpose.**# - i\_p0 is the parallel index**# - i\_f1 and i\_f2 are both free-dim indices, and will be used to transpose between the f1/f2 axes*

```
i_p0 = nl.arange(sz_p)[: , None, None]
```

```
i_f1 = nl.arange(sz_f1)[None, :, None]
```

```
i_f2 = nl.arange(sz_f2)[None, None, :]
```

*# Perform the transposition via a SBUF-to-SBUF copy, with access-pattern manipulation**# Note that we have 2D tensors and 3 indices, since we need to represent a 2D access pattern \*per partition\***# RHS traverses an F1 x F2 matrix in a row major manner**# LHS traverses an F2 x F1 (new) matrix in a row major manner*

```
out_tile = nl.ndarray(shape=(sz_p, sz_f2*sz_f1), dtype=out_tensor.dtype)
```

```
out_tile[i_p0, i_f2*sz_f1+i_f1] = nl.copy(in_tile[i_p0, i_f1*sz_f2+i_f2])
```

*# Finally, we store out\_tile to external memory*

```
nl.store(out_tensor, value=out_tile)
```

```
return out_tensor
```

## Launching kernel and testing correctness

To execute the kernel, we prepare tensors `a` and call `tensor_transpose2D_kernel_`:

```

1 import torch
2 from torch_xla.core import xla_model as xm
3
4 if __name__ == "__main__":
5     device = xm.xla_device()
6
7     P, X, Y = 5, 3, 4
8     a = torch.arange(P*X*Y, dtype=torch.int8).reshape((P, X*Y)).to(device=device)
9     a_t_nki = torch.zeros((P, Y*X), dtype=torch.int8).to(device=device)
10
11     a_t_nki = tensor_transpose2D_kernel_(a, (X, Y))
12
13     a_t_torch = torch.transpose(a.reshape(P, X, Y), 1, 2).reshape(P, X * Y)
14
15     print(a, a_t_nki, a_t_torch)
16
17     allclose = torch.allclose(a_t_torch, a_t_nki)
18     if allclose:
19         print("NKI and PyTorch match")
20     else:
21         print("NKI and PyTorch differ")
22
23     assert allclose

```

## JAX

### Compute kernel

We can reuse the same NKI compute kernel defined for PyTorch above.

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3
4
5 @nki.jit
6 def tensor_transpose2D_kernel_(in_tensor, shape2D):
7     """
8     NKI kernel to reorder the elements on axis[1] of the input tensor.
9
10    Every row of the input tensor is a flattened row-major 2D matrix.
11    The shape2D argument defines the dimensions of the flattened matrices (#rows,#cols).
12    Our goal in this kernel is to transpose these flattened 2D matrices, i.e. make them (
    ↪ #cols,#rows).
13
14    Example:
15        in_tensor = [a0,a1,a2,a3,b0,b1,b2,b3,c0,c1,c2,c3]
16        shape2D = (3,4)

```

(continues on next page)

(continued from previous page)

```

17  this means that in_tensor has 3 rows and 4 columns, i.e. can be represented as:
18      [a0,a1,a2,a3]
19      [b0,b1,b2,b3]
20      [c0,c1,c2,c3]
21  after transpose, we expect to get:
22      [a0,b0,c0]
23      [a1,b1,c1]
24      [a2,b2,c2]
25      [a3,b3,c3]
26  Thus, out_tensor is expected to be [a0,b0,c0,a1,b1,c1,a2,b2,c2,a3,b3,c3]
27
28  Args:
29      in_tensor: an input tensor
30      shape2D: tuple representing the dimensions to be transposed: (#rows, #cols)
31      out_tensor: an output (transposed) tensor
32  """
33  out_tensor = nl.ndarray(in_tensor.shape, dtype=in_tensor.dtype,
34                          buffer=nl.shared_hbm)
35  # Gather input shapes
36  sz_p, _ = in_tensor.shape
37
38  # Load input data from external memory to on-chip memory
39  in_tile = nl.load(in_tensor)
40
41  # Performing f1/f2 transpose
42  # =====
43  # The desired transpose pattern is provided as an input:
44  sz_f1, sz_f2 = shape2D
45
46  # We're going to need 3 indices to perform f1:f2 transpose.
47  # - i_p0 is the parallel index
48  # - i_f1 and i_f2 are both free-dim indices, and will be used to transpose between the
49  ↪ f1/f2 axes
50  i_p0 = nl.arange(sz_p)[: , None, None]
51  i_f1 = nl.arange(sz_f1)[None, :, None]
52  i_f2 = nl.arange(sz_f2)[None, None, :]
53
54  # Perform the transposition via a SBUF-to-SBUF copy, with access-pattern manipulation
55  # Note that we have 2D tensors and 3 indices, since we need to represent a 2D access
56  ↪ pattern *per partition*
57  # RHS traverses an F1 x F2 matrix in a row major manner
58  # LHS traverses an F2 x F1 (new) matrix in a row major manner
59  out_tile = nl.ndarray(shape=(sz_p, sz_f2*sz_f1), dtype=out_tensor.dtype)
60  out_tile[i_p0, i_f2*sz_f1+i_f1] = nl.copy(in_tile[i_p0, i_f1*sz_f2+i_f2])
61
62  # Finally, we store out_tile to external memory
63  nl.store(out_tensor, value=out_tile)
64
65  return out_tensor

```

## Launching kernel and testing correctness

To execute the kernel, we prepare array `a` and call `tensor_transpose2D_kernel_`:

```

1 import jax
2 import jax.numpy as jnp
3
4 if __name__ == "__main__":
5     P, X, Y = 5, 37, 44
6     a = jax.random.uniform(jax.random.PRNGKey(42), (P, X * Y))
7     a_t_nki = tensor_transpose2D_kernel_(a, shape2D=(X, Y))
8
9     a_t_jax = jnp.transpose(a.reshape(P, X, Y), axes=(0, 2, 1)).reshape(P, X * Y)
10    print(a, a_t_nki, a_t_jax)
11
12    allclose = jnp.allclose(a_t_jax, a_t_nki)
13    if allclose:
14        print("NKI and JAX match")
15    else:
16        print("NKI and JAX differ")
17
18    assert allclose

```

---

**Note:** We pass `shape2D` as kwargs to pass the shape as a compile-time constant to the kernel function.

---

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- NKI baremetal implementation: `transpose2d_nki_kernels.py`
- **PyTorch implementation: `transpose2d_torch.py`**
  - You must also download `transpose2d_nki_kernels.py` into the same folder to run this PyTorch script.
- **JAX implementation: `transpose2d_jax.py`**
  - You must also download `transpose2d_nki_kernels.py` into the same folder to run this JAX script.

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

Run NKI baremetal implementation:

```
python3 transpose2d_nki_kernels.py
```

Run PyTorch implementation:

```
python3 transpose2d_torch.py
```

Run JAX implementation:

```
python3 transpose2d_jax.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## AveragePool2D

In this tutorial, we examine a case of dimensionality reduction. We implement a 2D AveragePool operation, which is used in many vision neural networks. In doing so, we learn about:

- NKI syntax and programming model.
- multi-dimensional memory access patterns in NKI.

The 2D AveragePool operation takes  $C \times [H, W]$  matrices and reduces each matrix along the H and W axes. To leverage free-dimension flexible indexing, we can map the C (parallel) axis to the P dimension and H/W (contraction) axes to the F dimension. Performing such a 2D pooling operation requires a 4D memory access pattern in the F dimension, with reduction along two axes. [Figure](#) below illustrates the input and output tensor layouts.

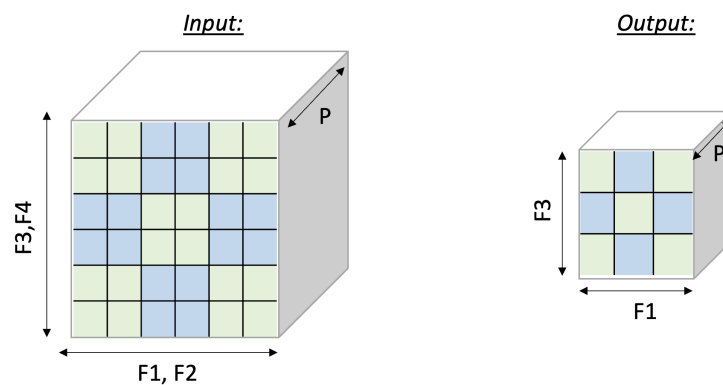


Fig. 7.68: 2D-Pooling Operation (reducing on axes F2 and F4)

## PyTorch

### Compute kernel

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3 from neuronxcc.nki.typing import tensor
4
5 @nki.jit
6 def tensor_avgpool_kernel(in_tensor, pool_size):
7     """NKI kernel to compute a 2D avg-pool operation
8
9     Args:
10         in_tensor: an input tensor, of shape C x H x W
11         pool_size: an integer representing a (square) pool-window size

```

(continues on next page)

(continued from previous page)

```

12
13 Return:
14     out_tensor: the resulting output tensor, of shape C x (H/pool_size) x (W/pool_size)
15     """
16
17     # Get input/output dimensions
18     sz_cin, sz_hin, sz_win = in_tensor.shape
19     sz_hout = sz_hin // pool_size
20     sz_wout = sz_win // pool_size
21     # Create output tensor shared between all SPMD instances as result tensor
22     out_tensor = nl.ndarray((sz_cin, sz_hout, sz_wout), dtype=in_tensor.dtype,
23                             buffer=nl.shared_hbm)
24
25     # Set relevant sizes
26     sz_p = sz_cin
27     sz_pool = pool_size
28
29     # Generate pool index patterns (requires two extra dimensions, for the pool window)
30     i0, i1, i2, i3, i4 = nl.mgrid[0:sz_p, 0:sz_hin//sz_pool, 0:sz_pool, 0:sz_win//sz_pool,
31     ↪ 0:sz_pool]
32
33     # Load input data from external memory to on-chip memory
34     in_tile: tensor[sz_p, sz_hin, sz_win] = nl.load(in_tensor)
35
36     # Perform the pooling operation:
37     # We use numpy's advanced indexing, in order to extend in_tile to 5D, and then reduce-
38     ↪ average two dimension.
39     # axis[0] is the index for p_dim, and thus doesn't participate in the reduction
40     ↪ operation.
41     # axis[1] and axis[2] together index the rows, with axis[2] responsible for inner
42     ↪ strides
43     # (i.e. inside a pooling window), and axis[1] responsible for the outer strides. As
44     ↪ such, we reduce over axis[2].
45     # Similarly, axis[3] and axis[4] together index the columns, and we thus reduce over
46     ↪ axis[4].
47     out_tile : tensor[sz_p, sz_hout, sz_wout] = nl.sum(in_tile[i0, sz_pool*i1+i2, sz_
48     ↪ pool*i3+i4],
49                                                         axis=[2,4]) / (pool_size*pool_size)
50
51     # Store the results back to hbm
52     nl.store(out_tensor, value=out_tile)
53
54     # Transfer the ownership of `out_tensor` to the caller
55     return out_tensor

```



## Launching kernel and testing correctness

To execute the kernel, we prepare tensors `in_tensor` and call `tensor_avgpool_kernel`:

```

1 import torch
2 from torch_xla.core import xla_model as xm
3
4 if __name__ == "__main__":
5     device = xm.xla_device()
6
7     # Now let's run the kernel
8     POOL_SIZE = 2
9     C, HIN, WIN = 2, 6, 6
10    HOUT, WOUT = HIN//POOL_SIZE, WIN//POOL_SIZE
11
12    in_tensor = torch.arange(C * HIN * WIN, dtype=torch.bfloat16).reshape(C, HIN, WIN).
13    ↪to(device=device)
14    out_nki = torch.zeros((C, HOUT, WOUT), dtype=torch.bfloat16).to(device=device)
15
16    out_nki = tensor_avgpool_kernel(in_tensor, POOL_SIZE)
17
18    out_torch = torch.nn.functional.avg_pool2d(in_tensor, POOL_SIZE, POOL_SIZE)
19
20    print(in_tensor, out_nki, out_torch) # an implicit XLA barrier/mark-step
21
22    if (out_nki == out_torch).all():
23        print("NKI and Torch match")
24    else:
25        print("NKI and Torch differ")

```

## JAX

### Compute kernel

Let's reuse the same NKI kernel implementation defined for PyTorch above:

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3 from neuronxcc.nki.typing import tensor
4
5 @nki.jit
6 def tensor_avgpool_kernel(in_tensor, pool_size):
7     """NKI kernel to compute a 2D avg-pool operation
8
9     Args:
10         in_tensor: an input tensor, of shape C x H x W
11         pool_size: an integer representing a (square) pool-window size
12
13     Return:
14         out_tensor: the resulting output tensor, of shape C x (H/pool_size) x (W/pool_size)
15     """

```

(continues on next page)

(continued from previous page)

```

16
17 # Get input/output dimensions
18 sz_cin, sz_hin, sz_win = in_tensor.shape
19 sz_hout = sz_hin // pool_size
20 sz_wout = sz_win // pool_size
21 # Create output tensor shared between all SPMD instances as result tensor
22 out_tensor = nl.ndarray((sz_cin, sz_hout, sz_wout), dtype=in_tensor.dtype,
23                          buffer=nl.shared_hbm)
24
25 # Set relevant sizes
26 sz_p = sz_cin
27 sz_pool = pool_size
28
29 # Generate pool index patterns (requires two extra dimensions, for the pool window)
30 i0, i1, i2, i3, i4 = nl.mgrid[0:sz_p, 0:sz_hin//sz_pool, 0:sz_pool, 0:sz_win//sz_pool,
31 ↪ 0:sz_pool]
32
33 # Load input data from external memory to on-chip memory
34 in_tile: tensor[sz_p, sz_hin, sz_win] = nl.load(in_tensor)
35
36 # Perform the pooling operation:
37 # We use numpy's advanced indexing, in order to extend in_tile to 5D, and then reduce-
38 ↪ average two dimension.
39 # axis[0] is the index for p_dim, and thus doesn't participate in the reduction
40 ↪ operation.
41 # axis[1] and axis[2] together index the rows, with axis[2] responsible for inner
42 ↪ strides
43 # (i.e. inside a pooling window), and axis[1] responsible for the outer strides. As
44 ↪ such, we reduce over axis[2].
45 # Similarly, axis[3] and axis[4] together index the columns, and we thus reduce over
46 ↪ axis[4].
47 out_tile : tensor[sz_p, sz_hout, sz_wout] = nl.sum(in_tile[i0, sz_pool*i1+i2, sz_
48 ↪ pool*i3+i4],
49                                                    axis=[2,4]) / (pool_size*pool_size)
50
51 # Store the results back to hbm
52 nl.store(out_tensor, value=out_tile)
53
54 # Transfer the ownership of `out_tensor` to the caller
55 return out_tensor

```

In order to pass `pool_size` as a compile time constant, we pass `pool_size` as kwargs.

```
out_nki = tensor_avgpool_kernel(in_array, pool_size=POOL_SIZE)
```

We write a reference JAX implementation of `AveragePool2D` as JAX does not have a primitive for it.

```

1 import jax.numpy as jnp
2
3 # Reference JAX implementation
4 def jax_average_pool_2D(in_tensor, pool_size):
5     c, h_in, w_in = in_tensor.shape

```

(continues on next page)

(continued from previous page)

```

6   reshaped = in_tensor.reshape(c, h_in // pool_size, pool_size, w_in // pool_size, pool_
   ↪size)
7   return jnp.nanmean(reshaped, axis=(2, 4))

```

## Launching kernel and testing correctness

To execute the kernel, we prepare array `in_array` and invoke the kernel caller function `tensor_avgpool_kernel`:

```

1  if __name__ == "__main__":
2      POOL_SIZE = 2
3      C, HIN, WIN = 2, 6, 6
4      HOUT, WOUT = HIN//POOL_SIZE, WIN//POOL_SIZE
5
6      in_array = jnp.arange(C * HIN * WIN, dtype=jnp.float32).reshape(C, HIN, WIN)
7
8      out_nki = tensor_avgpool_kernel(in_array, pool_size=POOL_SIZE)
9      out_jax = jax_average_pool_2D(in_array, pool_size=POOL_SIZE)
10
11     print(in_array, out_nki, out_jax)
12
13     if jnp.allclose(out_nki, out_jax):
14         print("NKI and JAX match")
15     else:
16         print("NKI and JAX differ")

```

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- NKI baremetal implementation: `average_pool2d_nki_kernels.py`
- **PyTorch implementation: `average_pool2d_torch.py`**
  - You must also download `average_pool2d_nki_kernels.py` into the same folder to run this PyTorch script.
- **JAX implementation: `average_pool2d_jax.py`**
  - You must also download `average_pool2d_nki_kernels.py` into the same folder to run this JAX script.

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

Run NKI baremetal implementation:

```
python3 average_pool2d_nki_kernels.py
```

Run PyTorch implementation:

```
python3 average_pool2d_torch.py
```

Run JAX implementation:

```
python3 average_pool2d_jax.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Matrix multiplication

In this tutorial, we will start with a simple NKI matrix multiplication kernel and optimize it step by step. In doing so, we learn about:

- The NKI syntax and programming model.
- Layout, tiling, and memory management considerations when performing matrix multiplication in NKI.
- Best practices for validating and benchmarking your custom kernel against a reference native torch implementation.

## Basic compute kernel

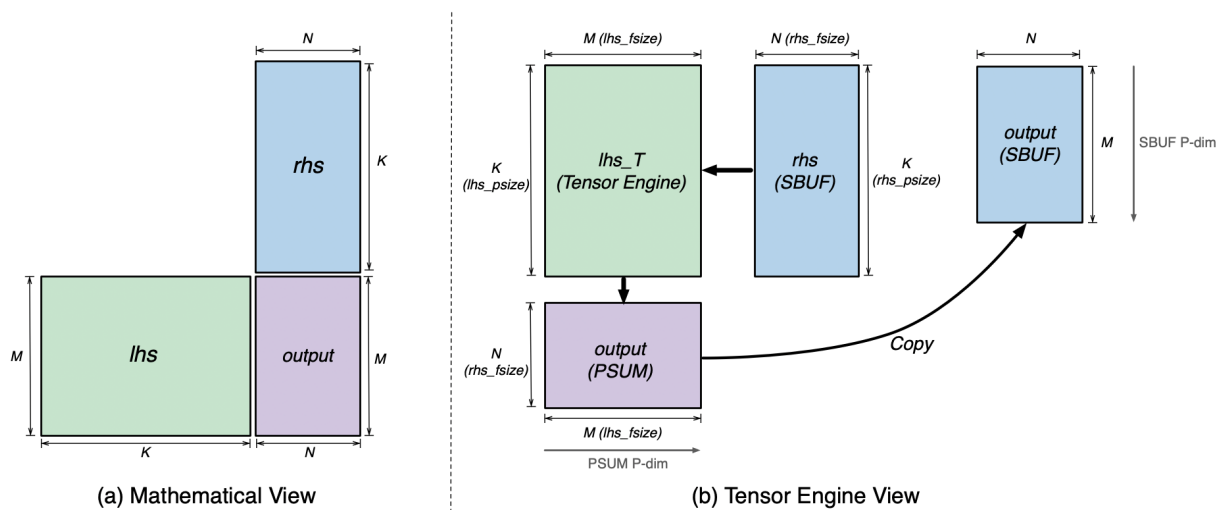


Fig. 7.69:  $M \times K \times N$  Matrix Multiplication Visualization

Fig. 7.69 illustrates how a simple matrix multiplication:  $\text{lhs} [M, K] * \text{rhs} [K, N] = \text{output} [M, N]$  would be mapped to the Tensor Engine (TensorE) and SRAMs from its original mathematical view. Note, the PSUM partition dimension is rotated 90 degrees from SBUF partition dimension solely for layout visualization. The copy preserves

the output tile layout from PSUM to SBUF, by copying data from each PSUM partition to the corresponding SBUF partition.

The NKI example below implements a compute kernel for a single-tile matrix multiplication. It computes a 64(M) x 128(K) x 512 (N) matrix multiplication operation.

```

1 @nki.jit
2 def nki_matmul_basic(lhsT, rhs):
3     """NKI kernel to compute a 64x128x512 matrix multiplication operation
4
5     Args:
6         lhsT: an input tensor of shape [128,64], a left hand side argument of the
7             matrix multiplication, delivered transposed for optimal performance
8         rhs: an input tensor of shape [128,512], a right hand side argument of the
9             matrix multiplication
10    Returns:
11        result: the resulting output tensor of shape [64,512]
12    """
13    result = nl.ndarray((64, 512), dtype=lhsT.dtype, buffer=nl.shared_hbm)
14
15    # Defining indexes for input LHS.T
16    # - Note: here we take LayoutConstraint #1 into account:
17    # "For MatMult, contraction axis must be mapped to P-dim"
18    i_lhsT_p, i_lhsT_f = nl.mgrid[0:128, 0:64]
19
20    # Defining indexes for input RHS
21    # - Note: here we take LayoutConstraint #1 into account:
22    # "For MatMult, contraction axis must be mapped to P-dim"
23    i_rhs_p, i_rhs_f = nl.mgrid[0:128, 0:512]
24
25    # Defining indexes for the output ([64,128]@[128,512] -> [64,512])
26    i_out_p, i_out_f = nl.mgrid[0:64, 0:512]
27
28    # Loading the inputs (HBM->SBUF)
29    # Note: here we take Tile dtype definition into account,
30    # which forces P-dim as the left most index
31    lhs_tile = nl.load(lhsT[i_lhsT_p, i_lhsT_f])
32    rhs_tile = nl.load(rhs[i_rhs_p, i_rhs_f])
33
34    # Perform the matrix-multiplication
35    # Note1: We set transpose_x to True, to indicate that the LHS input is transposed
36    # Note2: A NKI matmul instruction always writes to PSUM in float32 data-type
37    result_psum = nl.matmul(lhs_tile, rhs_tile, transpose_x=True)
38
39    # Copy the result from PSUM back to SBUF, and cast to expected output data-type
40    result_sbuf = nl.copy(result_psum, dtype=result.dtype)
41
42    # The result of a [64,128] x [128,512] matrix multiplication has a shape of [64, 512].
43    # This dictates which indices to use to address the result tile.
44    nl.store(result[i_out_p, i_out_f], value=result_sbuf)
45
46    return result

```

In this example, we define the NKI kernel as `nki_matmul_basic_`:

1. We define indices to access the LHS and RHS input tensors.
2. To adhere to NKI's layout considerations (*Layout Considerations*), we map the contraction axis of both LHS and RHS to the P-dimension, which means we load LHS in transposed form.
3. To adhere to NKI's tile size considerations (*Tile Size Considerations*), we limit the matmul instruction arguments to tiles of up to [128, 128] for LHS, and [128, 512] for RHS.
4. Using the `nl.load` operation, we load the inputs from HBM tensors to SBUF tiles.
5. We then use the `nl.matmul` operation to perform the matrix multiplication. Note that we set the `transpose_x` argument to `True`, since the LHS argument is transposed. Also note that the `64x128` dimension here actually under-utilizes the TensorE, but it helps to distinguish the M, K and N dimensions for education purposes in this first code example.
6. `nl.matmul` always writes its result to PSUM, and since `nl.store` only moves data from SBUF to HBM, we copy the multiplication result from PSUM back to SBUF using `nl.copy`.

We can then execute the kernel and verify correctness against the torch implementation as follows. Note that we use `torch.allclose` to tolerate numerical error inherent to floating-point arithmetic.

```

1 device = xm.xla_device()
2 cpu = torch.device('cpu')
3
4 # Test the small workload with basic kernel
5 lhs_small = torch.rand((64, 128), dtype=torch.bfloat16, device=device)
6 rhs_small = torch.rand((128, 512), dtype=torch.bfloat16, device=device)
7
8 # Run NKI kernel
9 output_small = nki_matmul_basic_(lhs_small.T, rhs_small)
10
11 # Run torch reference
12 output_small_torch = torch.matmul(lhs_small, rhs_small)
13
14 # Compare results
15 print("Checking correctness of nki_matmul_basic")
16 if torch.allclose(output_small_torch, output_small, atol=1e-4, rtol=1e-2):
17     print("NKI and Torch match")
18 else:
19     print("NKI and Torch differ")

```

## Tiling matrix multiplications

So far, we've limited our matrix multiplication to the tile sizes allowed by NKI's tile size and layout constraints. Next, we'll see how to handle larger matrix multiplications. Let's start with a pseudo-code for tiling an  $[M, K]$  @  $[K, N]$  matrix-multiplication. Note that we assume the left-hand-side matrix ( $[M, K]$ ) is already transposed to  $LHS\_T$  ( $[K, M]$ ) for optimal performance of the underlying TensorE.

```

# LHS_T: left-hand-side matmul argument (shape [K,M])
# RHS: right-hand-side matmul argument (shape [K,N])
# RES: matmul result (shape [M,N])

# Tile LHS_T free dimension
for m in range(0, M, 128):
    # Tile RHS free dimension

```

(continues on next page)

(continued from previous page)

```

for n in range(0, N, 512):
    # Zero-out the accumulator buffer
    accum = zeros((128, 512))
    # Tile contraction dimension
    for k in range(0, K, 128):
        lhsT_tile = LHS_T[m : m+128, k : k+128]
        rhs_tile = RHS[k : k+128, n : n+512]
        accum += dot(lhsT_tile, rhs_tile)
    RES[m : m+128, n : n+512] = accum

```

This form of tiling can be achieved in NKI as follows:

```

1 @nki.jit
2 def nki_matmul_tiled(lhsT, rhs):
3     """NKI kernel to compute a matrix multiplication operation in a tiled manner
4
5     Args:
6         lhsT: an input tensor of shape [K,M], where both K and M are multiples for
7             128. It is the left-hand-side argument of the matrix multiplication,
8             delivered transposed for optimal performance.
9         rhs: an input tensor of shape [K,N], where K is a multiple of 128, and N
10            is a multiple of 512. It is the right-hand-side argument of the matrix
11            multiplication.
12     Returns:
13         result: the resulting output tensor of shape [M,N]
14     """
15
16     K, M = lhsT.shape
17     K_, N = rhs.shape
18     assert K == K_, "lhsT and rhs must have the same contraction dimension"
19     result = nl.ndarray((M, N), dtype=lhsT.dtype, buffer=nl.shared_hbm)
20
21     TILE_M = nl.tile_size.gemm_stationary_fmax # 128
22     TILE_K = nl.tile_size.pmax # 128
23     TILE_N = nl.tile_size.gemm_moving_fmax # 512
24
25     # Use affine_range to loop over tiles
26     for m in nl.affine_range(M // TILE_M):
27         for n in nl.affine_range(N // TILE_N):
28             # Allocate a tensor in PSUM
29             res_psum = nl.zeros((TILE_M, TILE_N), nl.float32, buffer=nl.psum)
30
31             for k in nl.affine_range(K // TILE_K):
32                 # Declare the tiles on SBUF
33                 lhsT_tile = nl.ndarray((TILE_K, TILE_M), dtype=lhsT.dtype, buffer=nl.sbuf)
34                 rhs_tile = nl.ndarray((TILE_K, TILE_N), dtype=rhs.dtype, buffer=nl.sbuf)
35
36                 # Load tiles from lhsT and rhs
37                 lhsT_tile[...] = nl.load(lhsT[k * TILE_K:(k + 1) * TILE_K,
38                                         m * TILE_M:(m + 1) * TILE_M])
39                 rhs_tile[...] = nl.load(rhs[k * TILE_K:(k + 1) * TILE_K,
40                                         n * TILE_N:(n + 1) * TILE_N])

```

(continues on next page)

(continued from previous page)

```

41
42     # Accumulate partial-sums into PSUM
43     res_psum += nl.matmul(lhsT_tile[...], rhs_tile[...], transpose_x=True)
44
45     # Copy the result from PSUM back to SBUF, and cast to expected output data-type
46     res_sb = nl.copy(res_psum, dtype=result.dtype)
47     nl.store(result[m * TILE_M:(m + 1) * TILE_M, n * TILE_N:(n + 1) * TILE_N],
48             value=res_sb)
49
50     return result

```

A few notes about the above code example:

First, in current release of NKI, the following NKI code pattern is the only way to trigger PSUM accumulation for matmuls on TensorEngine reliably:

```

# condition 1: a psum buffer with zeros
psum_buf = nl.zeros(..., buffer=nl.psum)

# condition 2: an affine range loop
for i in nl.affine_range(N):
    # condition 3: add matmul results from TensorEngine
    psum_buf += nl.matmul(stationary_tile, moving_tile) # or nisa.nc_matmul

```

The `nki_matmul_tiled_` kernel meets all three conditions above, and so do the kernels in the rest of this tutorial. The use of *PSUM accumulation architecture feature* is critical to achieve good performance out of TensorEngine when the contraction dimension of the matmul is greater than 128.

Second, note the use of `nl.mgrid` to define indices, this is the same as the `mgrid` in NumPy. It is similar to the other way to define indexes through `nl.arange` but it enables a more concise way to introduce indexes from multiple dimensions. `nl.affine_range` is used to define loop-level iterators, which is the recommended iterator type when the loop does not have loop-carried dependency (Note, associative reductions are not considered loop carried dependencies in this context).

Finally, there is an alternative way to implement this tiled matrix multiplication kernel using the SPMD programming model. We can use the SPMD model to launch  $(M/128) \times (N/512)$  instances of the kernel to complete the innermost loop. For more details, refer to the *SPMD programming model*.

## Optimization 1: Removing Redundant Loads

Currently, every `nl.matmul` is accompanied with two `nl.load` calls in the inner loop, both of which move data from HBM to SBUF. Let's introduce a metric, arithmetic intensity, to help understand why this is problematic. The arithmetic intensity of a workload is defined as the number of computation operations performed per byte of data accessed from HBM on average. The reason why we do not consider data accessed from SBUF in this metric is because the SBUF bandwidth (~20x higher than HBM) is high enough to sustain the peak computation throughput in TensorE.

Fig. 7.70 shows the roofline model, which models the relationship between arithmetic intensity of a workload and its achievable performance on a given computing platform. To saturate TensorE in a NeuronCore-v2, the arithmetic intensity threshold of a workload is 222 Flops/Byte for `bfloat16` data type. Inside the inner loop of `nki_matmul_tiled_`, accessing `lhsT_tile` and `rhs_tile` requires 160 KB of data read from HBM, while the `nl.matmul` call involves 16 MFlops. This leads to an arithmetic intensity of 102, which is significantly lower than the saturation threshold of 222. Therefore, `nki_matmul_tiled_` operates in the memory bound region of the roofline model and under-utilizes TensorE. To make the best out of TensorE, we need to improve the arithmetic intensity of the matmul kernel.



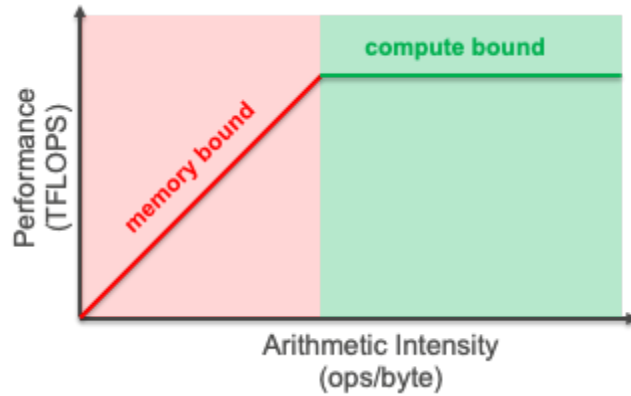


Fig. 7.70: Roofline Model: The Relationship Between Arithmetic Intensity and Performance

With NKI, programmers can control when and how to load data from HBM into SBUF and also perform computation. We will demonstrate in the upcoming steps how to increase the arithmetic intensity of the matmul kernel using NKI, thereby maximizing the utilization of TensorE.

First, we notice that in `nki_matmul_tiled_`, the same tiles from `lhsT` and `rhs` matrices are loaded more than once across different iterations of the inner loop. The following example reduces these redundant loads through hoisting them out of the innermost loop.

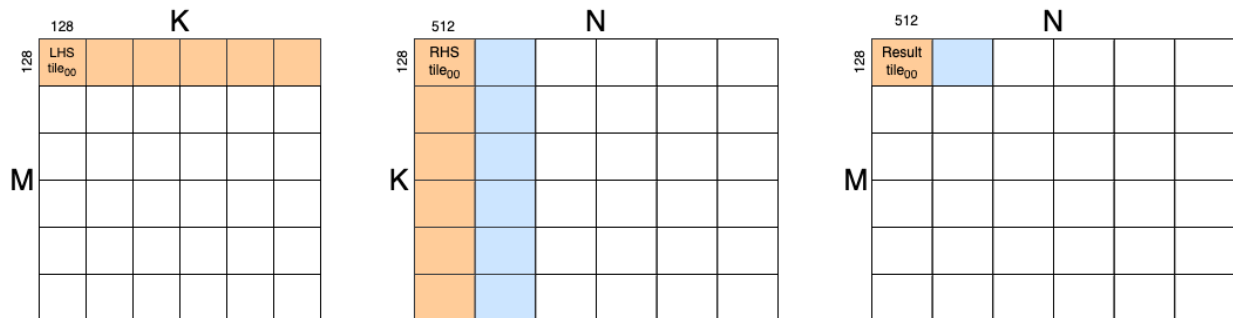


Fig. 7.71: Memory Pattern After Hoisting Loads Out of the Innermost Loop

```

1 @nki.jit
2 def nki_matmul_hoist_load(lhsT, rhs):
3     """NKI kernel to compute a matrix multiplication operation in a tiled manner
4     while hoisting the load of the lhsT and rhs to outer loops.
5
6     Args:
7         lhsT: an input tensor of shape [K,M], where both K and M are multiples for
8         128. It is the left-hand-side argument of the matrix multiplication,
9         delivered transposed for optimal performance.
10        rhs: an input tensor of shape [K,N], where K is a multiple of 128, and N
11        is a multiple of 512. It is the right-hand-side argument of the matrix
12        multiplication.
13    Returns:
14        result: the resulting output tensor of shape [M,N]
15    """
16

```

(continues on next page)

(continued from previous page)

```

17 K, M = lhsT.shape
18 K_, N = rhs.shape
19 assert K == K_, "lhsT and rhs must have the same contraction dimension"
20 result = nl.ndarray((M, N), dtype=lhsT.dtype, buffer=nl.shared_hbm)
21
22 TILE_M = nl.tile_size.gemm_stationary_fmax # 128
23 TILE_K = nl.tile_size.pmax # 128
24 TILE_N = nl.tile_size.gemm_moving_fmax # 512
25
26 # Define the indices (shape) of the tiles
27 i_lhsT = nl.mgrid[0:TILE_K, 0:TILE_M]
28 i_rhs = nl.mgrid[0:TILE_K, 0:TILE_N]
29 i_res = nl.mgrid[0:TILE_M, 0:TILE_N]
30
31 # Use affine_range to loop over tiles
32 for m in nl.affine_range(M // TILE_M):
33     # Load a whole column tiles from lhsT (with K * TILE_N numbers)
34     # This corresponds to the whole row in the original lhs
35     lhsT_tiles = nl.ndarray((K // TILE_K, nl.par_dim(TILE_K), TILE_N),
36                             dtype=lhsT.dtype,
37                             buffer=nl.sbuf)
38
39     for k in nl.affine_range(K // TILE_K):
40         # use `.p` for partition dimension and `.x` for the first free dimension
41         lhsT_tiles[k, i_lhsT.p, i_lhsT.x] = nl.load(lhsT[k * TILE_K + i_lhsT.p,
42                                                         m * TILE_M + i_lhsT.x])
43
44     for n in nl.affine_range(N // TILE_N):
45
46         # Load a whole column tiles from rhs (with K * TILE_M numbers)
47         rhs_tiles = nl.ndarray((K // TILE_K, nl.par_dim(TILE_K), TILE_N),
48                                 dtype=rhs.dtype,
49                                 buffer=nl.sbuf)
50
51         for k in nl.affine_range(K // TILE_K):
52             rhs_tiles[k, i_rhs.p, i_rhs.x] = nl.load(rhs[k * TILE_K + i_rhs.p,
53                                                         n * TILE_N + i_rhs.x])
54
55         # Allocate a tile in PSUM for the result
56         res_psum = nl.zeros((TILE_M, TILE_N), nl.float32, buffer=nl.psum)
57         for k in nl.affine_range(K // TILE_K):
58             # Accumulate partial-sums into PSUM
59             res_psum[...] += nl.matmul(lhsT_tiles[k, i_lhsT.p, i_lhsT.x],
60                                         rhs_tiles[k, i_rhs.p, i_rhs.x],
61                                         transpose_x=True)
62
63         # Copy the result from PSUM back to SBUF, and cast to expected output data-type
64         res_sb = nl.copy(res_psum, dtype=result.dtype)
65         nl.store(result[m * TILE_M + i_res.p, n * TILE_N + i_res.x], value=res_sb)
66
67 return result

```

## Optimization 2: Reuse More Load Through Blocking

While hoisting the load out of the innermost loop eliminates some redundant loads, we can push this further by re-ordering the computation and the associated memory accesses. The technique we are going to use is called *blocking*. Blocking explicitly improves temporal locality and reduces memory accesses. It is very similar to the tiling step we did earlier in spirit.

Note that we reserve the word “tile” for defining the granularity of computation and “tiling” for the previous optimization technique that maps the high-level computation onto multiple matrix multiplication instructions executed on the TensorE. TensorE processes a specific “tile size” in a single instruction, leveraging the inherent parallelism in matrix multiplication.

Here, we do blocking, by grouping the work associated with a set of tiles together at another loop nest level. Blocking effectively interleaves a set of compute instructions and loading (DMA) instructions. This optimization does not bring us additional parallelism in computation, but rather improve the arithmetic intensity. This shifts a memory-bound matrix multiplication implementation to a compute-bound one, in order to fully leverage the compute capabilities of TensorE.

Fig. 7.72 below visualizes the memory pattern after blocking both free dimensions.

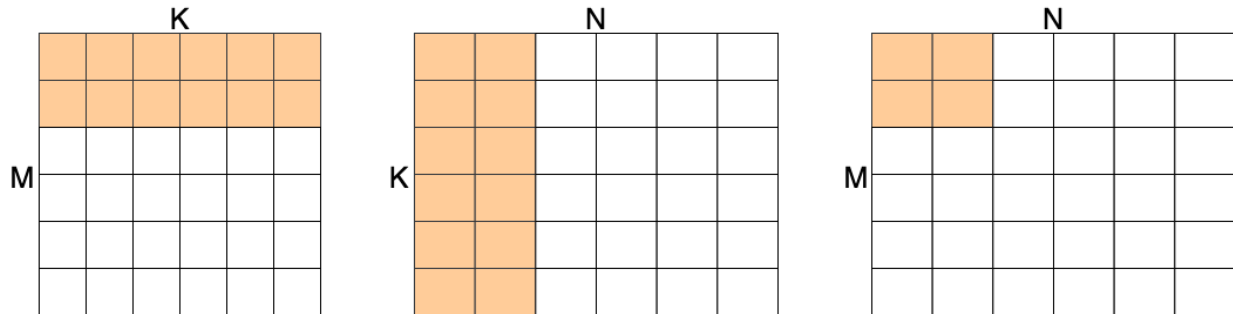


Fig. 7.72: Memory Pattern After Blocking Free Dimensions

```

1  @nki.jit
2  def nki_matmul_block_free_dimension(lhsT, rhs):
3      """NKI kernel to compute a matrix multiplication operation while blocking the
4          free dimensions of the LHS and RHS to improve memory access pattern.
5
6      Args:
7          lhsT: an input tensor of shape [K,M], where both K and M are multiples for
8              128. It is the left-hand-side argument of the matrix multiplication,
9              delivered transposed for optimal performance.
10         rhs: an input tensor of shape [K,N], where K is a multiple of 128, and N
11             is a multiple of 512. It is the right-hand-side argument of the matrix
12             multiplication.
13     Returns:
14         result: the resulting output tensor of shape [M,N]
15     """
16
17     K, M = lhsT.shape
18     K_, N = rhs.shape
19     assert K == K_, "lhsT and rhs must have the same contraction dimension"
20     result = nl.ndarray((M, N), dtype=lhsT.dtype, buffer=nl.shared_hbm)
21

```

(continues on next page)

(continued from previous page)

```

22  TILE_M = nl.tile_size.gemm_stationary_fmax # 128
23  TILE_K = nl.tile_size.pmax # 128
24  TILE_N = nl.tile_size.gemm_moving_fmax # 512
25
26  # Define the indices (shape) of the tiles
27  i_lhsT = nl.mgrid[0:TILE_K, 0:TILE_M]
28  i_rhs = nl.mgrid[0:TILE_K, 0:TILE_N]
29  i_res = nl.mgrid[0:TILE_M, 0:TILE_N]
30
31  # Configuring the blocking size for the free dimensions
32  TILES_IN_BLOCK_M = 2
33  TILES_IN_BLOCK_N = 2
34
35  BLOCK_M = TILE_M * TILES_IN_BLOCK_M # 256
36  BLOCK_N = TILE_N * TILES_IN_BLOCK_N # 1024
37
38  # the size has to be multiple of block size
39  assert M % BLOCK_M == 0
40  assert N % BLOCK_N == 0
41
42  # Loop over blocks over the M dimension
43  for m in nl.affine_range(M // BLOCK_M):
44      # Load TILES_IN_BLOCK_M columns tiles from lhsT
45      lhsT_tiles = nl.ndarray(
46          (TILES_IN_BLOCK_M, K // TILE_K, nl.par_dim(TILE_K), TILE_M),
47          dtype=lhsT.dtype,
48          buffer=nl.sbuf)
49      for bm in nl.affine_range(TILES_IN_BLOCK_M):
50          for k in nl.affine_range(K // TILE_K):
51              lhsT_tiles[bm, k, i_lhsT.p, i_lhsT.x] = nl.load(
52                  lhsT[k * TILE_K + i_lhsT.p,
53                      (m * TILES_IN_BLOCK_M + bm) * TILE_M + i_lhsT.x])
54
55      for n in nl.affine_range(N // BLOCK_N):
56          # Load TILES_IN_BLOCK_N columns from rhs
57          rhs_tiles = nl.ndarray(
58              (TILES_IN_BLOCK_N, K // TILE_K, nl.par_dim(TILE_K), TILE_N),
59              dtype=rhs.dtype,
60              buffer=nl.sbuf)
61          for bn in nl.affine_range(TILES_IN_BLOCK_N):
62              for k in nl.affine_range(K // TILE_K):
63                  rhs_tiles[bn, k, i_rhs.p, i_rhs.x] = nl.load(
64                      rhs[k * TILE_K + i_rhs.p,
65                          (n * TILES_IN_BLOCK_N + bn) * TILE_N + i_rhs.x])
66
67      for bm in nl.affine_range(TILES_IN_BLOCK_M):
68          for bn in nl.affine_range(TILES_IN_BLOCK_N):
69              # Allocate a tensor in PSUM
70              res_psum = nl.zeros((TILE_M, TILE_N), nl.float32, buffer=nl.psum)
71              for k in nl.affine_range(K // TILE_K):
72                  # Accumulate partial-sums into PSUM
73                  res_psum += nl.matmul(lhsT_tiles[bm, k, i_lhsT.p, i_lhsT.x],

```

(continues on next page)

(continued from previous page)

```

74         rhs_tiles[bn, k, i_rhs.p, i_rhs.x],
75         transpose_x=True)
76
77     # Copy the result from PSUM back to SBUF, and cast to expected output data-type
78     res_sb = nl.copy(res_psum, dtype=result.dtype)
79     nl.store(result[(m * TILES_IN_BLOCK_M + bm) * TILE_M + i_res.p,
80                   (n * TILES_IN_BLOCK_N + bn) * TILE_N + i_res.x],
81              value=res_sb)
82
83     return result

```

### Optimization 3: Further Blocking and DMA Efficiency Optimization

Next, let's also consider blocking the contraction dimension. Without blocking the contraction dimension, each block of computation leads to the final result of each output block directly, since the input blocks in both `lhs_T` and `rhs` cover the entire contraction dimension. After contraction dimension blocking, the accumulation is separated into different groups. We can accumulate the partial sum from each computation block back to an SBUF tensor for the final result. A small amount of HBM traffic might also be introduced if the partial sum cannot be kept in SBUF before being consumed. On the bright side, we can increase the block size for the free dimensions, which continues to improve the arithmetic intensity.

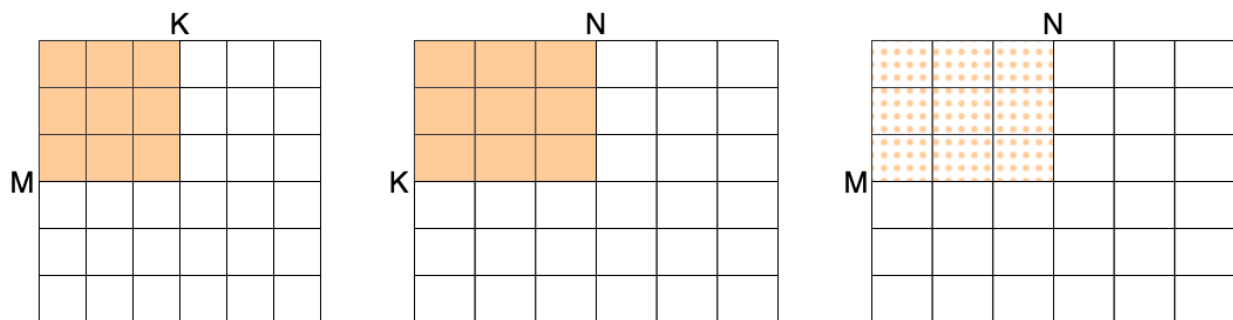


Fig. 7.73: Memory Pattern After Blocking All Dimensions

One final step we can do with NKI is to optimize the layout of the loaded tiles to improve DMA efficiency. This is done through arranging the order of dimensions in `nl.ndarray` and marking the partition dimension.

By putting all these optimizations together, we can use NKI to implement optimized matrix multiplication for different sizes. Note that different sizes of input matrices require different optimization plans. The following code optimizes for large matrix multiplication where the free dimensions of both input matrices are multiples of 2048 and the contraction dimension is a multiple of 512.

With the blocking configuration in the code (16 tiles or 2048 numbers in the `M` dimension; 2 tiles or 1024 numbers in the `N` dimension; and 8 tiles or 1024 numbers in the `K` dimension), this computation has an arithmetic intensity of 683 Flops/Byte ( $2048 \times 1024 \times 1024 / (2048 \times 1024 + 1024 \times 1024)$ ). This is certainly above the threshold of 222.

At the same time, this blocking configuration keeps all the tensors within the SBUF limit as much as possible. With all matrices in BF16 data type, the `lhs_tiles` requires 4MB and `rhs_tiles` requires 2MB SBUF memory. The `result_tiles` requires  $4 * \text{NUM\_BLOCK\_M}$  MB SBUF memory, where `NUM_BLOCK_M` is `M // 2048`. Thus, as long as `M <= 8192`, the required SBUF memory is under the 24 MB budget ( $4 + 2 + 4 * (8192 // 2048) == 22$  MB). When the `M` dimension becomes bigger, spilling and reloading of the `result_tiles` will happen, but because the frequency is relatively low, the computation can still be sufficient.

Since the K blocking loop is hand optimized for our ideal data locality, we do not actually want the compiler to rewrite this loop during its vectorization and other loop-level optimization passes. To communicate this we use `nl.sequential_range()` to construct the K blocking loop.

```

1 @nki.jit
2 def nki_matmul_fully_optimized_(
3     lhsT,
4     rhs,
5     # Meta-parameters
6     TILES_IN_BLOCK_M=16,
7     TILES_IN_BLOCK_N=2,
8     TILES_IN_BLOCK_K=8,
9 ):
10     """NKI kernel to compute a large matrix multiplication efficiently by
11         blocking all dimensions and doing layout optimization.
12
13     Args:
14         lhsT: an input tensor of shape [K,M], where K is a multiple of 128 *
15             TILES_IN_BLOCK_K and M is a multiple of 128 * TILES_IN_BLOCK_M. It is the
16             left-hand-side argument of the matrix multiplication, delivered transposed
17             for optimal performance.
18         rhs: an input tensor of shape [K,N], where K is a multiple of 128 *
19             TILES_IN_BLOCK_K and N is a multiple of 512 * TILES_IN_BLOCK_N. It is
20             the right-hand-side argument of the matrix multiplication.
21         TILES_IN_BLOCK_*: meta parameters to control blocking dimensions
22     Returns:
23         result: the resulting output tensor of shape [M,N]
24     """
25
26     K, M = lhsT.shape
27     K_, N = rhs.shape
28     assert K == K_, "lhsT and rhs must have the same contraction dimension"
29     result = nl.ndarray((M, N), dtype=lhsT.dtype, buffer=nl.shared_hbm)
30
31     TILE_M = nl.tile_size.gemm_stationary_fmax # 128
32     TILE_K = nl.tile_size.pmax # 128
33     TILE_N = nl.tile_size.gemm_moving_fmax # 512
34
35     BLOCK_M = TILE_M * TILES_IN_BLOCK_M
36     BLOCK_N = TILE_N * TILES_IN_BLOCK_N
37     BLOCK_K = TILE_K * TILES_IN_BLOCK_K
38
39     # the size has to be multiple of block size
40     assert M % BLOCK_M == 0
41     assert N % BLOCK_N == 0
42     assert K % BLOCK_K == 0
43
44     NUM_BLOCK_M = M // BLOCK_M
45     NUM_BLOCK_N = N // BLOCK_N
46     NUM_BLOCK_K = K // BLOCK_K
47
48     # Blocking N dimension (the RHS free dimension)
49     for n in nl.affine_range(NUM_BLOCK_N):

```

(continues on next page)

(continued from previous page)

```

50 result_tiles = nl.zeros((NUM_BLOCK_M, TILES_IN_BLOCK_M, TILES_IN_BLOCK_N,
51                          nl.par_dim(TILE_M), TILE_N),
52                          dtype=lhsT.dtype,
53                          buffer=nl.sbuf)
54
55 # Blocking K dimension (the contraction dimension)
56 # Use `sequential_range` because we do not want the compiler to change this loop by,
57 # for example, vectorizing it
58 for k in nl.sequential_range(NUM_BLOCK_K):
59     # Loading tiles from rhs
60     # setting the load tile to `TILE_K x BLOCK_SIZE_N` to optimize DMA performance
61     i_rhs = nl.mgrid[0:TILE_K, 0:BLOCK_N]
62     rhs_tiles = nl.ndarray((TILES_IN_BLOCK_K, nl.par_dim(TILE_K), BLOCK_N),
63                            dtype=rhs.dtype,
64                            buffer=nl.sbuf)
65
66     for bk_r in nl.affine_range(TILES_IN_BLOCK_K):
67         rhs_tiles[bk_r, i_rhs.p, i_rhs.x] = nl.load(
68             rhs[(TILES_IN_BLOCK_K * k + bk_r) * TILE_K + i_rhs.p,
69                 BLOCK_N * n + i_rhs.x])
70
71     # Blocking M dimension (the LHS free dimension)
72     for m in nl.affine_range(NUM_BLOCK_M):
73         # Loading tiles from lhsT
74         i_lhsT = nl.mgrid[0:TILE_K, 0:BLOCK_M]
75         lhsT_tiles = nl.ndarray((TILES_IN_BLOCK_K, nl.par_dim(TILE_K), BLOCK_M),
76                                 dtype=lhsT.dtype,
77                                 buffer=nl.sbuf)
78         for bk_l in nl.affine_range(TILES_IN_BLOCK_K):
79             lhsT_tiles[bk_l, i_lhsT.p, i_lhsT.x] = nl.load(
80                 lhsT[(TILES_IN_BLOCK_K * k + bk_l) * TILE_K + i_lhsT.p,
81                     BLOCK_M * m + i_lhsT.x])
82
83         # Do matmul with all tiles in the blocks
84         i_lhsT_mm = nl.mgrid[0:TILE_K, 0:TILE_M]
85         i_rhs_mm = nl.mgrid[0:TILE_K, 0:TILE_N]
86         i_res_mm = nl.mgrid[0:TILE_M, 0:TILE_N]
87         for bn in nl.affine_range(TILES_IN_BLOCK_N):
88             for bm in nl.affine_range(TILES_IN_BLOCK_M):
89                 res_tile = nl.zeros((TILE_M, TILE_N), dtype=nl.float32, buffer=nl.psum)
90
91                 for bk in nl.affine_range(TILES_IN_BLOCK_K):
92                     res_tile[...] += nisa.nc_matmul(
93                         lhsT_tiles[bk, i_lhsT_mm.p, bm * TILE_M + i_lhsT_mm.x],
94                         rhs_tiles[bk, i_rhs_mm.p, bn * TILE_N + i_rhs_mm.x])
95
96                 # Accumulate on corresponding SBUF tile
97                 result_tiles[m, bm, bn, i_res_mm.p,
98                             i_res_mm.x] += res_tile[i_res_mm.p, i_res_mm.x]
99
100     # Copying the result from SBUF to HBM
101     for m in nl.affine_range(NUM_BLOCK_M):

```

(continues on next page)

(continued from previous page)

```

102     for bm in nl.affine_range(TILES_IN_BLOCK_M):
103         i_res = nl.mgrid[0:TILE_K, 0:TILE_N]
104         i_res_packed = nl.mgrid[0:TILE_K, 0:BLOCK_N]
105         result_packed = nl.ndarray((TILE_K, BLOCK_N),
106                                     dtype=result_tiles.dtype,
107                                     buffer=nl.sbuf)
108
109         # coalesce result tiles for better DMA performance
110         for bn in nl.affine_range(TILES_IN_BLOCK_N):
111             result_packed[i_res.p,
112                           bn * TILE_N + i_res.x] = nl.copy(result_tiles[m, bm, bn,
113                                                                    i_res.p,
114                                                                    i_res.x])
115         nl.store(result[(TILES_IN_BLOCK_M * m + bm) * TILE_K + i_res_packed.p,
116                        BLOCK_N * n + i_res_packed.x],
117                 value=result_packed[i_res_packed.p, i_res_packed.x])
118
119     return result

```

## Testing Correctness and Benchmarking

To test the correctness of the kernels, we compare the result with the `torch.matmul` with `torch.allclose`.

```

1  # Test the large workload with tiled kernels
2  lhs = torch.rand((4096, 1024), dtype=torch.bfloat16, device=device)
3  rhs = torch.rand((1024, 2048), dtype=torch.bfloat16, device=device)
4
5  # Run torch reference
6  output_torch = torch.matmul(lhs, rhs).to(device=cpu)
7
8  def check_match(nki_func):
9      output = nki_func(lhs.T, rhs)
10     output_nki = output.to(device=cpu)
11     if torch.allclose(output_torch, output_nki, atol=1e-4, rtol=1e-2):
12         print("NKI and Torch match")
13     else:
14         print("NKI and Torch differ")
15
16  print("Checking correctness of nki_matmul_tiled")
17  check_match(nki_matmul_tiled_)
18
19  print("Checking correctness of nki_matmul_hoist_load")
20  check_match(nki_matmul_hoist_load_)
21
22  print("Checking correctness of nki_matmul_block_free_dimension")
23  check_match(nki_matmul_block_free_dimension_)
24
25  print("Checking correctness of nki_matmul_fully_optimized")
26  check_match(nki_matmul_fully_optimized_)

```

Output from the test:



```

Checking correctness of nki_matmul_tiled
NKI and Torch match
Checking correctness of nki_matmul_hoist_load
NKI and Torch match
Checking correctness of nki_matmul_block_free_dimension
NKI and Torch match
Checking correctness of nki_matmul_fully_optimized
NKI and Torch match

```

To test for performance of each kernel here, we can use NKI's benchmark capability to measure the performance of the four different kernels on [4096,8192] @ [8192,8192] matrix multiplication.

```

1 if __name__ == "__main__":
2     # Benchmarking with large matrices to show the differences more clearly
3     lhsT = nt.tensor([8192, 4096], nl.bfloat16)
4     rhs = nt.tensor([8192, 8192], nl.bfloat16)
5
6     def benchmark_nki(nki_func):
7         bench_func = nki.benchmark(warmup=5, iters=10)(nki_func)
8         bench_func(lhsT, rhs)
9         latency_res = bench_func.benchmark_result.nc_latency
10        p99 = latency_res.get_latency_percentile(99)
11        print("Latency: {:.2f} ms (P99)".format(p99 / 1000.0))
12
13    print("Benchmarking nki_matmul_tiled")
14    benchmark_nki(nki_matmul_tiled_)
15
16    print("Benchmarking nki_matmul_hoist_load")
17    benchmark_nki(nki_matmul_hoist_load_)
18
19    print("Benchmarking nki_matmul_block_free_dimension")
20    benchmark_nki(nki_matmul_block_free_dimension_)
21
22    print("Benchmarking nki_matmul_fully_optimized")
23    benchmark_nki(nki_matmul_fully_optimized_)

```

| Kernels        | Latency (ms) | Hardware FLOPs Utilization (HFU, %) |
|----------------|--------------|-------------------------------------|
| Original Tiled | 51.80        | 10.98                               |
| Optimization 1 | 42.96        | 13.24                               |
| Optimization 2 | 22.07        | 26.51                               |
| Optimization 3 | 6.97         | 85.24                               |

As shown in the table above, with all the optimizations, the matrix multiplication kernel is 7x faster comparing to the original tiled version. We also profile the four different kernel implementations for the HFU (hardware FLOPs utilization). With all the optimizations, the final version reaches a HFU of 85.2%. The performance numbers here are specific to input matrix sizes ([4096,8192] @ [8192,8192]), data types (BF16), and server instance (Trn1.32xlarge).

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- All matrix multiplication NKI kernels: `matrix_multiplication_nki_kernels.py`
- PyTorch implementation: `matrix_multiplication_torch.py`

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

Run benchmarking of different NKI kernels:

```
python3 matrix_multiplication_nki_kernels.py
```

Run PyTorch implementation to validate the NKI results against the PyTorch implementation:

```
python3 matrix_multiplication_torch.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## RMSNorm

In this tutorial, we implement a kernel to perform RMSNorm of a **2D tensor**, as described in [Root Mean Square Layer Normalization](#). In doing so, we learn about:

- The NKI syntax and programming model
- Broadcasting tensors in different axis
- Mapping embarrassingly parallel vector operations efficiently to the NeuronCore
- Disable ineffectual data movement or compute within a tile using an execution mask

Before diving into RMSNorm of 2D input, let's go over the RMSNorm operator for a **1D vector**  $\mathbf{a}$  defined as below:

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\mathbf{a})} g_i, \text{ where } \text{RMS}(\mathbf{a}) = \sqrt{\frac{1}{n} \sum_{i=0}^n a_i^2}$$

Note,  $\mathbf{g}$  is the RMSNorm weight, which has the same shape as the input vector  $\mathbf{a}$ . The function  $\text{RMS}(\mathbf{a})$  produces a single scalar element, and we divide every element in the input vector  $\mathbf{a}$  by the  $\text{RMS}(\mathbf{a})$  scalar (i.e., a broadcast divide).

In Transformer models, we typically perform RMSNorm on a 2D input tensor instead (with shape `[sequence length, embedding size]`). 2D-RMSNorm simply performs 1D-RMSNorm as discussed above for *every row* of the input 2D tensor. The  $\mathbf{g}$  RMSNorm weight vector is shared (i.e., broadcasted) across the rows for the multiplication. [Figure](#) below visualizes the tensor shapes involved in 2D-RMSNorm, where  $\mathbf{a\_tensor}$  is the 2D input tensor and  $\mathbf{g\_tensor}$  is the 1D RMSNorm weight:

We are going to map the rows ( `$\mathbf{a\_tensor}.\text{shape}[0]$` ) to the partition dimension of the SBUF once we load the tensor from HBM. This is a natural layout choice since each SBUF partition has a one-to-one mapping to a parallel vector lane in the compute engines for calculating  $\text{RMS}(\mathbf{a\_tensor})$ .

Note, the division of  $\text{RMS}(\mathbf{a\_tensor})$  requires broadcasting of one scalar across all elements of  $\mathbf{a\_tensor}$  within each partition, which is considered a free-axis broadcast and supported by the flexible memory access pattern in hardware. On the other hand, the multiplication with  $\mathbf{g\_tensor}$  requires broadcasting of a vector across all partitions, which is

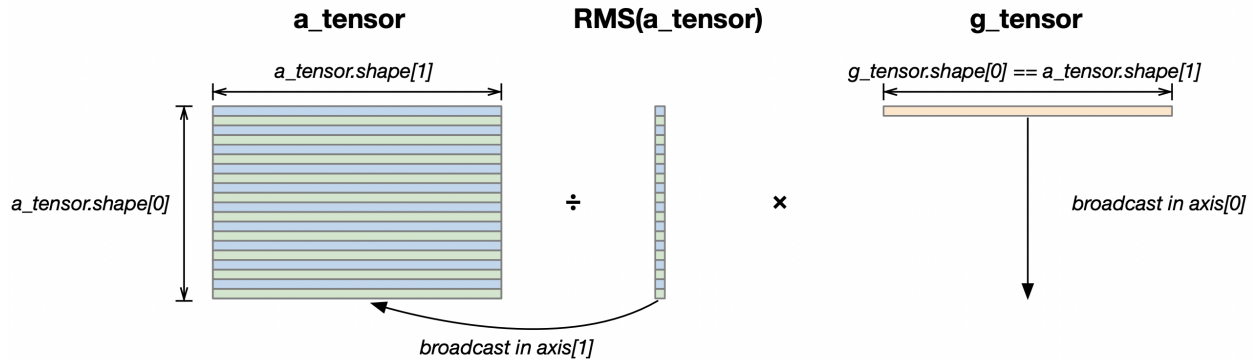


Fig. 7.74: RMSNorm tensor shapes

considered a partition-axis broadcast and must invoke another instruction for the broadcasting (`broadcast_to()` API, details see below implementation) .

### Compute kernel

```

1  import math
2  import neuronxcc.nki as nki
3  import neuronxcc.nki.language as nl
4
5
6  @nki.jit
7  def nki_rmsnorm_kernel(a_tensor, g_tensor):
8      # Calculate out_tensor = a_tensor/RMS(a_tensor) * g_tensor
9      # Where RMS(a_tensor) = sqrt((1/N) * sum(a_tensor * a_tensor))
10     # and N = a_tensor.shape[1]
11     # Reduction (mean) is performed in the free (2nd) dimension
12     out_tensor = nl.ndarray(a_tensor.shape, dtype=a_tensor.dtype,
13                             buffer=nl.shared_hbm)
14
15     # Make sure shapes match
16     assert a_tensor.shape[1] == g_tensor.shape[0]
17
18     # Generate tensor indices to index input tensor
19     ix = nl.arange(128)[: , None]
20     iw = nl.arange(1)[: , None]
21     iy = nl.arange(a_tensor.shape[1])[None, :]
22
23     num_rows = a_tensor.shape[0]
24
25     # Load RMSNorm weight once, reused by rows/tiles of a_tensor
26     g_tile = nl.load(g_tensor.reshape((1, g_tensor.shape[0]))[iw, iy])
27
28     # Process 128 rows at a time due to 128-partition tile size limitation
29     # Since we're not reducing across the first dimension
30     # Tiles can be processed independently
31     for i in nl.affine_range(math.ceil(a_tensor.shape[0]/128)):

```

(continues on next page)

(continued from previous page)

```

32
33 # Load input data from external memory to on-chip memory
34 a_tile = nl.load(a_tensor[i * 128 + ix, iy],
35                 mask=(i * 128 + ix < num_rows))
36
37 # Compute element-wise square of a_tensor
38 in_square = nl.square(a_tile)
39
40 # Calculate sum of squared elements, along last dimension
41 square_sum = nl.sum(in_square, axis=[1])
42
43 # Scale and get a reciprocal
44 mean = square_sum / a_tensor.shape[1]
45
46 # Take square root of mean and then reciprocal with
47 # rsqrt API (one ISA instruction)
48 rms_reciprocal = nl.rsqrt(mean)
49
50 # Scale the input tensor
51 out_tile = nl.multiply(a_tile, rms_reciprocal)
52
53 # Broadcast weight along first axis to match tensor shape
54 # num_rows_active = min(num_rows - i * 128, 128)
55 g_bcast = g_tile.broadcast_to((128, g_tensor.shape[0]))
56
57 # Multiply with the RMSNorm weight
58 out_tile[...] = nl.multiply(out_tile, g_bcast,
59                             mask=(i * 128 + ix < num_rows))
60
61 # store the addition results back to external memory (out_tensor)
62 nl.store(out_tensor[i * 128 + ix, iy], value=out_tile,
63         mask=(i * 128 + ix < num_rows))
64
65 return out_tensor

```

In this example, we implement RMSNorm for a 2D input tensor in `nki_rmsnorm_kernel`:

- We assume each SBUF partition is large enough to fit at least one row of `a_tensor` and one copy of `g_tensor` simultaneously.
- We load `g_tensor` once into the SBUF outside the main loop that iterates over tiles of `a_tensor` to achieve maximum reuse. The `g_tensor` is reshaped into a 2D tensor because SBUF is a two-dimensional memory and hence expects at least two dimension for any SBUF tensor. A reshape of an HBM tensor without changing the underlying storage format is in fact a no-op with no performance cost in the final compiled executable.
- To adhere to NKI's tile-size considerations (*Tile Size Considerations*), we limit the partition axis size of `g_tensor` tile to be 128.
- The trip count of the compute loop is `math.ceil(a_tensor.shape[0]/128)`. In cases where `a_tensor.shape[0]` is not a multiple of 128, we can disable ineffectual data movement or compute in the last iteration using the mask field (discussions below).
- Within the compute loop:
  - We load one tile of `g_tensor` with shape `(128, g_tensor.shape[1])` using `nl.load` API. We guard

the loading boundary by specifying `mask=(i * 128 + ix < num_rows)`, which ensures we don't access out-of-bound memory when the number of rows in `a_tensor` is not a multiple of 128.

- We perform the free-axis broadcast multiply (division of `RMS(a)`) using `nl.multiply(a_tile, rms_reciprocal)`, which is lowered into `nki.isa.tensor_scalar` instruction under the hood.
- To broadcast multiply with the RMSNorm weight `g_tensor`, we need to perform a partition-axis broadcast of the `g_tensor`. The number of partitions to broadcast to depends on how many active rows are being normalized in the current loop iteration: `min(num_rows - i * 128, 128)`. Next, we can do element-wise multiplication of the broadcasted `g_tensor` and the intermediate normalized tile `out_tile`, which is lowered into `nki.isa.tensor_tensor` instruction under the hood.
- Finally, we store the normalized tile back into HBM using the `nl.store` API. We guard the store boundary similar to load boundary using the mask field.

## Launching kernel and testing correctness

### PyTorch

Below we write a reference PyTorch implementation of RMSNorm and verify our NKI kernel output against the reference in the same script as the kernel.

```

1  # Reference torch implementation
2  def torch_rmsnorm_kernel(a_tensor, g_tensor):
3      # Square the tensor (element-wise)
4      in_square = a_tensor.pow(2)
5      # Calculate means in the free dimension
6      mean = in_square.mean(dim=1, keepdim=True)
7      # Scale by reciprocal of sqrt(mean)
8      tensor = a_tensor * torch.rsqrt(mean)
9
10     # Scale the output by the weight
11     return tensor * g_tensor
12
13 from torch_xla.core import xla_model as xm
14 device = xm.xla_device()
15
16 a_tensor = torch.rand((250, 512), dtype=torch.float32).to(device=device)
17 g_tensor = torch.rand((512), dtype=torch.float32).to(device=device)
18
19 output_nki = nki_rmsnorm_kernel(a_tensor, g_tensor)
20 print(f"output_nki={output_nki}")
21
22 output_torch = torch_rmsnorm_kernel(a_tensor, g_tensor)
23 print(f"output_torch={output_torch}")
24
25 if torch.allclose(output_torch, output_nki, atol=1e-5, rtol=1e-3):
26     print("NKI and Torch match")
27 else:
28     print("NKI and Torch differ")

```

Output:

```

2024-07-27 15:22:50.000670: 7592 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2024-07-27 15:22:50.000672: 7592 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd:
↳neuronx-cc compile --target=trn1 --framework=XLA /tmp/ubuntu/neuroncc_compile_workdir/
↳54c8e689-108c-433e-832a-f9282acdf114/model.MODULE_7170924315921358669+d41d8cd9.hlo_
↳module.pb --output /tmp/ubuntu/neuroncc_compile_workdir/54c8e689-108c-433e-832a-
↳f9282acdf114/model.MODULE_7170924315921358669+d41d8cd9.neff --verbose=35
DGE ON Levels: {'scalar_dynamic_offset', 'io'}

.
Compiler status PASS
output_nki=tensor([[0.8418, 1.3092, 0.7372, ..., 0.1458, 0.8831, 0.2339],
[0.1745, 0.3416, 0.1519, ..., 0.3358, 0.1832, 0.4795],
[0.0111, 1.1799, 0.8628, ..., 0.3107, 0.8328, 0.5663],
...,
[1.1213, 0.5449, 0.3020, ..., 0.4050, 0.4838, 0.0834],
[0.8246, 0.5027, 0.2745, ..., 0.4069, 1.0456, 1.0978],
[0.6415, 0.3637, 0.1462, ..., 0.2441, 1.0535, 0.4138]],
device='xla:0')
2024-07-27 15:22:51.000907: 7592 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2024-07-27 15:22:51.000908: 7592 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd:
↳neuronx-cc compile --target=trn1 --framework=XLA /tmp/ubuntu/neuroncc_compile_workdir/
↳6d2046fc-c02d-4d3d-8746-50399ad50832/model.MODULE_18272098496972694952+d41d8cd9.hlo_
↳module.pb --output /tmp/ubuntu/neuroncc_compile_workdir/6d2046fc-c02d-4d3d-8746-
↳50399ad50832/model.MODULE_18272098496972694952+d41d8cd9.neff --verbose=35
DGE ON Levels: {'scalar_dynamic_offset', 'io'}

.
Compiler status PASS
output_torch=tensor([[0.8418, 1.3092, 0.7372, ..., 0.1458, 0.8831, 0.2339],
[0.1745, 0.3416, 0.1519, ..., 0.3358, 0.1832, 0.4795],
[0.0111, 1.1799, 0.8628, ..., 0.3107, 0.8328, 0.5663],
...,
[1.1213, 0.5449, 0.3020, ..., 0.4050, 0.4838, 0.0834],
[0.8246, 0.5027, 0.2745, ..., 0.4069, 1.0456, 1.0978],
[0.6415, 0.3637, 0.1462, ..., 0.2441, 1.0535, 0.4138]],
device='xla:0')
2024-07-27 15:22:53.000466: 7592 INFO ||NEURON_CACHE||: Compile cache path: /var/tmp/
↳neuron-compile-cache
2024-07-27 15:22:53.000467: 7592 INFO ||NEURON_CC_WRAPPER||: Call compiler with cmd:
↳neuronx-cc compile --target=trn1 --framework=XLA /tmp/ubuntu/neuroncc_compile_workdir/
↳32c983cd-2c40-4723-8342-d4422107708c/model.MODULE_968738949480579147+d41d8cd9.hlo_
↳module.pb --output /tmp/ubuntu/neuroncc_compile_workdir/32c983cd-2c40-4723-8342-
↳d4422107708c/model.MODULE_968738949480579147+d41d8cd9.neff --verbose=35
DGE ON Levels: {'io', 'scalar_dynamic_offset'}

.
Compiler status PASS
NKI and Torch match

```

## JAX

Below we write a reference JAX implementation of RMSNorm and verify our NKI kernel output against the reference in the same script as the kernel.

```

1  # Reference JAX implementation
2  def jax_rms_norm(a_tensor, g_tensor):
3      # Square the tensor (element-wise)
4      in_square = jnp.square(a_tensor)
5      # Calculate means in the free dimension
6      mean = in_square.mean(axis=1, keepdims=True)
7      # Scale by reciprocal of sqrt(mean)
8      tensor = a_tensor * jnp.reciprocal(jnp.sqrt(mean))
9
10     # Scale the output by the weight
11     return tensor * g_tensor
12
13     a_key, g_key = jax.random.split(jax.random.PRNGKey(42))
14     a_tensor = jax.random.uniform(a_key, (250, 512))
15     g_tensor = jax.random.uniform(g_key, (512,))
16
17     output_nki = nki_rmsnorm_kernel(a_tensor, g_tensor)
18
19     print(a_tensor)
20
21     print(f"output_nki={output_nki}")
22
23     output_jax = jax_rms_norm(a_tensor, g_tensor)
24     print(f"output_jax={output_jax}")
25
26     if jnp.allclose(output_jax, output_nki, atol=1e-5, rtol=1e-3):
27         print("NKI and JAX match")
28     else:
29         print("NKI and JAX differ")

```

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- NKI baremetal implementation: `rmsnorm_nki_kernels.py`
- PyTorch reference implementation: `rmsnorm_torch.py`
- JAX reference implementation: `rmsnorm_jax.py`

You can also view the source code in the GitHub repository [nki\\_samples](#)

### Example usage of the scripts:

Run NKI baremetal implementation:

```
python3 rmsnorm_nki_kernels.py
```

Run PyTorch implementation:

```
python3 rmsnorm_torch.py
```

Run JAX implementation:

```
python3 rmsnorm_jax.py
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## LayerNorm

In this tutorial, we implement a kernel to perform LayerNorm of a **2D tensor**, as described in [Layer Normalization](#). LayerNorm is a common normalization mechanism used in Transformer models, similar to [RMSNorm](#). However, LayerNorm requires more vector operations to optimize compute efficiency in Vector Engine. In doing so, we will revisit the key concepts we learned in the [RMSNorm](#) and additionally learn about:

- Using nki.isa APIs to efficiently compute mean and variance, and minimize the number of traversals over input data by combining multiple vector instructions into one
- Take surrounding compute into consideration when deciding tensor layouts

Before diving into LayerNorm for a 2D tensor, let's go over the LayerNorm operator for a **1D vector**  $y$  defined as below:

$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{var}[x] + \epsilon}} * \gamma + \beta$$

The parameters are:

- $x$ : Input 1D vector
- $y$ : Output 1D vector, same shape as  $x$
- $\mathbb{E}[x]$ : Mean of  $x$
- $\text{var}[x]$ : Variance of  $x$
- $\epsilon$ : A small constant scalar for numerical stability
- $\gamma, \beta$ : LayerNorm affine transform parameters, each has the same shape as  $x$

In Transformer models, we typically need to perform LayerNorm on a 2D input tensor (with shape: `[sequence_length, hidden_size]`), where the first dimension is `sequence_length` long corresponding to the number of tokens currently being processed, and the second dimension is the embedding dimension of each token.

Different tokens (i.e., rows in the `[sequence_length, hidden_size]` 2D vector) undergo different 1D LayerNorm independently. Therefore, we need to calculate different mean and variance for different rows and broadcast (i.e., share) the same  $\gamma, \beta$  parameters across the rows.

Figure below visualizes the tensor shape involved in 2D-LayerNorm, where `input_tensor` is 2D input vector and `gamma_vector` and `beta_vector` are affine transform parameters:



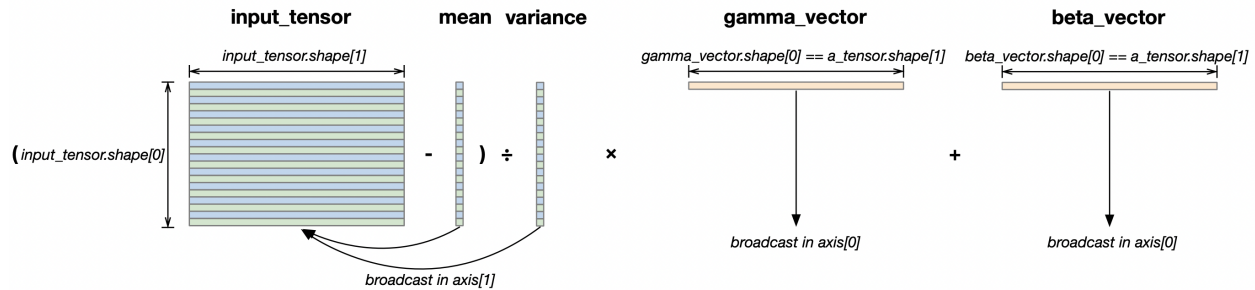


Fig. 7.75: LayerNorm tensor shapes

Compared to *RMSNorm*, LayerNorm requires calculations of mean and variance, instead of simple square and summation. Also, LayerNorm performs two instances of free-axis broadcast and two instances of partition-axis broadcast, while *RMSNorm* requires one instance of each. Therefore, LayerNorm involves way more computation (vector operations in particular) than *RMSNorm*.

## Implement NKI kernel

Next, we will present two versions of LayerNorm implementation, starting from a naive version using `nki.language` APIs and ending with an optimized version using `nki.isa` APIs.

### Version 1: `nki.language` APIs only

```

1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3 import neuronxcc.nki.isa as nisa
4 import numpy as np
5 import math
6
7 @nki.jit
8 def nki_layernorm_kernel_v1(input_tensor, epsilon, gamma_vector, beta_vector):
9     """Computes LayerNorm.
10     Used nki.language APIs only.
11     """
12     output_tensor = nl.ndarray(input_tensor.shape, dtype=input_tensor.dtype,
13                               buffer=nl.shared_hbm)
14
15     # Ensure that the shapes of tensors match
16     assert input_tensor.shape[1] == gamma_vector.shape[0] == beta_vector.shape[0]
17
18     # Generate tile indices for loading/storing data
19     i_p_io = nl.arange(nl.tile_size.pmax)[: , None]
20     i_f_io = nl.arange(input_tensor.shape[1])[None, :]
21     i_p_param = nl.arange(1)[: , None]
22
23     # Number of rows in the input tensor
24     num_rows = input_tensor.shape[0]
25
26     # Load gamma and beta, which will be reused across rows/tiles of input_tensor

```

(continues on next page)

(continued from previous page)

```

27 gamma_sb = nl.load(gamma_vector.reshape((1, gamma_vector.shape[0]))[i_p_param, i_f_io])
28 beta_sb = nl.load(beta_vector.reshape((1, beta_vector.shape[0]))[i_p_param, i_f_io])
29
30 # Broadcast the gamma and beta to match the dimensions of the tiles
31 gamma_sb_bcast = gamma_sb.broadcast_to((nl.tile_size.pmax, gamma_vector.shape[0]))
32 beta_sb_bcast = beta_sb.broadcast_to((nl.tile_size.pmax, beta_vector.shape[0]))
33
34 # Tile partition dimension of the input tensor by nl.tile_size.pmax
35 for i in nl.affine_range(math.ceil(input_tensor.shape[0]/nl.tile_size.pmax)):
36     # Load input tile
37     input_sb = nl.load(input_tensor[i * nl.tile_size.pmax + i_p_io, i_f_io],
38                         mask=(i * nl.tile_size.pmax + i_p_io < num_rows))
39
40     # Compute mean and variance
41     mean = nl.mean(input_sb, axis=1)
42     # Trick to calculate var with mean: mean(x^2) - mean(x)^2
43     var = nl.mean(nl.square(input_sb), axis=1) - mean * mean
44
45     # Normalize the input by shifting with the mean
46     # and scaling with rsqrt of variance and epsilon
47     shift_scale_tensor = (input_sb - mean) * nl.rsqrt(var + epsilon)
48
49     # Scale the normalized tile using gamma and add beta
50     output_sb = shift_scale_tensor * gamma_sb_bcast + beta_sb_bcast
51
52     nl.store(output_tensor[i * nl.tile_size.pmax + i_p_io, i_f_io], value=output_sb,
53             mask=(i * nl.tile_size.pmax + i_p_io < num_rows))
54
55 return output_tensor

```

- To adhere to NKI's tile-size considerations (*Tile Size Considerations*), we limit the partition axis size of `input_tensor` tile to be 128 (`nl.tile_size.pmax`).
- **Load gamma and beta, and perform the partition-axis broadcast:**
  - The multiplication with `shift_scale_tensor` requires broadcasting of gamma and beta across all partitions(`broadcast_to()` API)
- The trip count of the compute loop is `math.ceil(input_tensor.shape[0]/nl.tile_size.pmax)`. In cases where `input_tensor.shape[0]` is not a multiple of `nl.tile_size.pmax`, we can disable ineffectual data movement or compute in the last iteration using the mask field.
- **Within the compute loop:**
  - We load one tile of `input_tensor` with shape `(nl.tile_size.pmax, input_tensor.shape[1])` using `nl.load` API. We guard the loading boundary by specifying `mask=(i * nl.tile_size.pmax + i_p_io < input_tensor.shape[0])`, which ensures we don't access out-of-bound memory when the number of rows in `input_tensor` is not a multiple of `nl.tile_size.pmax`.
  - Compute the mean and variance using *nki.language.mean*
  - Normalize one tile of `input_tensor` using mean and variance. The variance is preprocessed using *nki.language.rsqrt*
  - Scale the normalized tile using *gamma* and add *beta*

- Finally, we store the normalized tile back into HBM using the `nl.store` API. We guard the store boundary similar to load boundary using the mask field.

Next, we will optimize the above implementation using `nki.isa` APIs in version 2

## Version 2: `nki.isa` APIs to calculate mean/variance and perform shift/scale

```

1 @nki.jit
2 def nki_layernorm_kernel_v2(input_tensor, epsilon, gamma_vector, beta_vector):
3     """Computes LayerNorm.
4     Used nki.isa APIs to calculate mean/variance and perform shift/scale.
5     """
6     output_tensor = nl.ndarray(input_tensor.shape, dtype=input_tensor.dtype,
7                                buffer=nl.shared_hbm)
8
9     # Ensure that the shapes of tensors match
10    assert input_tensor.shape[1] == gamma_vector.shape[0] == beta_vector.shape[0]
11
12    # Generate tile indices for loading/storing data
13    i_p_io = nl.arange(nl.tile_size.pmax)[: , None]
14    i_f_io = nl.arange(input_tensor.shape[1])[None, :]
15    i_p_param = nl.arange(1)[: , None]
16
17    # Number of rows in the input tensor
18    num_rows = input_tensor.shape[0]
19
20    # Load gamma and beta, which will be reused across rows/tiles of input_tensor
21    gamma_sb = nl.load(gamma_vector.reshape((1, gamma_vector.shape[0]))[i_p_param, i_f_io])
22    beta_sb = nl.load(beta_vector.reshape((1, beta_vector.shape[0]))[i_p_param, i_f_io])
23
24    # Broadcast the gamma and beta to match the dimensions of the tiles
25    gamma_sb_bcast = gamma_sb.broadcast_to((nl.tile_size.pmax, gamma_vector.shape[0]))
26    beta_sb_bcast = beta_sb.broadcast_to((nl.tile_size.pmax, beta_vector.shape[0]))
27
28    # Tile partition dimension of the input tensor by nl.tile_size.pmax
29    for i in nl.affine_range(math.ceil(input_tensor.shape[0]/nl.tile_size.pmax)):
30        # Load input tile
31        input_sb = nl.load(input_tensor[i * nl.tile_size.pmax + i_p_io, i_f_io],
32                             mask=(i * nl.tile_size.pmax + i_p_io < num_rows))
33
34        # Tile free dimension of the input tensor by nl.tile_size.bn_stats_fmax,
35        # as bn_stats has a free dimension size limit
36        i_f_bn = nl.arange(nl.tile_size.bn_stats_fmax)[None, :]
37        i_f_stats = nl.arange(6)[None, :]
38        num_bn_stats = math.ceil(input_tensor.shape[1]/nl.tile_size.bn_stats_fmax)
39        stats_results = nl.ndarray((nl.tile_size.pmax, 6*num_bn_stats), dtype=np.float32)
40        for j in nl.affine_range(num_bn_stats):
41            stats_results[i_p_io, j * 6 + i_f_stats] = nisa.bn_stats(
42                input_sb[i_p_io, j * nl.tile_size.bn_stats_fmax + i_f_bn],
43                mask=(j * nl.tile_size.bn_stats_fmax + i_f_bn < input_tensor.shape[1]),
44                dtype=np.float32)
45

```

(continues on next page)

(continued from previous page)

```

46 # Aggregate bn_stats results to compute mean and var
47 i_f_aggr = nl.arange(6*num_bn_stats)[None, :]
48 mean_var = nisa.bn_aggr(stats_results[i_p_io, i_f_aggr])
49 mean = mean_var[i_p_io, 0]
50 var = mean_var[i_p_io, 1]
51
52 # Get reciprocal of sqrt(var + epsilon)
53 scale_var = nl.rsqrt(var + epsilon)
54
55 # Putting the shift and scale together in one line to trigger two alu_op tensor_
↪vector instruction
56 # shift_scale_tensor = (input_sb - mean_var[i_p_stats, i_f_mean]) * scale_var
57 shift_scale_tensor = nisa.tensor_scalar(data=input_sb, op0=np.subtract,
58                                         operand0=mean,
59                                         op1=np.multiply,
60                                         operand1=scale_var)
61
62 # Scale the normalized tile using gamma and add beta
63 output_sb = shift_scale_tensor * gamma_sb_bcast + beta_sb_bcast
64
65 nl.store(output_tensor[i * nl.tile_size.pmax + i_p_io, i_f_io], value=output_sb,
66          mask=(i * nl.tile_size.pmax + i_p_io < num_rows))
67
68 return output_tensor

```

- Considering the free dimension size limit of *nki.isa.bn\_stats*, which is 512(*nl.tile\_size.bn\_stats\_fmax*), the trip count of *bn\_stats* compute loop is  $\text{math.ceil}(\text{input\_tensor.shape}[1]/\text{nl.tile\_size.bn\_stats\_fmax})$ .
- Used *nki.isa.bn\_stats* and *nki.isa.bn\_aggr* to calculate the mean and variance
- Used *nki.isa.tensor\_scalar* to do shift and scale of mean and variance in a single instruction

## Performance in Version 1 and Version 2

Let's assume the data type for the kernel is float32 and that the SBUF partition is sufficiently large to hold the intermediate data simultaneously without significant spilling. Define the variable  $N = \text{input\_tensor.shape}[1]$ .

- **Compute mean and variance:**
  - Version 1 : The performance cost of the mean calculation is  $N$  Vector Engine cycles, and the variance calculation is  $N$  Scalar Engine +  $2N$  Vector Engine cycles.
  - Version 2 : By replacing these calculations with *bn\_stats* and *bn\_aggr* APIs, the cost is roughly reduced to  $N$  Vector Engine cycles, ignoring the cost of *nki.isa.bn\_aggr*, assuming  $N$  is sufficiently large.
- **Perform shift and scale of mean and variance in a single instruction:**
  - Version 1 : The performance cost of the shift/scale calculation requires two small instructions (*nl.rsqrt(var + epsilon)*) and two instructions with each iterating over  $N$  elements per partition (shift and scale,  $2N$ ).
  - Version 2 : By replacing these calculations with the *tensor\_scalar* API, the cost is reduced to  $N$  Vector Engine cycles

The latency measured on trn1 using an input tensor of (300, 1000) shows a 14.9% improvement.

```
>>>> Running version v1.
```

```
Latency results are:
```

```
NCLatency:
```

```
p0 = 2306us
```

```
p1 = 2306us
```

```
p10 = 2308us
```

```
p25 = 2309us
```

```
p50 = 2311us
```

```
p90 = 2313us
```

```
p99 = 2314us
```

```
p100 = 2314us
```

```
>>>> Running version v2.
```

```
Latency results are:
```

```
NCLatency:
```

```
p0 = 1963us
```

```
p1 = 1963us
```

```
p10 = 1965us
```

```
p25 = 1966us
```

```
p50 = 1969us
```

```
p90 = 1972us
```

```
p99 = 1974us
```

```
p100 = 1975us
```

## Launching kernel and testing correctness

Below is a reference PyTorch implementation of LayerNorm, which we use to verify our NKI kernel output against the reference output

```
1 import torch
2 from torch_xla.core import xla_model as xm
3 import argparse
4 import os
5
6 os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
7
8 # Reference torch implementation
9 def layernorm_layer(input_tensor, epsilon, gamma_vector, beta_vector):
10     # Compute the mean and variance of the input tensor along the last dimension
11     mean = input_tensor.mean(dim=-1, keepdim=True)
12     variance = input_tensor.var(dim=-1, keepdim=True, unbiased=False)
13     # Subtract the mean from the input and divide by the square root of the variance,
14     # plus epsilon
15     normalized_input = (input_tensor - mean) / torch.sqrt(variance + epsilon)
16     # Apply the affine transformation
17     normalized_input = normalized_input * gamma_vector + beta_vector
18     return normalized_input
19
20 def parse_args():
21     parser = argparse.ArgumentParser(
```

(continues on next page)

(continued from previous page)

```

21     """Run LayerNorm pytorch implementation.
22     """
23     parser.add_argument("--nrows",
24                         default=4*1024,
25                         type=int,
26                         help="""The number of input rows""")
27     parser.add_argument("--ncols",
28                         default=8*1024,
29                         type=int,
30                         help="""The number of input columns""")
31     parser.add_argument("--version",
32                         default="v1",
33                         choices=["v1", "v2"],
34                         help="Test versions")
35     args = parser.parse_args()
36     return args
37
38
39 from neuronxcc.nki.docs.examples.layer_norm.layer_norm_nki_kernel import nki_layer_norm_
40 ↪kernel_v1, \
41     nki_layer_norm_kernel_v2
42
43 if __name__ == "__main__":
44     args = parse_args()
45     func_dict = {"v1": nki_layer_norm_kernel_v1,
46                 "v2": nki_layer_norm_kernel_v2,
47                 }
48
49     device = xm.xla_device()
50     num_rows = args.nrows
51     num_cols = args.ncols
52
53     # Generate toy example
54     input_tensor = torch.rand((num_rows, num_cols), dtype=torch.float32)
55     gamma_vector = torch.rand((num_cols), dtype=torch.float32)
56     beta_vector = torch.rand((num_cols), dtype=torch.float32)
57     epsilon = 1e-5
58
59     # Compute torch layernorm layer in cpu
60     output_torch = layer_norm_layer(input_tensor, epsilon, gamma_vector, beta_vector)
61
62     # Copy tensors to NeuronDevice
63     input_tensor = input_tensor.to(device=device)
64     gamma_vector = gamma_vector.to(device=device)
65     beta_vector = beta_vector.to(device=device)
66
67     print(f">>> Running version {args.version}.")
68     func = func_dict[args.version]
69
70     # add nki_jit decorator
71
72     # Compute NKI layernorm kernel in NeuronDevice

```

(continues on next page)

(continued from previous page)

```

72     xm.mark_step()
73     output_nki = func(input_tensor, epsilon, gamma_vector, beta_vector)
74     xm.mark_step()
75     output_nki = output_nki.to(device='cpu')
76
77     # Accuracy check : Compare the output tensors
78     allclose = torch.allclose(output_torch, output_nki, atol=1e-3, rtol=1e-2)
79     if allclose:
80         print("NKI and Torch match")
81     else:
82         print("NKI and Torch differ")

```

## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- PyTorch reference implementation: `layernorm_torch.py`
- Two versions of NKI kernels: `layernorm_nki_kernel.py`

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts

### Performance mode

Check the performance numbers for `nki_layernorm_kernel_v1` and `nki_layernorm_kernel_v2`, and generate NEFF files for profiling:

```
python3 layernorm_nki_kernel.py --mode perfs
```

### Accuracy mode

Check NKI kernel accuracy against PyTorch implementation:

```
python3 layernorm_torch.py --version v1
python3 layernorm_torch.py --version v2
```

Check optimized Layernorm kernel(`nki_layernorm_kernel_v2`) accuracy against `nki_layernorm_kernel_v1`:

```
python3 layernorm_nki_kernel.py --mode accuracy
```

### Input tensor size

```
python3 layernorm_torch.py --nrows 4096 --ncols 8192
python3 layernorm_nki_kernel.py --nrows 4096 --ncols 8192
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Fused Self Attention

In this tutorial, we implement a kernel to perform the self attention seen in Stable Diffusion 2.1(SD2.1) from Stability AI. The model is available [here](#). In doing so, we learn about:

- The NKI syntax and programming model
- Layout, tiling, and memory management considerations when performing attention computation in NKI
- Fusion techniques for implementing efficient attention kernel

## Background

In SD2.1, the core computation of the self attention is the following. Given

- Q: (seqlen, d\_head)
- K: (seqlen, d\_head)
- V: (seqlen, d\_head)

where d\_head and seqlen represents the head dimension and the seqlen length of the model. The batch dimensions have been removed for simplicity. We would like to compute,

$$S = Q * K.T$$
$$R = \text{softmax}(S) * V$$

When generating images of size 512x512,

$$\text{seqlen} = 4096$$
$$d_{head} = 64$$

We assume the data type of all inputs and outputs to be bfloat16.

## Naive Algorithm

Fig. 7.76 shows the scenario if we compute the attention naively. We would first compute  $S=Q * K.T$ , which has a size of [4096, 4096]. Since the result is in bfloat16, this intermediate matrix has a size of 4096 \* 4096 \* 2 bytes = 32MB, far exceeding the total space available in the SBUF(24MB on NeuronCore-v2). This means that we have to spill data from SBUF to HBM after S is computed, and load it back into SBUF when we compute softmax. This leads to lots of data movements between HBM and SBUF, degrading performance.

## Fusion to Save SBUF Space

To avoid exhausting SBUF space, we would like to avoid computing the entirety of the multiplication of Q and K.T at once. One way is to fuse the softmax computation with the second matrix multiplication.

As shown in the Fig. 7.77, in order to produce one block of the final result, we only need to compute the highlighted strip S1 to compute the block r1 in the final result.

Recall the TensorEngine on NeuronCore-v2 can process a maximum 128 contraction dimension, and the free dimension of the left hand side matrix has a maximum of 128. In the matrix multiplication  $S1 = q1 * K.T$ , as labeled in Fig. 7.77, the size of the free dimension of q1 should be 128 and S1 has a shape of [128, 4096]. Therefore, the size of S1 is 128 \* 4096 \* 2 bytes=1MB, which is 32 times smaller than computing the full intermediate matrix.



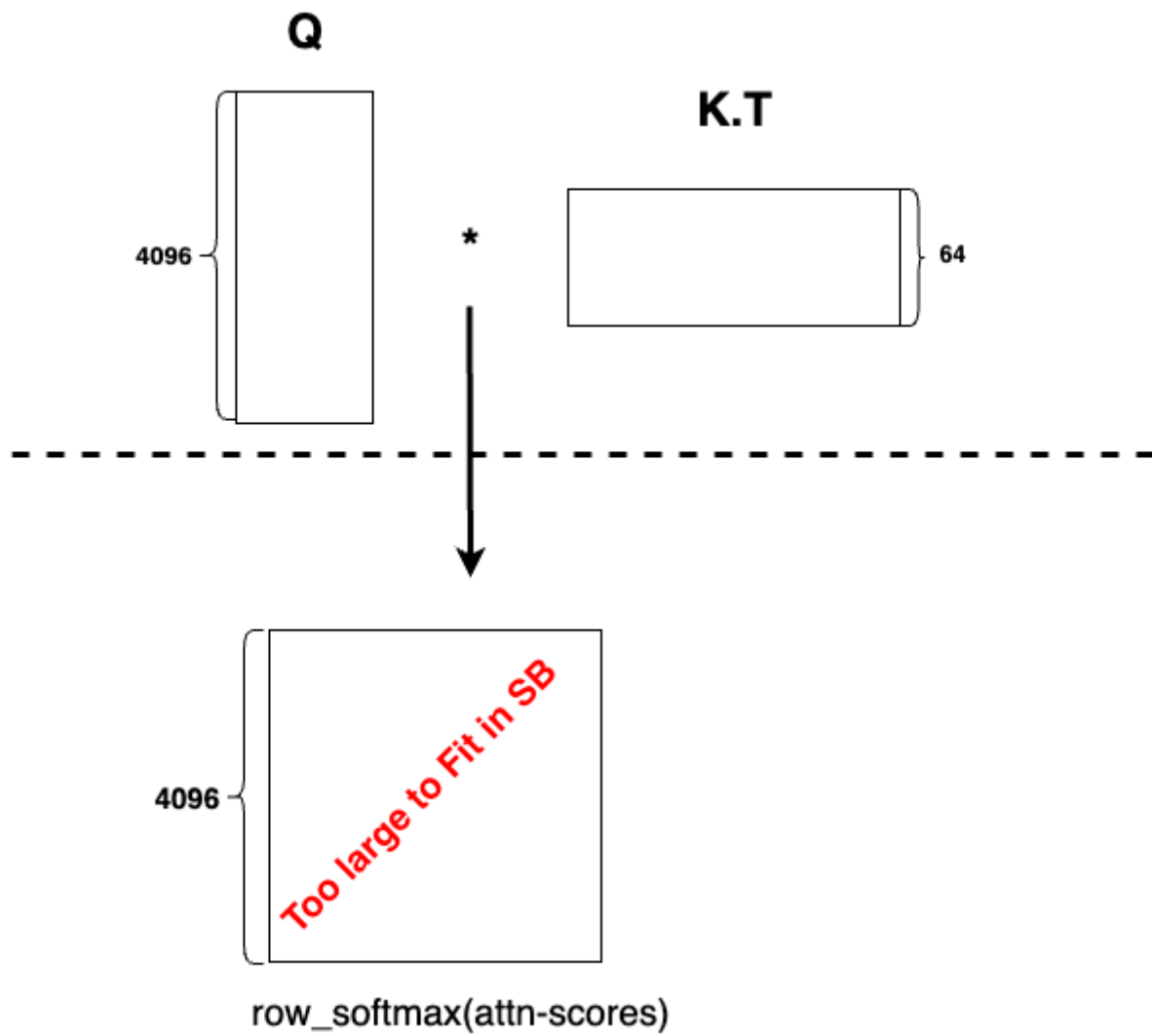


Fig. 7.76: Naively multiple  $Q$  and  $K.T$  produces a large intermediate matrix

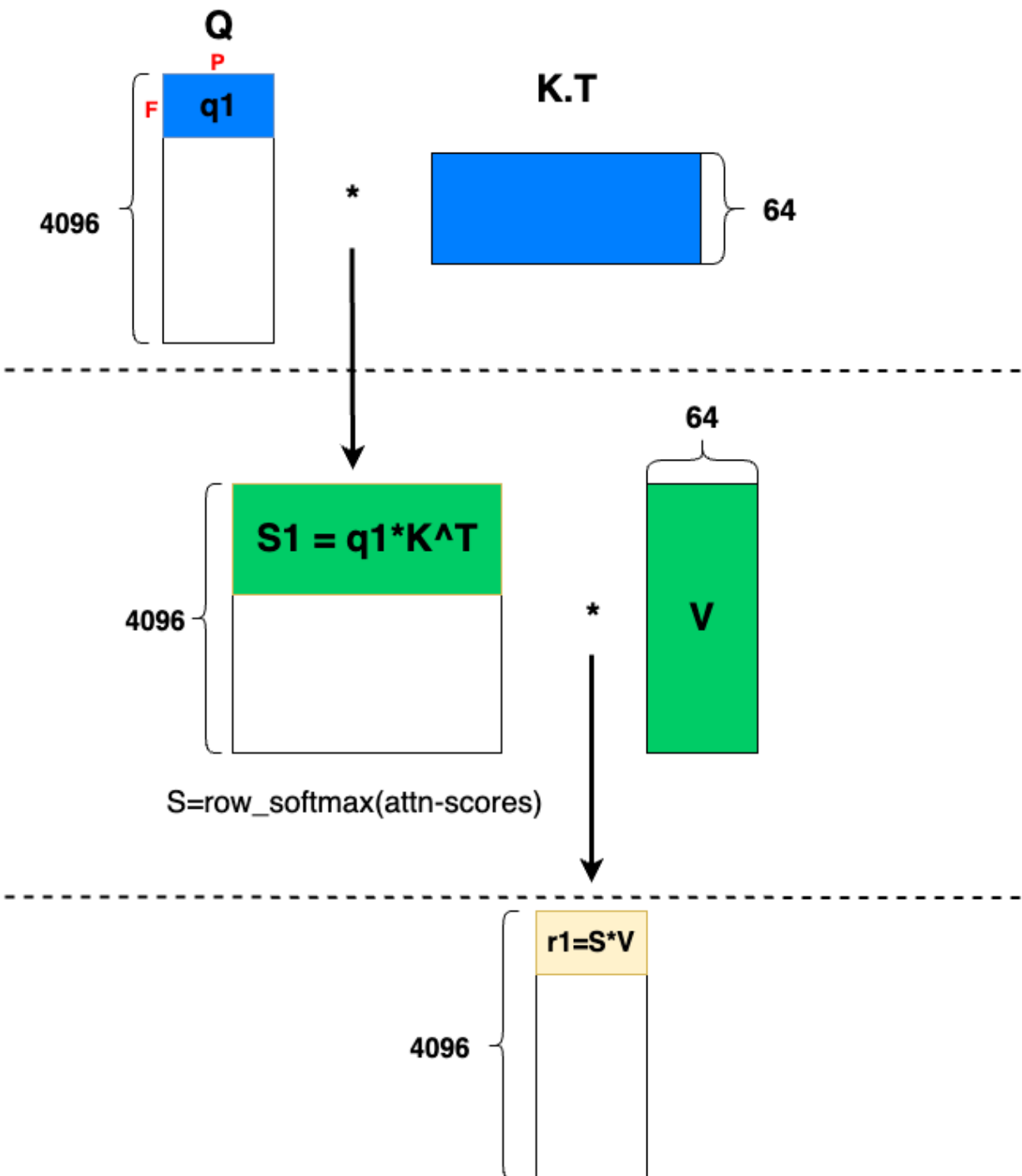


Fig. 7.77: We only need to compute  $S1$  to produce  $r1$

We can then produce the entire result by looping over the tiles in  $Q$ .

## Softmax implementation

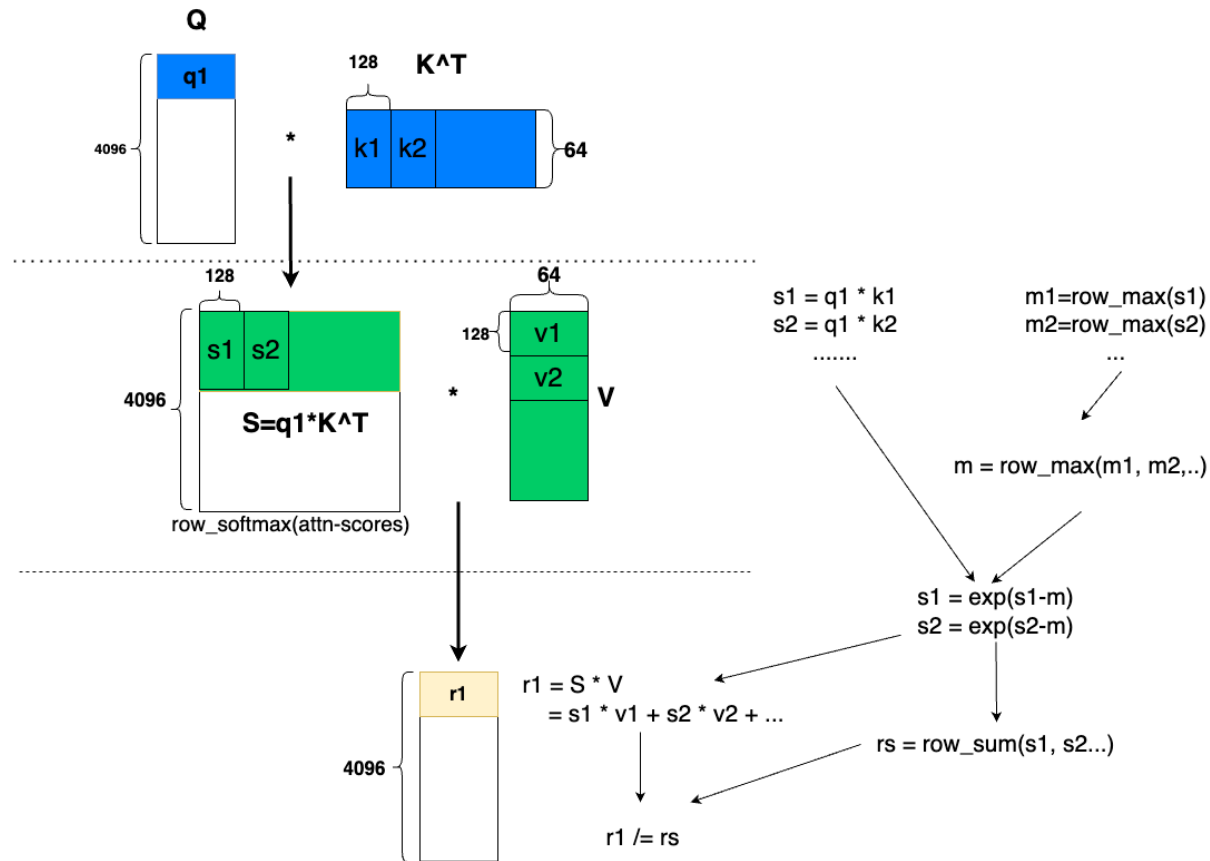


Fig. 7.78: Softmax implementation

We need to perform softmax activation on  $Q * K.T$ , the scheme is shown in the Fig. 7.78. We first compute partial row-wise maximum on each  $s_i$  tile to produce  $m1, m2, \dots$ , then we find the global row-wise maximum  $m$  of  $S$  by computing row-wise maximum on  $m1, m2, \dots$ . After subtracting  $m$  from  $s1, s2, \dots$ , we compute the natural exponential and sum them together to find the row-wise sum  $rs$ .

In a regular softmax, we would divide each  $s1, s2, \dots$  with  $rs$ , however, here we can delay the division to after we compute  $r1$  due to the associativity of scalar-matrix multiplication. Since  $rs$  is smaller than  $r1$ , we save FLOPS by delaying the division. This is also a major optimization deployed in [FlashAttention-v2](#).

We finally multiply  $s_i$  and  $v_i$ , and sum them together to get  $r1$ . By looping over tiles in  $Q$ , we produce the entire result  $R$ .

## Compute kernel

```

1 @nki.jit
2 def fused_self_attn_for_SD_small_head_size(q_ref, k_ref, v_ref, use_causal_mask=False,
3                                           mixed_precision=True):
4     """
5     Fused self attention kernel for small head dimension Stable Diffusion workload,
6     simplified for this tutorial.
7
8     Computes softmax(QK^T)V. Decoder model can optionally include a causal mask
9     application. Does not include QKV projection, output projection, dropout,
10    residual connection, etc.
11
12    This kernel is designed to be used for Stable Diffusion models where the
13    d_head is smaller or equal to 128. Assertion is thrown if `d_head` does
14    not satisfy the requirement.
15
16    IO tensor layouts:
17    - q_ptr: shape (seq_q, d_head)
18    - k_ptr: shape (seq_k, d_head)
19    - v_ptr: shape (seq_v, d_head)
20    - out_ptr: shape (seq_q, d_head)
21    - We use seq_q and seq_k and seq_v just for clarity, this kernel requires
22      seq_q == seq_k == seq_v
23
24    IO tensor dtypes:
25    - This kernel assumes all IO tensors have the same dtype
26    - If mixed_precision is True, then all Tensor Engine operation will be performed in
27      bfloat16 and accumulation will be performed in float32. Otherwise the intermediates
28      will be in the same type as the inputs.
29    """
30    # Use q_ref dtype as the intermediate tensor dtype
31    # Assume all IO tensors have the same dtype
32    kernel_dtype = q_ref.dtype
33    pe_in_dt = nl.bfloat16 if mixed_precision else np.float32
34    assert q_ref.dtype == k_ref.dtype == v_ref.dtype
35
36    # Shape checking
37    seqlen, d_head = q_ref.shape
38    assert d_head <= 128, "Cannot use this kernel for d_head > 128"
39    assert tuple(q_ref.shape) == (seqlen, d_head), 'Input shape mismatch!'
40    assert tuple(k_ref.shape) == (seqlen, d_head), 'Input shape mismatch!'
41    assert tuple(v_ref.shape) == (seqlen, d_head), \
42    f'Input shape mismatch! Expected: {(seqlen, d_head)} Actual: {tuple(v_ref.shape)}'
43    out_ref = nl.ndarray((seqlen, d_head), dtype=q_ref.dtype, buffer=nl.shared_hbm)
44
45    # Softmax scaling factor, multiplied onto Q
46    softmax_scale = 0.125
47
48    q_seq_n_tiles, q_seq_tile_size = seqlen // 128, 128
49    k_seq_n_tiles, k_seq_tile_size = seqlen // 128, 128
50    # No tiling on d_head dimension since the dimension of d_head fits in SB

```

(continues on next page)

(continued from previous page)

```

51 d_head_tile_size = d_head
52 v_seq_n_tiles, v_seq_tile_size = seqlen // 128, 128
53
54 #####
55 # Step 1. transpose(tensor_v)
56 #####
57 # Buffer for v matrix transposed
58 # Pre-fetch and keep it in SBUF throughout different softmax tiles
59 trans_v = nl.ndarray((par_dim(v_seq_tile_size), v_seq_n_tiles, d_head), dtype=pe_in_dt)
60
61 for i_k_seq_tile in nl.affine_range(k_seq_n_tiles):
62     ip_v = nl.arange(v_seq_tile_size)[: , None]
63     if_v = nl.arange(d_head_tile_size)[None, :]
64     trans_v[ip_v, i_k_seq_tile, if_v] = nl.load(
65         v_ref[i_k_seq_tile * k_seq_tile_size + ip_v, if_v],
66         dtype=pe_in_dt)
67
68 q_local = nl.ndarray((q_seq_n_tiles, par_dim(d_head_tile_size), q_seq_tile_size),
69 dtype=pe_in_dt)
70 ip_q = nl.arange(d_head_tile_size)[: , None]
71 if_q = nl.arange(q_seq_tile_size)[None, :]
72 for i_q_seq_tile in nl.affine_range(q_seq_n_tiles):
73     q_local[i_q_seq_tile, ip_q, if_q] = nl.load_transpose2d(
74         q_ref[i_q_seq_tile * q_seq_tile_size + nl.arange(q_seq_tile_size)[: , None],
75             nl.arange(d_head_tile_size)[None, :]]
76         ],
77         dtype=pe_in_dt) * softmax_scale
78
79 k_local = nl.ndarray((k_seq_n_tiles, par_dim(d_head_tile_size), k_seq_tile_size),
80 dtype=pe_in_dt)
81 ip_k = nl.arange(d_head_tile_size)[: , None]
82 if_k = nl.arange(k_seq_tile_size)[None, :]
83 for i_k_seq_tile in nl.affine_range(k_seq_n_tiles):
84     k_local[i_k_seq_tile, ip_k, if_k] = nl.load_transpose2d(
85         k_ref[i_k_seq_tile * k_seq_tile_size + nl.arange(k_seq_tile_size)[: , None],
86             nl.arange(d_head_tile_size)[None, :]]
87         ],
88         dtype=pe_in_dt)
89
90 for i_q_seq_tile in nl.affine_range(q_seq_n_tiles): # indent = 2
91     # A SBUF buffer for an independent softmax tile
92     qk_res_buf = nl.ndarray((par_dim(q_seq_tile_size), seqlen), dtype=kernel_dtype)
93
94     neg_max_res = nl.ndarray((par_dim(q_seq_tile_size), k_seq_n_tiles), dtype=kernel_
95 dtype)
96     ip_max = nl.arange(q_seq_tile_size)[: , None]
97     if_max = nl.arange(k_seq_n_tiles)[None, :]
98
99     # Loop over RHS free of matmul(stationary=tensor_q, moving=tensor_k, contract=d_head)
100     for i_k_seq_tile in nl.affine_range(k_seq_n_tiles): # indent = 4
101
102         # Since the K^T tile is the RHS, the q_seq_len dimension will be P in the result
103         # PSUM buffer shape: [q_seq_tile_size P, k_seq_tile_size F]

```

(continues on next page)

(continued from previous page)

```

100     qk_psum = nl.zeros((par_dim(q_seq_tile_size), k_seq_tile_size),
101                        dtype=np.float32, buffer=nl.psum)
102
103     # Tensor indices for accessing qk result in k_seq_tile_size
104     ip_qk = nl.arange(q_seq_tile_size)[: , None]
105     if_qk = nl.arange(k_seq_tile_size)[None, :]
106
107     #####
108     # Step 2. matmul(stationary=tensor_q, moving=tensor_k, contract=d_head)
109     #####
110     qk_psum[ip_qk, if_qk] += nisa.nc_matmul(moving=k_local[i_k_seq_tile, ip_k, if_k],
111                                             stationary=q_local[i_q_seq_tile, ip_q, if_
112 ↪ q])
113
114     #####
115     # Step 3. Apply optional causal mask
116     #####
117     if use_causal_mask:
118         # Magic number -9984.0 to replace -inf similar to what neuronx-cc uses
119         qk_res_buf[ip_qk, i_k_seq_tile * k_seq_tile_size + if_qk] = nisa.affine_select(
120 ↪ + if_qk),
121         on_true_tile=qk_psum[ip_qk, if_qk], on_false_value=-9984.0, dtype=kernel_dtype)
122     else:
123         # Simply send psum result back to sbuf
124         qk_res_buf[ip_qk, i_k_seq_tile * k_seq_tile_size + if_qk] = nl.copy(qk_psum[ip_
125 ↪ qk, if_qk],
126                                     ↪
127                                     ↪ dtype=kernel_dtype)
128
129     #####
130     # Step 4. Softmax
131     #####
132     neg_max_res[ip_max, i_k_seq_tile] = nisa.tensor_reduce(
133         np.max, data=qk_res_buf[ip_qk, i_k_seq_tile * k_seq_tile_size + if_qk],
134         axis=(1,), dtype=kernel_dtype, negate=True)
135
136     neg_max_res_final = nisa.tensor_reduce(
137         np.min, data=neg_max_res[ip_max, if_max],
138         axis=(1,), dtype=kernel_dtype, negate=False)
139
140     ip_softmax = nl.arange(q_seq_tile_size)[: , None]
141     if_softmax = nl.arange(seq_len)[None, :]
142     ip_sum_res = nl.arange(q_seq_tile_size)[: , None]
143     if_sum_res = nl.arange(d_head_tile_size)[None, :]
144
145     softmax_res = nl.ndarray((par_dim(q_seq_tile_size), seq_len), dtype=pe_in_dt)
146     sum_divisor = nl.ndarray((par_dim(q_seq_tile_size), d_head_tile_size), dtype=kernel_
147 ↪ dtype)
148
149     # Simply use a large tile of seq_len in size since this is a "blocking" instruction
150     # Assuming the compiler will merge exp and reduce_add into a single instruction on_

```

(continues on next page)

(continued from previous page)

↪ACT

```

147     exp_res = nisa.activation(np.exp,
148                             data=qk_res_buf[ip_softmax, if_softmax],
149                             bias=neg_max_res_final, scale=1.0)
150
151     sum_res = nisa.tensor_reduce(np.add, data=exp_res, axis=(1,),
152                                 dtype=kernel_dtype)
153     softmax_res[ip_softmax, if_softmax] = nl.copy(exp_res, dtype=pe_in_dt)
154
155     sum_reciprocal_broadcast = (1.0 / sum_res).broadcast_to((q_seq_tile_size, d_head_
156 ↪tile_size))
157     sum_divisor[ip_sum_res, if_sum_res] = nl.copy(sum_reciprocal_broadcast, dtype=kernel_
158 ↪dtype)
159
160     # Buffer for transposed softmax results (FP32 in PSUM)
161     trans_softmax_res = nl.ndarray(
162         (par_dim(k_seq_tile_size), k_seq_n_tiles, q_seq_tile_size),
163         dtype=pe_in_dt)
164
165     # Result psum buffer has the hidden dim as P
166     attn_res_psum = nl.zeros((par_dim(d_head_tile_size), q_seq_tile_size),
167                             dtype=np.float32, buffer=nl.psum)
168
169     ip_scores_t = nl.arange(k_seq_tile_size)[: , None]
170     if_scores_t = nl.arange(q_seq_tile_size)[None, :]
171     # Loop over matmul_1 contraction
172     for i_k_seq_tile in nl.affine_range(k_seq_n_tiles):
173         #####
174         # Step 5. transpose(softmax_res)
175         #####
176         ip_scores = nl.arange(q_seq_tile_size)[: , None]
177         if_scores = nl.arange(k_seq_tile_size)[None, :]
178
179         trans_softmax_res[ip_scores_t, i_k_seq_tile, if_scores_t] = nisa.nc_transpose(
180             softmax_res[ip_scores, i_k_seq_tile * k_seq_tile_size + if_scores])
181
182         ip_out = nl.arange(d_head_tile_size)[: , None]
183         if_out = nl.arange(q_seq_tile_size)[None, :]
184         for i_k_seq_tile in nl.affine_range(k_seq_n_tiles):
185             #####
186             # Step 6. matmul_1(stationary=trans_v, moving=trans_softmax_res, contract=seqlen_
187 ↪v=seqlen_k)
188             #####
189             ip_v_t = nl.arange(k_seq_tile_size)[: , None]
190             if_v_t = nl.arange(d_head_tile_size)[None, :]
191             attn_res_psum[ip_out, if_out] += \
192                 nisa.nc_matmul(moving=trans_softmax_res[ip_scores_t, i_k_seq_tile, if_scores_t],
193                               stationary=trans_v[ip_v_t, i_k_seq_tile, if_v_t])
194
195         attn_res_sbuf = nl.copy(attn_res_psum[ip_out, if_out], dtype=kernel_dtype)
196
197         attn_res_div = attn_res_sbuf * nisa.nc_transpose(sum_divisor[ip_sum_res, if_sum_res])

```

(continues on next page)

(continued from previous page)

```

195
196     nl.store(
197         out_ref[i_q_seq_tile * q_seq_tile_size + if_out, ip_out],
198         value=attn_res_div)
199
200     return out_ref

```

## Launching kernel and testing correctness

Below we write a reference PyTorch implementation of the attention and verify our NKI kernel output against the reference in the same script as the kernel.

```

1  import torch
2  from torch_xla.core import xla_model as xm
3
4  from sd_attention_nki_kernels import fused_self_attn_for_SD_small_head_size
5
6
7  if __name__ == "__main__":
8
9      device = xm.xla_device()
10
11     def cpu_golden_attn(q, k, v):
12         softmax_scale = 0.125
13         q_scaled = q * softmax_scale
14         raw_score = torch.matmul(q_scaled, k.transpose(1, 0))
15
16         norm_score = torch.nn.functional.softmax(raw_score, dim=-1)
17
18         return torch.matmul(norm_score, v)
19
20     q_tensor = torch.rand((4096, 64), dtype=torch.float32).to(device=device)
21     k_tensor = torch.rand((4096, 64), dtype=torch.float32).to(device=device)
22     v_tensor = torch.rand((4096, 64), dtype=torch.float32).to(device=device)
23
24     output_nki = fused_self_attn_for_SD_small_head_size(q_tensor, k_tensor, v_tensor)
25
26     output_torch = cpu_golden_attn(q_tensor, k_tensor, v_tensor)
27
28     allclose = torch.allclose(output_torch, output_nki, atol=1e-5, rtol=1e-3)
29
30     if allclose:
31         print("NKI and Torch match")
32     else:
33         print("NKI and Torch differ")
34
35     assert allclose

```



## Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- Kernel Definition, accuracy testing and performance benchmark using baremetal mode: `sd_attention_nki_kernels.py`
- Use the kernel in PyTorch: `sd_attention_torch.py`

You can also view the source code in the GitHub repository [nki\\_samples](#)

## Example usage of the scripts:

### Performance mode

Check performance numbers of the attention kernel

```
python3 sd_attention_nki_kernels.py --mode perf
```

### Accuracy mode

Run PyTorch reference implementation and check correctness:

```
python3 sd_attention_torch.py
```

Run baremetal mode and check correctness:

```
python3 sd_attention_nki_kernels.py --mode accuracy
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## Fused Mamba

In this tutorial, we implement a NKI kernel for the [Mamba Large Language Model](#), a State Space Model (SSM) which replaces the attention of a regular Transformer model with a custom layer inspired by Recurrent Neural Networks. We will walk through the core computation step-by-step and map it to NKI APIs to form a functional kernel. Next, by scaling the input shapes of the kernel (both channel size and sequence length), we will iterate on a more hardware-efficient kernel implementation to improve the scaling efficiency.

In this tutorial, we learn about:

- Mapping different vector operations efficiently to NeuronCore compute engines, such as associative scan and element-wise operations between tensors
- Leveraging data reuse and tiling to reduce excessive data movement and keep compute engines busy
- Using [neuron-profile](#) to identify performance bottlenecks and opportunities

## PyTorch Reference Implementation

Before jumping to NKI, let's examine the compute definition of a Mamba-v1 layer using the below PyTorch script (mamba\_torch.py):

```

1  import torch
2  import torch_neuronx
3  import torch_xla.core.xla_model as xm
4  import os
5  import argparse
6
7  os.environ["NEURON_FRAMEWORK_DEBUG"] = "1"
8  os.environ["NEURON_CC_FLAGS"] = " --model-type=transformer --disable-dge "
9
10
11 def associative_scan(deltaA, deltaB_u):
12     """
13     Args:
14         deltaA: [batch_size, channels, state_size, seq_len]
15         deltaB_u: [batch_size, channels, state_size, seq_len]
16
17         Mamba uses an associative scan operator to aggregate information across
18         time sequentially (sequence length, e.g. sequence of tokens),
19         from the past to the present.
20     """
21     batch_size, channels, state_size, seq_len = deltaA.shape
22     out = torch.empty(batch_size, channels, state_size, seq_len,
23                       device=deltaA.device, dtype=deltaA.dtype)
24     for i in range(seq_len):
25         prev_state = out[..., i - 1] if i > 0 else 0
26         out[..., i] = deltaA[..., i] * prev_state + deltaB_u[..., i]
27     return out
28
29
30 def mamba_layer(delta, A, B, u, C):
31     """
32     Args:
33         delta: [batch, channels, seq_len]
34         u: [batch, channels, seq_len]
35         A: [channels, state_size]
36         B: [batch, state_size, seq_len]
37         C: [batch, state_size, seq_len]
38     """
39     # expand the tensors so they all have the same dimensions and compute elementwise
40     # products (with broadcast)
41     # deltaA and deltaB_u have shape [batch_size, channels, state_size, seq_len]
42     deltaA = torch.exp(delta[:, :, None, :] * A[None, :, :, None])
43     deltaB_u = delta[:, :, None, :] * B[:, None, :, :] * u[:, :, None, :]
44     scan_res = associative_scan(deltaA, deltaB_u)
45     # y sums over the `state_size` axis and has shape [batch_size, channels, seq_len]
46     mamba_out = (C[:, None, :, :] * scan_res).sum(dim=-2)
47     return mamba_out

```

(continues on next page)

(continued from previous page)

```

48
49 def parse_args():
50     parser = argparse.ArgumentParser(
51         """Run Mamba PyTorch implementation. Hard-coded small example only since
52         PyTorch implementation is very slow for larger configs.
53         """)
54     parser.add_argument("--mode",
55                         choices=["accuracy", "perf"],
56                         default="accuracy",
57                         help="""Do accuracy test or perf test.
58                             Accuracy test compares mamba_v1 kernel against PyTorch_
59                             implementation.
60                             Perf test will generate a NEFF for the PyTorch_
61                             implementation in local directory
62                             for a manual run of neuron-profile.
63                             """)
64     args = parser.parse_args()
65     return args
66
67 if __name__ == "__main__":
68     args = parse_args()
69
70     # Toy example
71     batch = 1
72     seq_len = 512
73     channels = 256
74     state_size = 16
75
76     dtype = torch.float32
77
78     device = xm.xla_device()
79
80     delta = torch.ones(batch, channels, seq_len, dtype=dtype, device=device)
81     u = torch.ones(batch, channels, seq_len, dtype=dtype, device=device)
82
83     # For numerical accuracy testing purposes, we choose negative numbers for A on_
84     # purpose.
85     # Otherwise, the associative scan will integrate too fast and overflow, which would
86     # mask any real numerical issues in our computation.
87     # A negative A will ensure we catch numerical issues when we have them.
88     A = -torch.ones(channels, state_size, dtype=dtype, device=device)
89     B = torch.ones(batch, state_size, seq_len, dtype=dtype, device=device)
90
91     C = torch.ones(batch, state_size, seq_len, dtype=dtype, device=device)
92
93     xm.mark_step()
94     torch_out = mamba_layer(delta, A, B, u, C)
95     xm.mark_step()
96     print(torch_out)

```

The input tensor shapes are as follows:

- delta: [batch, channels, seq\_len]
- u: [batch, channels, seq\_len]
- A: [channels, state\_size]
- B: [batch, state\_size, seq\_len]
- C: [batch, state\_size, seq\_len]

The key model parameters are:

- batch: batch size of the model.
- seq\_len: sequence length of the model.
- channels: hidden size of a token.
- state\_size: number of model states.

We use [batch=1, seq\_len=512, channels = 256, state\_size = 16] as a simple test case for initial performance evaluation.

Running the above Python script will compile the PyTorch compute graph using Neuron Compiler and generate a Neuron executable file (NEFF) in the same directory. We can then profile the NEFF on a single NeuronCore using *neuron-profiler*. Figure below is a screenshot of the profile. We see this initial PyTorch implementation takes **151.83 ms** to execute *on device*.



Fig. 7.79: Profile of Mamba PyTorch Implementation

Zooming into a portion of the profile, we notice the compute activities on different engines (TensorE/VectorE/ScalarE/GpSimdE) are quite sparse compared to data movement activities (the qSyncIO0 and qVectorSpillReload rows):

In this seemingly “memory-bound” execution trace, the achieved DMA throughput is also extremely low, hovering around 0.33% utilization throughout execution. Therefore, we are stressing neither the compute nor the memory subsystem, hinting the workload is running at low efficiency on the NeuronCore. In the rest of this tutorial, we will showcase how to re-write the above computation using NKI to achieve a device execution latency of **172.93 usec**, which is a **878x speedup** compared to the PyTorch reference implementation.



Fig. 7.80: Profile of Mamba PyTorch Implementation (Zoomed-in)

## Mapping Mamba Layer to NeuronCore

In this section, we will discuss how the computation can be mapped onto the NeuronCore architecture. We will also highlight the importance of choosing appropriate data layouts to achieve good compute efficiency.

Recall we have the following input tensor shapes in device memory:

- `delta`: [batch\_size, channels, seq\_len]
- `u`: [batch\_size, channels, seq\_len]
- `A`: [channels, state\_size]
- `B`: [batch\_size, state\_size, seq\_len]
- `C`: [batch\_size, state\_size, seq\_len]

In fact, the above tensor layout has been chosen carefully based on the computation done in NeuronCore, which we will discuss in more detail below.

In Mamba models, both `seq_len` and `channels` are typically in the thousands (such as `seq_len=16K`, `channels=4K`), while `batch_size` and `state_size` are much smaller by 2-3 order of magnitudes (such as `batch_size=4`, `state_size=16`). To simplify visualization of computation on multi-dimensional tensors, let's hold batch and `state_size` dimension constant and focus on computation per batch per state. Note, the `batch_size` dimension is considered a fully parallel axis in a Mamba layer, while `state_size` is only a partial parallel axis where results from different states will be accumulated together.

By extracting batch and `state_size` dimensions, we get the following input tensor shapes in device memory:

- `delta_i`: [channels, seq\_len]
- `u_i`: [channels, seq\_len]
- `A_i`: [channels]
- `B_i`: [seq\_len]
- `C_i`: [seq\_len]

Next, let's visualize the data flow and computation using 2D matrices or vectors step-by-step.

## Step 1: Element-wise multiplication of `delta_i` and `A_i`

We have the following PyTorch reference code for Step 1:

```
# delta[batch, channels, seq_len]
# A      [channels, state_size]
delta[:, :, None, :] * A[None, :, :, None]

# Holding batch and state_size constant
# delta_i: [channels, seq_len]
# A_i:      [channels]
delta_i[:, :] * A_i[:, :]
```

After the above transformation, the multiplication between `delta_i` and `A_i` involves a **broadcasting** across the `seq_len` dimension of `delta_i`. In NKI, free-dimension broadcast can often be folded into the actual computation instruction at no additional performance cost, while partition-dim broadcast often requires a separate instruction on TensorE (see TensorE alternative use case in [Trainium/Inferentia2 Architecture Guide](#)). As a result, we have two options for executing Step 1.

**Option 1: Map ``seq\_len`` to free dimension.** Element-wise multiplication of `delta_i` and `A_i` on NeuronCore can be done through `nisa.tensor_scalar` on either VectorE or ScalarE, which automatically broadcast `A_i` along the free dimension to match the `seq_len` dimension in `A_i`.

Note, the channels dimension is mapped to SBUF partition dimension. Since the input channels dimension has a size of 256 in our initial setup, which exceeds the architectural limitation of `nl.tile_size.pmax=128`, we must **tile** `delta_i` in the channels dimension (tiled dimension denoted as `channels_tiled`) and feed one tile into `nisa.tensor_scalar` at a time. Figure below illustrates the computation done for Option 1.

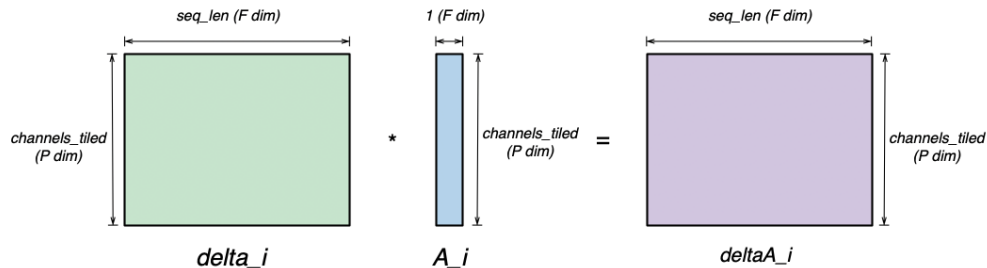


Fig. 7.81: Step 1, Option 1: `nisa.tensor_scalar`

As an example, the associated NKI code for batch `i_batch`, state `i_state` and tile `i_tile_channels` in channels is:

```
# Input shape in device memory matches the computation layout
# Device memory layout:
# delta_i: [channels, seq_len]
# A_i:      [channels]

# Computation layout in SBUF:
# delta_i: [par_dim(channels), seq_len]
# A_i:      [par_dim(channels)]

deltaA_i = nisa.tensor_scalar(delta_i, op0=nl.multiply, operand0=A_i)
```

Note, with this compute layout option, the `delta_i` tensor shape `[channels, seq_len]` in device memory can be

loaded into SBUF efficiently with `seq_len` as the free dimension and fed into VectorE/ScalarE for computation. No extra transposes are needed.

**Option 2: Map ``seq\_len`` to partition dimension.** Alternatively, if we choose a transposed layout for `delta_i` in SBUF for computation, we will need a partition-dimension broadcast of `A_i` using a separate instruction on TensorE (`A_i.broadcast_to(...)`) and then a `nisa.tensor_tensor` operation between `delta_i` and the broadcast `A_i` on VectorE. As a reminder, we need to tile the `seq_len` dimension to meet the tile size constraint `nl.tile_size.pmax=128`. Figure below illustrates the computation done for Option 2.

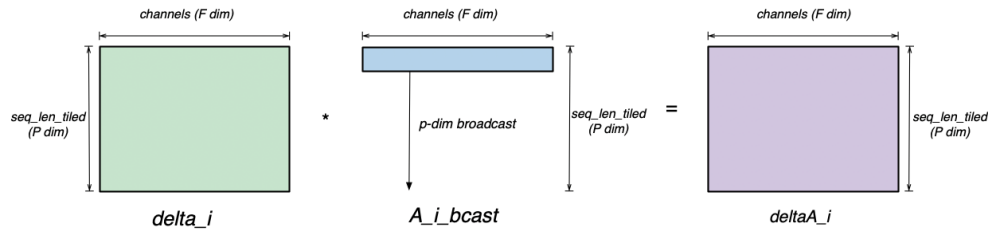


Fig. 7.82: Step 1, Option 2: p-dim broadcast + `nisa.tensor_tensor`

The associated NKI code is as follows:

```
# Input shape in device memory does NOT match the computation layout
# Device memory layout:
# delta_i: [channels, seq_len]
# A_i:      [channels]

# Computation layout in SBUF:
# delta_i: [par_dim(seq_len_tiled), channels]
# A_i:      [par_dim(1), channels]

A_i_bcast = A_i.broadcast_to((nl.tile_size.pmax, channels))
deltaA_i = nisa.tensor_tensor(delta_i, A_i_bcast, op=ml.multiply)
```

Assuming the same `delta_i` device memory layout `[channels, seq_len]`, before performing the `nisa.tensor_tensor` instruction, we will need to either:

- Do a regular load of `delta_i` into SBUF using `nl.load` and an explicit transpose on the loaded `delta_i` using `nl.transpose` to make `seq_len` lie in the free dimension, or
- Do a transposed load of `delta_i` using `nl.load_transpose2d`, which is significantly less efficient in memory bandwidth usage compared to `nl.load`

If Option2 was chosen as the compute layout, we would have incentives to define the `delta` input tensor shape as `[seq_len, channels]` in device memory instead.

From computation perspectives, Option 2 is less efficient than Option 1 because:

1. Option 2 needs an extra TensorE instruction performing partition dimension broadcast.
2. `nisa.tensor_tensor` is 2x slower than `nisa.tensor_scalar` for our input data type FP32 (see API doc for instruction cost estimates).

Therefore, for Step 1 only, Option 1 is the winner compared to Option 2. Let's continue with the rest of the steps to see if we need to revise this selection due to surrounding operator layout preferences.

## Step 2: Exponential of $\delta A_i$ .

Step 2 is evaluating exponential on  $\delta A_i$  from the previous step:

```
torch.exp(...)
```

In NeuronCore, evaluating an exponential function on a tensor is considered a scalar operation, which runs on ScalarE. This operation can be invoked through `nl.exp` or `nisa.activation`. However, ScalarE is able to perform a “pipelined multiply-add” on the input before evaluating a non-linear function (detail see [Trainium/Inferentia2 Architecture Guide](#)). In other words, we can fold Step 1 (Option 1) `nisa.tensor_scalar` and Step 2 into a single ScalarE instruction at no additional cost. This functionality is only exposed in the `nisa.activation` API. This folding is not feasible if we chose Option 2 `nisa.tensor_tensor` in Step 1. Figure below illustrates our new execution plan to combine Step 1 and 2 into `nisa.activation`:

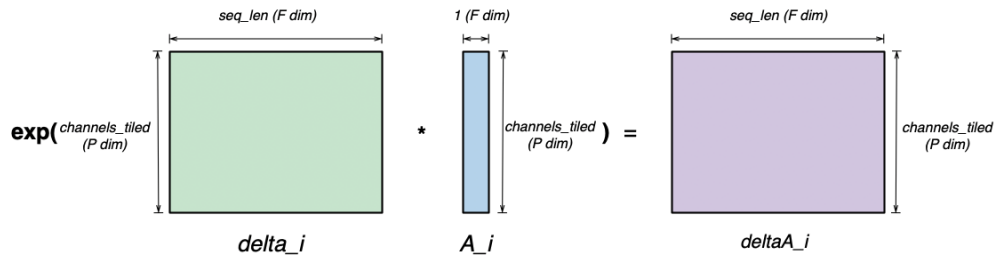


Fig. 7.83: Step 1&2: `nisa.activation`

The associated NKI code is as follows:

```
# Input shape in device memory matches the computation layout
deltaA_i = nisa.activation(op=nl.exp, data=delta_i, scale=A_i)
```

## Step 3: Element-wise multiplication of $\delta A_i$ , $B_i$ and $u_i$ .

PyTorch reference code for Step 3 is:

```
# delta[batch, channels, seq_len]
# B:   [batch, state_size, seq_len]
# u:   [batch, channels, seq_len]
delta[:, :, None, :] * B[:, None, :, :] * u[:, :, None, :]

# Holding batch and state_size constant
# delta_i: [channels, seq_len]
# B_i:     [seq_len]
# u_i:     [channels, seq_len]
delta_i[:, :] * B_i[None, :] * u_i[:, :]
```

This step involves similar compute layout and instruction choices as Step 1:

- `channels` is either partition or free dimension for both  $\delta A_i$  and  $u_i$
- multiplication with  $B_i$  is either through `nisa.tensor_tensor` or `nisa.tensor_scalar`

Since we preferred Step 1 to consume  $\delta A_i$  using `channels` as the partition dimension in previous steps, it is wise to follow the same layout choice here for  $\delta A_i$  to avoid any transposes. Given this layout choice, the multiplication with  $B_i$  will have to be a `nisa.tensor_tensor`. Figure below visualizes the computation in Step 3:



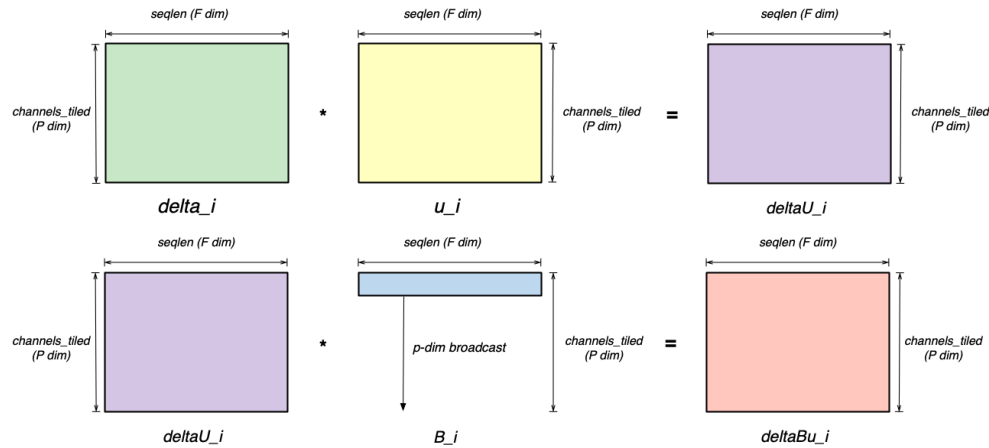


Fig. 7.84: Step 3: p-dim broadcast + 2x nisa.tensor\_tensor

The associated NKI code is as follows:

```
# Input shape in device memory does NOT match the computation layout
# Device memory layout:
# delta_i: [channels, seq_len]
# u_i:      [channels, seq_len]
# B_i:      [seq_len]

# Computation layout in SBUF:
# delta_i: [par_dim(channels_tiled), seq_len]
# u_i:      [par_dim(channels_tiled), seq_len]
# B_i:      [par_dim(1), seq_len]

deltaU_i = nisa.tensor_tensor(delta_i, u_i, op=ml.multiply)
B_i_bcast = B_i.broadcast_to((nl.tile_size.pmax, seq_len))
deltaBu_i = nisa.tensor_tensor(deltaU_i, B_i_bcast, op=ml.multiply)
```

#### Step 4: Associative scan between deltaA\_i and deltaBu\_i

In this step, we use an associative scan operator between deltaA and deltaBu to aggregate information across time sequentially (sequence length, e.g. sequence of tokens), from the past to the present. Here is a PyTorch reference implementation:

```
# deltaA: [batch_size, channels, state_size, seq_len]
# deltaB_u: [batch_size, channels, state_size, seq_len]
out = torch.empty(batch_size, channels, state_size, seq_len,
                  device=deltaA.device, dtype=deltaA.dtype)

for i in range(seq_len):
    # starting state is 0
    prev_state = out[..., i - 1] if i > 0 else 0
    # multiply deltaA by the previous time step state and then add deltaB_u
    out[..., i] = deltaA[..., i] * prev_state + deltaB_u[..., i]
```

By holding batch and state\_size dimensions constant, we get deltaA\_i and deltaBu\_i both with [channels\_tiled,

`seq_len]`, where `channels_tiled` is the partition dimension. The associative scan between these two tile shapes can be implemented in NKI naively through the following loop:

```
scan_i = nl.ndarray((channels_tiled, seq_len), ...)

# Peeling the first iteration out, which is
# equivalent to loop iterator dependent control flow within the loop
scan_i[0:channels_tiled, 0] = deltaBu[0:channels_tiled, 0]

for i in nl.sequential_range(seq_len - 1):
    scan_i[0:channels_tiled, i+1] = deltaA_i[0:channels_tiled, i+1] * scan_i[0:
    ↪ channels_tiled, i]
                                + deltaBu_i[0:channels_tiled, i+1]
```

Within the loop, the current implementation invokes one instruction for multiplication and another for addition. Since both instructions are performed among tiles of shape `[channels_tiled, 1]`, we can combine these two instructions using [\*nisa.tensor\\_scalar\*](#) which supports two operators in a pipelined fashion within an instruction at the same cost as a single operator. Below is a new implementation that could provide 2x speedup compared to the above:

```
scan_i = nl.ndarray((channels_tiled, seq_len), dtype=deltaA.dtype, buffer=nl.sbuf)
scan_i[0:channels_tiled, 0] = deltaBu[i_p, 0]

for i in nl.sequential_range(seq_len - 1):
    scan_i[0:channels_tiled, i+1] = nisa.tensor_scalar(
        deltaA[0:channels_tiled, i+1],
        op0=nl.multiply,
        operand0=scan_i[0:channels_tiled, i],
        op1=nl.add,
        operand1=deltaBu[0:channels_tiled, i+1])
```

However, the above loop nest will turn into `seq_len` many instructions with input tiles that have a single element per partition in SBUF. In addition, every `nisa.tensor_scalar` instruction has a data dependency on the output of the previous instruction. As discussed in the [\*Trainium/Inferentia2 Architecture Guide\*](#), these two traits combined in the instruction sequence is considered extremely *inefficient* on ScalarE/VectorE, where the static instruction overhead instead of the useful execution time would be dominating the engine timeline.

Conveniently, NKI exposes another instruction [\*nisa.tensor\\_tensor\\_scan\*](#) on VectorE, which can perform the above loop nest in a *single* instruction by caching the intermediate scan result from the previous time step internally in VectorE without going through SBUF.

```
scan_i = nisa.tensor_tensor_scan(deltaA_i, deltaBu_i, initial=0,
                                op0=np.multiply, op1=np.add)
```

Note, the shape of `scan_i` is exactly the same as the input `deltaA_i/deltaBu_i`: `[channels_tiled, seq_len]`.

## Step 5: Element-wise multiplication of C\_i and scan\_i

The PyTorch reference implementation is:

```
# scan_res: [batch_size, channels, state_size, seq_len]
# C:        [batch_size, state_size, seq_len]
scanC = C[:, None, :, :] * scan_res

# Holding batch and state constant
# scan_i: [channels_tiled, seq_len]
# C_i:    [seq_len]
scanC_i = C_i[None, :] * scan_i[:, :]
```

You know the drill - Since `channels_tiled` is the partition dimension in `scan_i` from the previous step, we need to perform a partition-dimension broadcast on `C_i` before invoking `nisa.tensor_tensor`:

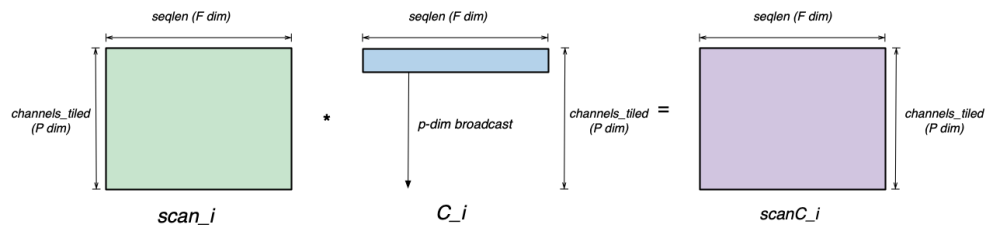


Fig. 7.85: Step 5: p-dim broadcast + `nisa.tensor_tensor`

The corresponding NKI code is:

```
C_i_bcast = C_i.broadcast((nl.tile_size.pmax, seq_len))
scanC_i = nisa.tensor_tensor(scan_i, C_i_bcast, op=ml.multiply)
```

## Step 6: Accumulation of scanC\_i along state\_size dimension

So far in Step 1-5, all the computation is logically parallel across the `state_size` dimension in a Mamba layer. The next step of computation introduces data dependency along the `state_size` dimension for the first time. The PyTorch reference implementation is:

```
# scan_res: [batch_size, channels, state_size, seq_len]
# C:        [batch_size, state_size, seq_len]
# -2 dim is state_size
scanC.sum(dim=-2)

# Holding batch constant only.
# scan_i_states: [channels_tiled, state_size, seq_len]
(scanC_i).sum(dim=-2)
```

In NKI, we can accumulate the `scanC_i` results across states element-wise using `state_size-1` number of `nisa.tensor_tensor` instructions:

Since we will be looping over different states, we can also declare an empty accumulation buffer `scanC_accum` of shape `[channels_tiled, seq_len]` outside of the loop structure and accumulate into this buffer at the end of the every loop iteration using `+=` operator. The use of a single accumulation buffer avoids allocating memory for `scanC_i` across all states in SBUF. The corresponding NKI code is:

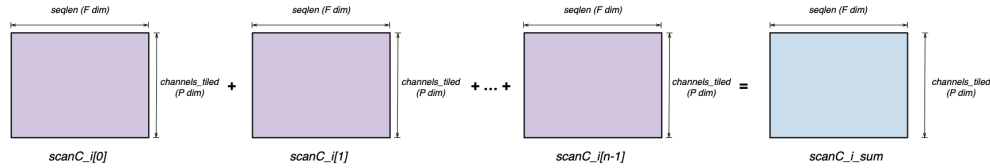


Fig. 7.86: Step 6: state\_size-1 number of nisa.tensor\_tensor

```
scanC_accum = nl.zeros(...)

for i_state in nl.affine_range(state_size):
    scanC_i = ...
    scanC_accum += scanC_i
```

### Initial NKI Kernel

Putting all the pieces together from the previous section, we can arrive at the below kernel implementation `mamba_v1`:

```
1 import neuronxcc.nki as nki
2 import neuronxcc.nki.language as nl
3 import neuronxcc.nki.isa as nisa
4 import numpy as np
5
6 @nki.jit
7 def mamba_v1(delta, u, A, B, C):
8     """Computes the SSM operation in the Mamba model.
9
10    :param delta: (batch_size, channels, seq_len)
11    :param u: (batch_size, channels, seq_len)
12    :param A: (channels, state_size)
13    :param B: (batch_size, state_size, seq_len)
14    :param C: (batch_size, state_size, seq_len)
15    :return: (batch_size, channels, seq_len)
16    """
17    batch_size, channels, seq_len = delta.shape
18    output = nl.ndarray((batch_size, channels, seq_len), dtype=delta.dtype,
19                        buffer=nl.shared_hbm)
20
21    _, state_size = A.shape
22
23    # We can relax this using mask paramters in all the NKI API calls
24    assert channels % 128 == 0
25
26    # Map channels to the partition dimension
27    # Tile channels to comply with NKI tile size constraints
28    channel_psize = nl.tile_size.pmax
29    n_channel_tile = channels // channel_psize
30
31    # Most outer loop with batch_size, parallel_for
32    for i_batch in nl.affine_range(batch_size):
```

(continues on next page)

(continued from previous page)

```

33     # partial accumulated scanC result with processed states
34     scanC_accum = nl.zeros((n_channel_tile, nl.par_dim(channel_psize), seq_len),
    dtype=delta.dtype)
35
36     # Second outer loop with state_size, partial parallel
37     for i_state in nl.affine_range(state_size):
38
39         # Inner loop: tiling channels
40         for i_channel_tile in nl.affine_range(n_channel_tile):
41             channel_start = i_channel_tile * channel_psize
42
43             # Load the relevant tile from delta and A
44             delta_i = nl.load(delta[i_batch, channel_start:channel_start+channel_
    psize, 0:seq_len])
45             A_i = nl.load(A[channel_start:channel_start+channel_psize, i_state])
46
47             # Step 1&2: Element-wise multiplication of delta_i and A_i and then
    exponential
48             deltaA = nisa.activation(op=nl.exp, data=delta_i, scale=A_i)
49
50             # Load the relevant tile from u and B
51             u_i = nl.load(u[i_batch, channel_start:channel_start+channel_psize, 0:
    seq_len])
52             B_i = nl.load(B[i_batch, i_state:i_state+1, 0:seq_len])
53
54             # Step 3: Element-wise multiplication of delta_i, B_i and u_i
55             deltaU = nisa.tensor_tensor(delta_i, u_i, op=nl.multiply)
56             B_i_bcast = B_i.broadcast_to((channel_psize, seq_len))
57             deltaBu = nisa.tensor_tensor(deltaU, B_i_bcast, op=nl.multiply)
58
59             # Step 4: Associative scan between deltaA and deltaBu
60             scan_res = nki.isa.tensor_tensor_scan(deltaA, deltaBu, initial=0,
61             op0=np.multiply, op1=np.add)
62
63             # Load the relevant tile from C
64             C_i = nl.load(C[i_batch, i_state:i_state+1, 0:seq_len])
65
66             # Step 5: Element-wise multiplication of scan_res and C_i
67             C_i_bcast = C_i.broadcast_to((channel_psize, seq_len))
68             scanC = nisa.tensor_tensor(scan_res, C_i_bcast, op=nl.multiply)
69
70             # Step 6: Accumulation of scanC along state_size dimension
71             # scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len] = nisa.tensor_
    tensor(
72             # scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len], scanC,
73             op=nl.add)
74             scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len] += scanC
75
76             # Store scanC_accum for a single batch to output
77             for i_channel_tile in nl.affine_range(n_channel_tile):
78                 channel_start = i_channel_tile * channel_psize
79                 nl.store(output[i_batch, channel_start:channel_start+channel_psize, 0:seq_

```

(continues on next page)

(continued from previous page)

```

79     len],
80         scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len])
81     return output

```

In the above code example,

- **We have three levels of loop nests. From the outer-most to inner-most:**
  - Iterating over `batch`: Different batch samples perform completely different computation. A tensor is the only input parameter that is shared among batch samples.
  - Iterating over `state_size`: Different states perform parallel computation until Step 6 as discussed in the previous section. Both `delta` and `u` tensors are shared across different states.
  - Iterating over `channels`: This is the most-inner dimension where we tile the input channels dimension into `nl.tile_size.pmax=128` chunks. Both `B` and `C` tensors are shared across different channels.
- The kernel above assumes `channels` is a multiple of `nl.tile_size.pmax=128`. We can relax this by adding a `mask` parameter in all the NKI API call in the kernel. To simplify the code example, we omit this change. See [NKI API Masking](#) for more information.
- We declare an empty intermediate tensor `scanC_accum` to hold partial summation from every state.
- **Within the inner loop, we process data for `nl.tile_size.pmax=128` channels for one batch sample in one state.**
  - We use the [slicing syntax](#) to index a tensor. For example, `delta[i_batch, channel_start:channel_start+channel_psize, 0:seq_len]` grabs data from the input `delta` tensor for the current range of channels at the current batch sample.
  - Note, in tensor slicing, the first index dimension from the left with a slicing range will be chosen as the partition dimension. When loading `B`, since we intend to load only one state's worth of data into one partition of SBUF (discussed in Step 3), we need to explicitly slice the state using: `nl.load(B[i_batch, **i_state:i_state+1**, 0:seq_len])`. Otherwise, `nl.load(B[i_batch, **i_state**, 0:seq_len])` will treat `seq_len` as the partition dimension, which is not what we planned for in Step 3 and would also trigger a NKI compilation error since `seq_len` exceeds `nl.tile_size.pmax`.
  - We accumulate partial `scanC_i` results into the accumulation buffer using the `+=` operator. This creates a loop-carried dependency for `scanC_accum` on the `i_state` loop.

## Performance Check

Let's re-run neuron-profile on the above NKI kernel:

Hooray! This NKI kernel implementation now takes 172.93 usec, which is **878x** speedup compared to the reference PyTorch implementation. Based on the profile, VectorE is the busiest compute engine in the Mamba layer. This makes sense because the bulk of computation in the kernel is in `nisa.tensor_tensor`, which can only run on VectorE.

Therefore, our goal is to keep VectorE as busy as possible throughout execution. Note, every NEFF execution involves certain start-up and tear-down overhead. We can use the **Selection Summary** feature in neuron-profile to find out the percentage of time VectorE is busy during the actual execution period:

As indicated by the above profile, VectorE is active over **98.71%** of the time, which is rather impressive. However, remember we used small input shapes as a toy example to get started: `[batch=1, seq_len=512, channels = 256, n = 16]`. Next, let's increase the `channels` and `seq_len` dimensions one by one and observe how VectorE efficiency changes.



Fig. 7.87: Profile of initial Mamba kernel implementation mamba\_v1



Fig. 7.88: Profile of initial Mamba kernel implementation mamba\_v1 (zoomed in)

## Increasing input channels size

Let's increase the size of channels by 16x, from 256 to a more realistic value 4096. We obtain the following profile:

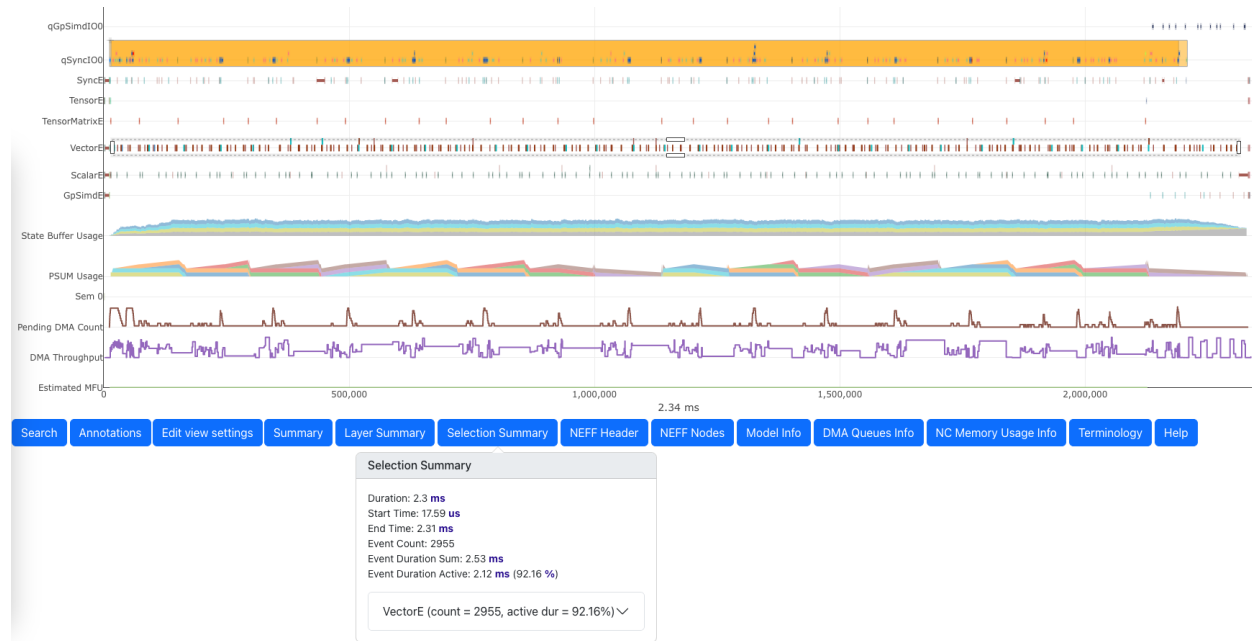


Fig. 7.89: Profile of mamba\_v1 kernel with 4K channels

The new device execution time with increased channels is now **2.34 ms**. We can see that `VectorE` active duration has dropped to **92.16%** during the core execution period, compared to **98.71%** previously with the toy example. Let's zoom into an arbitrary region of the profile to see what could be causing `VectorE` to go idle:

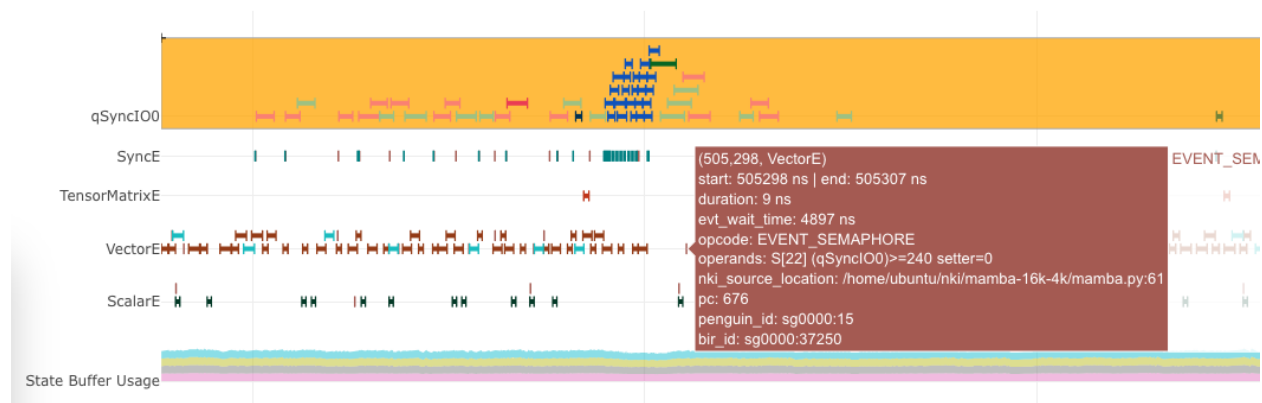


Fig. 7.90: mamba\_v1 kernel blocking on input tensor loading

By identifying a gap where `VectorE` is completely idle, we can hover over the first executed instruction after the gap to find out what's the reason for idleness in the instruction semaphore wait condition. In the above screenshot, the instruction is pending on `S[22]` to reach a value of 240, which is set by `qSyncIO0` activities. This means `VectorE` has been waiting for input tensors to be loaded before performing more computation. If you hover over `qSyncIO0` activities during the `VectorE` idle period, you can also see the exact input tensor name defined in NKI being loaded in the DMA:

We can find similar `VectorE` gaps through the execution trace. At this point, we can conclude one of the reasons why we have a lower `VectorE` active time percentage is due to *blocking* input tensor loading (`n1.load`) activities in the DMA.



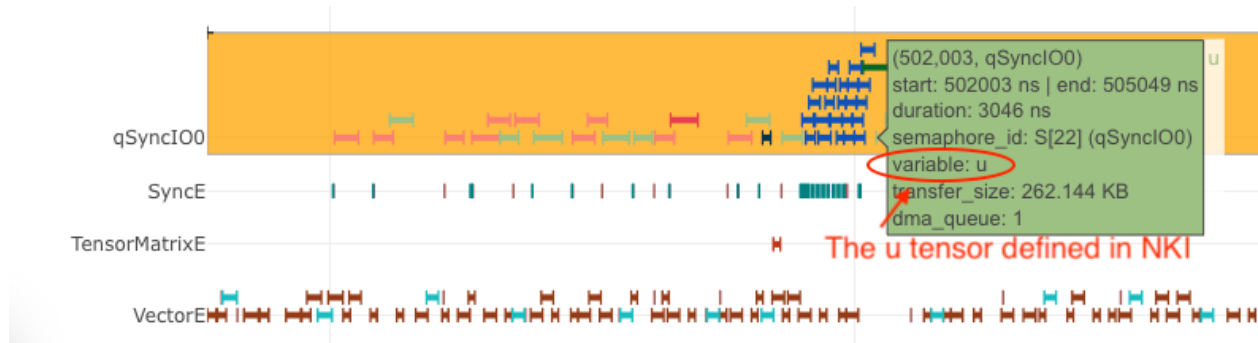


Fig. 7.91: DMA loading tensor u in mamba\_v1 profile

Next, let's spend some time analyzing DMA efficiency.

Zooming out, we can make several observations. First, we see two orange boxes around the qSyncIO0 row. Hovering over the top left corners of the boxes shows two similar performance warnings for loading IO tensors:

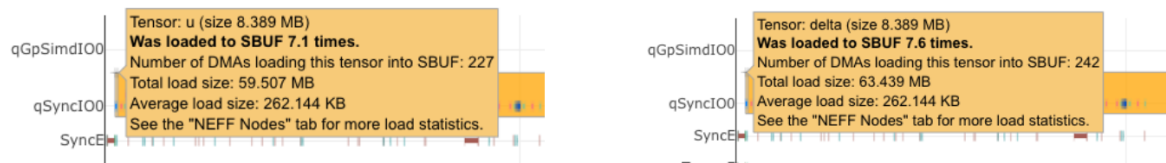


Fig. 7.92: Performance warnings for reloading u and delta tensors

This indicates we reload both the input u and delta tensors around 7 times. This could be inevitable when we don't have sufficient on-chip memory (SBUF) to allow full reuse of the input data tensors. However, the profiler shows we are only hitting around 50% capacity usage throughout execution:



Fig. 7.93: Low SBUF usage

Therefore, the input tensor reloading is likely not justified, and we should investigate whether we can optimize the NKI kernel to avoid it.

## Minimizing data reloading by loop reordering

To understand why delta and u are being reloaded, let's revisit our input tensor shapes:

- delta: [batch\_size, channels, seq\_len]
- u: [batch\_size, channels, seq\_len]
- A: [channels, state\_size]
- B: [batch\_size, state\_size, seq\_len]
- C: [batch\_size, state\_size, seq\_len]

Let's hold `batch_size` constant since the majority of input tensors have completely different slices for different batch samples:

- `delta`: [`channels`, `seq_len`]
- `u`: [`channels`, `seq_len`]
- `A`: [`channels`, `state_size`]
- `B`: [`state_size`, `seq_len`]
- `C`: [`state_size`, `seq_len`]

`delta` and `u` tensors have the same shape with `channels` as the outer dimensions, while `B` and `C` have the same shape with `state_size` as the outer dimension. All four of these input tensors have `seq_len` as the inner dimension. Therefore, we say `delta/u` is reused across different states, while `B/C` are reused across different channels. Given this conflicting reuse dimensions, we further say it is more important to **prioritize reuse of ``delta/u``** because the expected size of `channels` is much higher than `state_size`:

- `state_size` is now 16 and typically stay small
- `channels` is now 4096 and typically in the thousands

In NKI, we can prioritize `delta/u` reuse through loop ordering. Recall in the initial NKI kernel implementation, we have the following inner loops:

```
...
for i_state in nl.affine_range(state_size):
    for i_channel_tile in nl.affine_range(n_channel_tile):
        # step 1-6
...
```

Since these two loops are executed serially within a single NeuronCore, the loop instances will be unrolled by Neuron Compiler. With the channel dimension in the fastest dimension, we will need to load `delta/u` across all channels in the first state, and then likely reload them again in the later states due to a large total memory size in `delta` and `u` (16MB in this case).

To prioritize reuse of `delta/u`, we should reorder the above loop nests. To further enforce the reuse, we can hoist the `nl.load` calls for `delta/u` outside of the `i_state` inner loop:

```
...
for i_channel_tile in nl.affine_range(n_channel_tile):
    delta_i = nl.load(...)
    u_i = nl.load(...)

    for i_state in nl.affine_range(state_size):
        # step 1-6
...
```

As a side effect of this loop re-ordering, we can also spot a loop fusion opportunity since we have two `i_channel_tile` loop nests at the same level now:

```
scanC_accum = nl.zeros((n_channel_tile, nl.par_dim(channel_psize), seq_len), ...)
...

# First i_channel_tile loop
for i_channel_tile in nl.affine_range(n_channel_tile):
    delta_i = nl.load(...)
    u_i = nl.load(...)
```

(continues on next page)

(continued from previous page)

```

    for i_state in nl.affine_range(state_size):
        # step 1-6

# Second i_channel_tile loop
for i_channel_tile in nl.affine_range(n_channel_tile):
    nl.store(..., scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len])

...

```

By fusing the two `i_channel_tile` loop nests into a single loop nest, we can pull the declaration of `scanC_accum` inside the `i_channel_tile` loop and further reduce the `scanC_accum` size requirement by a factor of `n_channel_tile`:

```

...

# First i_channel_tile loop
for i_channel_tile in nl.affine_range(n_channel_tile):
    scanC_accum = nl.zeros((nl.par_dim(channel_psize), seq_len), ...)

    delta_i = nl.load(...)
    u_i = nl.load(...)

    for i_state in nl.affine_range(state_size):
        # step 1-6

    nl.store(..., scanC_accum[i_channel_tile, 0:channel_psize, 0:seq_len])

...

```

Let's modify our initial NKI kernel implementation accordingly to get `mamba_v2`:

```

1 @nki.jit
2 def mamba_v2(delta, u, A, B, C):
3     """Computes the SSM operation in the Mamba model.
4
5     :param delta: (batch_size, channels, seq_len)
6     :param u: (batch_size, channels, seq_len)
7     :param A: (channels, state_size)
8     :param B: (batch_size, state_size, seq_len)
9     :param C: (batch_size, state_size, seq_len)
10    :return: (batch_size, channels, seq_len)
11    """
12    batch_size, channels, seq_len = delta.shape
13    output = nl.ndarray((batch_size, channels, seq_len), dtype=delta.dtype,
14                        buffer=nl.shared_hbm)
15    _, state_size = A.shape
16
17    assert channels % 128 == 0
18
19    # Map channels to the partition dimension
20    # Tile channels to comply with NKI tile size constraints

```

(continues on next page)

(continued from previous page)

```

21 channel_psize = nl.tile_size.pmax
22 n_channel_tile = channels // channel_psize
23
24 # Most outer loop with batch_size, parallel_for
25 for i_batch in nl.affine_range(batch_size):
26
27     # Second outer loop: tiling channels
28     for i_channel_tile in nl.affine_range(n_channel_tile):
29         channel_start = i_channel_tile * channel_psize
30
31         # partial accumulated scanC result with processed states
32         scanC_accum = nl.zeros((nl.par_dim(channel_psize), seq_len), dtype=delta.
↪dtype)
33
34         # Load delta/u once to be reused across states
35         delta_i = nl.load(delta[i_batch, channel_start:channel_start+channel_psize,
↪0:seq_len])
36         u_i = nl.load(u[i_batch, channel_start:channel_start+channel_psize, 0:seq_
↪len])
37
38         # Inner loop with state_size, partial parallel
39         for i_state in nl.affine_range(state_size):
40             # Load the relevant tile from A
41             A_i = nl.load(A[channel_start:channel_start+channel_psize, i_state])
42
43             # Step 1&2: Element-wise multiplication of delta_i and A_i and then
↪exponential
44             deltaA = nisa.activation(op=nl.exp, data=delta_i, scale=A_i)
45
46             # Load the relevant tile from B
47             B_i = nl.load(B[i_batch, i_state:i_state+1, 0:seq_len])
48
49             # Step 3: Element-wise multiplication of delta_i, B_i and u_i
50             deltaU = nisa.tensor_tensor(delta_i, u_i, op=nl.multiply)
51             B_i_bcast = B_i.broadcast_to((channel_psize, seq_len))
52             deltaBu = nisa.tensor_tensor(deltaU, B_i_bcast, op=nl.multiply)
53
54             # Step 4: Associative scan between deltaA and deltaBu
55             scan_res = nki.isa.tensor_tensor_scan(deltaA, deltaBu, initial=0,
56             op0=np.multiply, op1=np.add)
57
58             # Load the relevant tile from C
59             C_i = nl.load(C[i_batch, i_state:i_state+1, 0:seq_len])
60
61             # Step 5: Element-wise multiplication of scan_res and C_i
62             C_i_bcast = C_i.broadcast_to((channel_psize, seq_len))
63             scanC = nisa.tensor_tensor(scan_res, C_i_bcast, op=nl.multiply)
64
65             # Step 6: Accumulation of scanC along state_size dimension
66             scanC_accum[0:channel_psize, 0:seq_len] += scanC
67
68             # Store scanC_accum for a single batch to output

```

(continues on next page)

(continued from previous page)

```

69         nl.store(output[i_batch, channel_start:channel_start+channel_psize, 0:seq_
70         len],
71               scanC_accum[0:channel_psize, 0:seq_len])
72     return output

```

We recapture the profile for the new kernel implementation:



Fig. 7.94: Profile of mamba\_v2 kernel with loop reordering optimization

The device execution time is now **1.61 ms**, which is a **31%** reduction in latency compared to our initial kernel implementation. We can also see VectorE active duration is back up to 99.63% and the performance warnings on input tensor reloading are now gone. In case you are curious, the above loop reordering optimization alone provides around 30% of latency reduction, while the loop fusion optimization contributes the remaining 1% performance boost. This makes sense because the loop reordering addresses our key performance concern around input data reloading, while reducing intermediate tensor size is only a nice-to-have given we were quite low on SBUF usage to begin with.

## Increasing input seq\_len size

Next, let's increase the input seq\_len by **16x**, from 512 to 8192 and recompile the above NKI kernel. Below is the associated performance profile:

The new profile now takes **53.33 ms**, which is **33x longer** than the previous profile. VectorE active duration has dropped down to a new low: 58.93%. Compared to the profile captured with a smaller seq\_len, we notice new DMA activity rows qSyncSpillReload0 and qVectorSpillReload0, which are associated with data movement traffic for intermediate data spill from SBUF into device memory or reload back to SBUF. Zooming into a smaller portion of the profile:

We can see VectorE enters idle states due to a blocking semaphore wait for qSyncSpillReload0 activities, which indicates the extra spill/reload is indeed degrading overall computation performance. In addition, we can see low SBUF usage peaking at merely 50%. Computation and data movement are also not overlapped properly, leading to low average utilization in both compute engines and DMA throughput in the overall timeline.



Fig. 7.95: Profile of mamba\_v2 kernel with 8K seq\_len

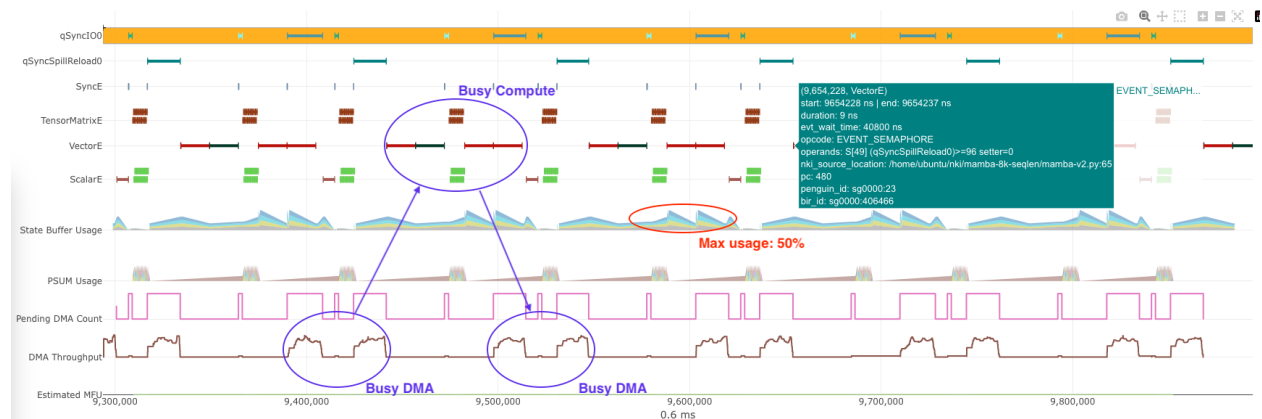


Fig. 7.96: Poor overlap of computation and data movement

Intuitively, increasing `seq_len` of the kernel increases the active tile sizes of input and intermediate tensors in the free dimension, which could cause severe fragmentation in SBUF and excessive data movements to spill/reload tensors in SBUF. To mitigate these inefficiencies, we must **tile** the `seq_len` dimension in our NKI kernel through a new loop level.

### Mitigate spilling by tiling `seq_len`

We have **three** key considerations when adding this new loop level:

1. tile size selection,
2. loop-carried dependency handling
3. loop ordering with other loop nests.

**Tile size of ``seq\_len``.** Since previously with `seq_len=512` in our toy example, we were able to achieve close to 100% VectorE utilization, let's set the tile size `seq_len_fsize` to 512 as a starting point. We can revisit this decision as needed once we obtain a new profile.

**Loop-carried dependency.** Splitting `seq_len` into chunks is straightforward for all computation steps except for Step 4. In the associative scan operation, the next loop iteration requires results from the previous iteration for computation. As a result, we will introduce another loop-carried dependency here with the scan tiles. This dependency can be handled through the `initial` input parameter:

```
scan_init = nl.zeros((channel_psize, 1), ...)

for i_seq_len_tile in static_range(seq_len // seq_len_fsize):
    scan_i = nisa.tensor_tensor_scan(deltaA, deltaBu, initial=scan_init,
                                     op0=np.multiply, op1=np.add)
    scan_init = scan_i[0:channel_psize, seq_len_fsize-1]
```

Note, we choose to use `static_range` instead of `affine_range` due to the new loop-carried dependencies.

**Loop ordering.** Recall from our latest NKI kernel implementation, we have the following loop nest:

```
...
for i_batch in nl.affine_range(batch_size):

    for i_channel_tile in nl.affine_range(n_channel_tile):
        scanC_accum = nl.zeros((nl.par_dim(channel_psize), **seq_len**), ...)

        delta_i = nl.load(delta[i_batch, channel_start:channel_start+channel_psize, 0:
↪ **seq_len**])
        u_i = nl.load(u[i_batch, channel_start:channel_start+channel_psize, 0:**seq_
↪ len**])

        for i_state in nl.affine_range(state_size):
            A_i = nl.load(A[channel_start:channel_start+channel_psize, i_state])

            B_i = nl.load(B[i_batch, i_state:i_state+1, 0:**seq_len**])
            C_i = nl.load(C[i_batch, i_state:i_state+1, 0:**seq_len**])

            deltaA = ...
            deltaBu = ...
            scanC = ...
        ...
```

(continues on next page)

(continued from previous page)

```

scanC_accum += ...

nl.store(..., scanC_accum[i_channel_tile, 0:channel_psize, 0:**seq_len**])

...

```

Let's denote the above loop ordering as `[batch_size, n_channel_tile, state_size]`, and our key question here is where to insert `seq_len` in this list.

Appending `seq_len` to the above list, that is, making `seq_len` the new inner-most loop, would involve the least amount of code changes to our current NKI kernel. However, it will lead to the least amount of SBUF usage reduction, since this loop ordering won't be tiling `scanC_accum`, `delta_i` and `u_i` tensors. Given `seq_len=8192` and FP32 data types, these three tensors will occupy  $8192 \times 3 = 24576$  KiB/partition, half of the available SBUF capacity. Let's go ahead and experiment this loop ordering in a new kernel `mamba_v3`:

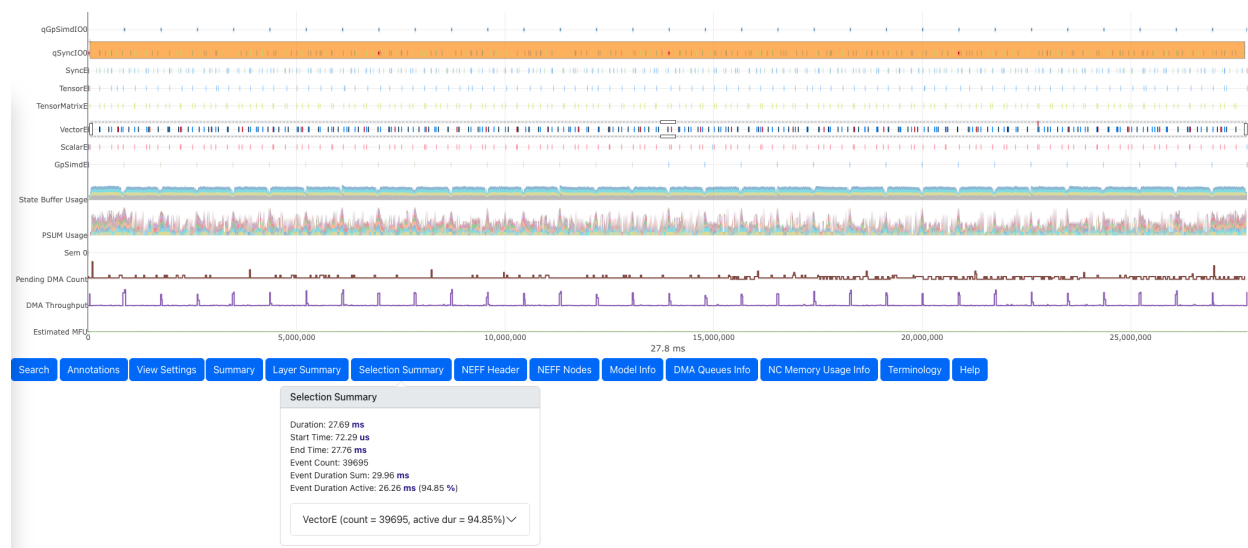


Fig. 7.97: Profile of `mamba_v3` kernel with `seq_len` tiling optimization

With the above profile, the kernel now takes **27.8 ms**, which is **48%** reduction in latency compared to no `seq_len` tiling. `VectorE` is now 94.85% active, and we no longer have spilling related DMA activities.

Finally, since the key advantage of Mamba compared to Transformer models is Mamba's computation and latency should scale linearly with respect to `seq_len`, instead of quadratically in Transformers, let's plot the measured kernel latencies across different `seq_len` up to 8K (what we have optimized so far) and compare it against "perfect latencies" assuming linear scaling from `seq_len=512`. We evaluate scaling efficiency using `perfect latency / measured latency`, which is a higher the better metric. Finally, to showcase the importance of the last `seq_len` tiling optimization for scaling `seq_len`, we also compare scaling efficiency for `mamba_v2` (no `seq_len` tiling) and `mamba_v3` (`seq_len` tiling).



| seq_len | Perfect Latency (ms) | mamba_v2 Measured Latency (ms) | mamba_v2 Scaling Efficiency | mamba_v3 Measured Latency (ms) | mamba_v3 Scaling Efficiency |
|---------|----------------------|--------------------------------|-----------------------------|--------------------------------|-----------------------------|
| 512     | N/A                  | 1.6                            | N/A                         | 1.6                            | N/A                         |
| 1024    | 3.2                  | 4.4                            | 72.73%                      | 3.3                            | 96.97%                      |
| 2048    | 6.4                  | 8.9                            | 71.91%                      | 6.6                            | 96.97%                      |
| 3072    | 9.6                  | 13.1                           | 73.28%                      | 10.1                           | 95.05%                      |
| 4096    | 12.8                 | 17.6                           | 72.73%                      | 13.3                           | 96.24%                      |
| 5120    | 16                   | 23.7                           | 67.51%                      | 17.3                           | 92.49%                      |
| 6144    | 19.2                 | 27.5                           | 69.82%                      | 19.6                           | 97.96%                      |
| 7168    | 22.4                 | 41.3                           | 54.24%                      | 24.2                           | 92.56%                      |
| 8192    | 25.6                 | 52.2                           | 49.04%                      | 27.8                           | 92.09%                      |

The above data shows the last NKI kernel implementation mamba\_v3 can reach 90%+ scaling efficiency up to 8K seq\_len. To support even larger seq\_len, we will need more aggressive tiling by pulling the seq\_len loop level further towards the outer-loop level to tile more input/intermediate tensors to keep spilling low and VectorE busy.

### Download All Source Code

Click the links to download source code of the kernels and the testing code discussed in this tutorial.

- PyTorch reference implementation: `mamba_torch.py`
- Three versions of NKI kernels: `mamba_nki_kernels.py`

You can also view the source code in the GitHub repository [nki\\_samples](#)

### Example usage of the scripts:

#### Performance mode

Run PyTorch reference implementation to generate a NEFF for profiling:

```
python3 mamba_torch.py --mode perf
```

Check performance numbers of mamba\_v1/mamba\_v2/mamba\_v3:

```
python3 mamba_nki_kernels.py --mode perf --version v1 v2 v3 --batch 1 --seq_len 2048 --
↪ channels 512 --state_size 16
```

#### Accuracy mode

Check mamba\_v1 NKI kernel accuracy against PyTorch implementation:

```
python3 mamba_torch.py --mode accuracy
```

Check optimized Mamba kernel (mamba\_v2, mamba\_v3) accuracy against mamba\_v1:

```
python3 mamba_nki_kernels.py --mode accuracy --version v1 v2 v3 --batch 1 --seq_len 2048_
↪ --channels 512 --state_size 16
```

This document is relevant for: Inf2, Trn1, Trn2

- [SPMD Tensor addition](#)

- *SPMD Tensor addition with multiple Neuron Cores*
- *Transpose2D*
- *AveragePool2D*
- *Matrix multiplication*
- *RMSNorm*
- *Fused Self Attention*
- *LayerNorm*
- *Fused Mamba*

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## NKI Kernels

### nki.kernels

The source code of the kernels in the `neuronxcc.nki.kernels` namespace is available at the GitHub Repository [nki-samples](#). They are optimized kernels from the Neuron Team serving as samples. The repository also contains numeric tests, performance benchmarks, as well as scripts to use them in real models.

You are welcome to customize them to fit your unique workloads, and contributing to the repository by opening a PR. Note that these kernels are already being deployed as part of the Neuron stack. With flash attention as an example, [compiling Llama models with transformers-neuronx](#) will automatically invoke the `flash_fwd` kernel listed here. Therefore, replacing the framework operators with these NKI kernels likely won't result in extra performance benefit.

See the [README](#) page of the GitHub Repository [nki-samples](#) for more details.

The documentation of the kernels is available at <https://aws-neuron.github.io/nki-samples>.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Misc

*This document is relevant for:* Inf2, Trn1, Trn2

## NKI FAQ

### When should I use NKI?

NKI enables customers to self serve, onboard novel deep learning architectures, and implement operators currently unsupported by traditional ML Framework operators. With NKI, customers can experiment with models and operators and can create unique differentiation. Additionally, in cases where the compiler's optimizations are too generalized for a developers' particular use case, NKI enables customers to program directly against the Neuron primitives and therefore optimize performance of existing operators that are not being compiled efficiently.

## Which AWS chips does NKI support?

NKI supports all families of chips included in AWS custom-built machine learning accelerators, Trainium and Inferentia. This includes the second generation chips and beyond, available in the following instance types: Inf2, Trn1, Trn1n and Trn2.

## Which compute engines are supported?

The following AWS Trainium and Inferentia compute engines are supported: Tensor Engine, Vector Engine, Scalar Engine, and GpSimd Engine. For more details, see the *NeuronDevice Architecture Guide*, and refer to *nki.isa* APIs to identify which engines are utilized for each instruction.

## How do I launch a NKI kernel onto a logical NeuronCore with Trainium2 from NKI?

A logical NeuronCore (LNC) can consist of multiple physical NeuronCores. In the current Neuron release, an LNC on Trainium2 can have up to two physical NeuronCores (subject to future changes).

For more details on NeuronCore configurations, see [Logical NeuronCore configurations](#).

In NKI, users can launch a NKI kernel onto multiple physical NeuronCores within a logical NeuronCore using single program, multiple data (SPMD) grids.

For a step-by-step guide, refer to the tutorial here: *SPMD Tensor addition with multiple NeuronCores*.

## What ML Frameworks support NKI kernels?

NKI is integrated with *PyTorch* and *JAX* frameworks. For more details, see the *NKI Kernel as a Framework Custom Operator*.

## What Neuron software does not currently support NKI?

NKI does not currently support integration with Neuron Custom C++ Operators, Transformers NeuronX, and Neuron Collective Communication.

## Where can I find NKI sample kernels?

NKI hosts an open source sample repository [nki-samples](#) which includes a set of reference kernels and tutorial kernels built by the Neuron team and external contributors. For more information, see [nki.kernels](#) and *NKI tutorials*.

## What should I do if I have trouble resolving a kernel compilation error?

Refer to *NKI Error Manual* for a detailed guidance on how to resolve some of the common NKI compilation errors.

If you encounter compilation errors from Neuron Compiler that you cannot understand or resolve, you may check out NKI sample [GitHub issues](#) and open an issue if no similar issues exist.

## How can I debug numerical issues in NKI kernels?

We encourage NKI programmers to build kernels incrementally and verify output of small operators one at a time. NKI also provides a CPU simulation mode that supports printing of kernel intermediate tensor values to the console. See [\*nki.simulate\*](#) for a code example.

## How can I optimize my NKI kernel?

To learn how to optimize your NKI kernel, see the [\*NKI Performance Guide\*](#).

## Does NKI support entire Neuron instruction set?

Neuron will iteratively add support for the Neuron instruction set through adding more [\*nki.isa\*](#) (Instruction Set Architecture) APIs in upcoming Neuron releases.

## Will NKI APIs guarantee backwards compatibility?

The [\*NKI APIs\*](#) follow the Neuron Software Maintenance policy for Neuron APIs. For more information, see the [\*SDK Maintenance Policy\*](#).

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

## Neuron Kernel Interface (NKI) release notes

### Neuron Kernel Interface (NKI) (Beta) [2.24]

Date: 06/24/2025

- `sqrt` valid data range extended for accuracy improvement with wider numerical values support.
- [\*nki.language.gather\\_flattened\*](#) new API
- [\*nki.isa.nc\\_match\\_replace8\*](#) additional param `dst_idx`
- improved docs/examples on [\*nki.isa.nc\\_match\\_replace8\*](#), [\*nki.isa.nc\\_stream\\_shuffle\*](#)
- improved error messages

### Neuron Kernel Interface (NKI) (Beta) [2.23]

Date: 05/20/2025

- [\*nki.isa.range\\_select\*](#) (for `trn2`) new instruction
- `abs`, `power` ops supported on to [\*nki.isa\*](#) tensor instruction
- `abs` op supported on [\*nki.isa.activation\*](#) instruction
- [\*GpSIMD engine\*](#) support added to `add`, `multiply` in 32bit integer to [\*nki.isa\*](#) tensor operations
- [\*nki.isa.tensor\\_copy\\_predicated\*](#) support for reversing predicate.
- [\*nki.isa.tensor\\_copy\\_dynamic\\_src\*](#), [\*tensor\\_copy\\_dynamic\\_dst\*](#) engine selection.

- *nki.isa.dma\_copy* additional support with *dge\_mode*, *oob\_mode*, and in-place add *rmw\_op*.
- *+=*, *-=*, */=*, *\*=* operators now work consistently across loop types, PSUM, and SBUF,
- fixed simulation for instructions: *nki.language.rand*, *random\_seed*, *nki.isa.dropout*
- fixed simulation masking behavior
- Added warning when the block dimension is used for SBUF and PSUM tensors, see: *NKI Block Dimension Migration Guide*

## Neuron Kernel Interface (NKI) (Beta) [2.22]

Date: 04/03/2025

- New modules and APIs:
  - *nki.profile*
  - *nki.isa* new APIs:
    - \* *tensor\_copy\_dynamic\_dst*
    - \* *tensor\_copy\_predicated*
    - \* *max8, nc\_find\_index8, nc\_match\_replace8*
    - \* *nc\_stream\_shuffle*
  - *nki.language* new APIs: *mod*, *fmod*, *reciprocal*, *broadcast\_to*, *empty\_like*
- Improvements:
  - *nki.isa.nc\_matmul* now supports PE tiling feature
  - *nki.isa.activation* updated to support reduce operation and *reduce commands*
  - *nki.isa.engine* enum
  - *engine* parameter added to more *nki.isa* APIs that support engine selection (ie, *tensor\_scalar*, *tensor\_tensor*, *memset*)
  - Documentation for *nki.kernels* have been moved to the Github: <https://aws-neuron.github.io/nki-samples>. The source code can be viewed at <https://github.com/aws-neuron/nki-samples>.
    - \* These kernels are still shipped as part of Neuron package in *neuronxcc.nki.kernels* module
- Documentation updates:
  - Kernels public repository <https://aws-neuron.github.io/nki-samples>
  - Updated *profiling guide* to use *nki.profile* instead of *nki.benchmark*
  - NKI ISA Activation functions table now have *valid input data ranges* listed
  - NKI ISA Supported Math operators now have *supported engine* listed
  - Clarify *+=* syntax support/limitation

## Neuron Kernel Interface (NKI) (Beta) [2.21]

Date: 12/16/2024

- New modules and APIs:
  - *nki.compiler* module with Allocation Control and Kernel decorators, see guide for more info.
  - *nki.isa*: new APIs (`activation_reduce`, `tensor_partition_reduce`, `scalar_tensor_tensor`, `tensor_scalar_reduce`, `tensor_copy`, `tensor_copy_dynamic_src`, `dma_copy`), new activation functions (`identity`, `silu`, `silu_dx`), and target query APIs (`nc_version`, `get_nc_version`).
  - *nki.language*: new APIs (`shared_identity_matrix`, `tan`, `silu`, `silu_dx`, `left_shift`, `right_shift`, `ds`, `spmd_dim`, `nc`).
  - New *datatype*: `float8_e5m2`
  - New *kernels* (`allocated_fused_self_attn_for_SD_small_head_size`, `allocated_fused_rms_norm_qkv`) added, kernels moved to public repository.
- Improvements:
  - Semantic analysis checks for *nki.isa* APIs to validate supported ops, dtypes, and tile shapes.
  - Standardized naming conventions with keyword arguments for common optional parameters.
  - Transition from function calls to kernel *decorators* (`jit`, `benchmark`, `baremetal`, `simulate_kernel`).
- Documentation updates:
  - New *Direct Allocation Developer Guide*
  - Tutorial for *SPMD usage with multiple Neuron Cores on Trn2*

## Neuron Kernel Interface (NKI) (Beta)

Date: 12/03/2024

- NKI support for Trainium2, including full integration with Neuron Compiler. Users can directly shard NKI kernels across multiple Neuron Cores from an SPMD launch grid. See *tutorial* for more info. See *Trainium2 Architecture Guide* for an initial version of the architecture specification (more details to come in future releases).
- New calling convention in NKI kernels, where kernel output tensors are explicitly returned from the kernel instead of pass-by-reference. See any *NKI tutorial* for code examples.

## Neuron Kernel Interface (NKI) (Beta) [2.20]

Date: 09/16/2024

- This release includes the beta launch of the Neuron Kernel Interface (NKI) (Beta). NKI is a programming interface enabling developers to build optimized compute kernels on top of Trainium and Inferentia. NKI empowers developers to enhance deep learning models with new capabilities, performance optimizations, and scientific innovation. It natively integrates with PyTorch and JAX, providing a Python-based programming environment with Triton-like syntax and tile-level semantics offering a familiar programming experience for developers. Additionally, to enable bare-metal access precisely programming the instructions used by the chip, this release includes a set of NKI APIs (*nki.isa*) that directly emit Neuron Instruction Set Architecture (ISA) instructions in NKI kernels.

- In addition to documentation, we've included many of the innovative kernels used with-in the neuron-compiler such as [mamba](#) and [flash attention](#) as open-source samples in a new [nki-samples](#) GitHub repository. New kernel contributions are welcome via GitHub Pull-Requests as well as feature requests and bug reports as GitHub Issues. For more information see the [latest documentation](#). Included in this initial beta release is an in-depth [getting started](#), [architecture](#), [profiling](#), and [performance guide](#), along with multiple [tutorials](#), [api reference documents](#), documented [known issues](#) and [frequently asked questions](#).

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## NKI Known Issues

This document outlines some of the known issues and limitations for the NKI beta release.

### Unsupported Syntax:

1. Top-level tensors must be on HBM. The input and output tensors of the top-level NKI kernel (the kernel function decorated with `nki_jit/nki.baremetal` or called by JAX `nki_call`) must be located in HBM. We currently do not support using tensors stored in SBUF or PSUM as the input or output of the top-level kernel. Tensors must be loaded from HBM into SBUF before use, and output tensors must be stored from SBUF back into HBM. See [nl.load](#) and [nl.store](#).
2. Indexing:
  - Tile on SBUF/PSUM must have at least 2 dimensions as described [here](#). If using a 1D tile on SBUF/PSUM, users may get an “Insufficient rank” error. Workaround this by creating a 2D tile, e.g.,

```
buf = nl.zeros((128, ), dtype=dtype, buffer=nl.sbuf) # this won't work
buf = nl.zeros((128, 1), dtype=dtype, buffer=nl.sbuf) # this works
```

  - Users must index their [N, 1] or [1, M] shaped 2D buffers with both indices, do `my_sbuf[0:N, 0]` or `my_sbuf[0, 0:M]` to access them, since accessing in 1D `my_sbuf[0:N]` won't work.
  - Use `nl.arange` for indirect load/store access indexing, `nl.mgrid` won't work. See code examples in [nl.load](#) and [nl.store](#).
  - If indexing with `[0, 0]` gets internal errors, try using `[0:1, 0:1]` or `nl.mgrid[0:1, 0:1]` instead.
  - If indexing with `[0:1, ...]` gets internal errors, try using `[0, ...]` instead.
3. Masks conjunction: Use `&` to combine masks. We do not support using `and` for masks. See examples in [NKI API Masking](#).
4. [nisa.bn\\_stats](#) does not support mask on the reduce dimension, the mask sent to `bn_stats` could not contain any indices from the reduction dimension.
5. Partition dimension broadcasting is not supported on operator overloads (i.e, `+`, `-`, `*`, `/`, `<<`, `>>`, etc), use `nki` language APIs instead (i.e, `nl.add`, `nl.multiply`, ...).
6. When direct allocation API is used, non-IO HBM tensors are not supported.
  - All tensors declared with `buffer=nl.shared_hbm` must be returned as the result of the kernel.
  - Tensors declared with `buffer=nl.hbm` or `buffer=nl.private_hbm` are not allowed.
  - An error “[NKI005] (float32 [128, 512] %'<name of the hbm tensor>':5)0: DRAM location of kind Internal mapping failed. Only input/output/const DRAM location is supported!” will be thrown when such tensor is encountered.

### Unexpected Behavior:

1. Simulation using `nki.simulate_kernel`:
  - Custom data types like `nl.float32r`, `nl.bfloat16`, `nl.float8_e4m3`, and `nl.float8_e5m2` simulate in fp32 precision. Also, NumPy API calls outside of the NKI kernel, such as `np.allclose` may not work with the above types.
2. Execution:
  - Unwritten output tensor will have garbage data. See detail [here](#).
  - `nl.invert` (aka `bitwise_not`) produces incorrect result with `bool` input type, use `int8` type instead.
3. Profiler:
  - When using `neuron-profile` use the flag `--disable-dge` to workaround a temporary issue with DMA information. See the [Profile using neuron-profile](#) section for more details.
4. Optimization:
  - Users need to declare their NKI buffers as small as possible to avoid buffer overflow errors. An error “[GCA046] Some infinite-cost nodes remain” may mean there’s a buffer overflow, workaround this by creating smaller local buffers.
5. Compiler passes:
  - NKI ISA API may not be one-to-one with generated hardware ISA instructions. The compiler may aid in the support of these instruction calls by adding additional instructions.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 7.3 Neuron Custom C++ Operators [Beta]

Neuron Custom C++ Operators enable developers to write C++ Custom Operators (“CustomOps”) that run on NeuronCores. This enables developers to extend operator support beyond what is officially supported by Neuron.

Developers can use standard PyTorch custom operators programming interfaces to leverage Neuron Custom C++ Operators feature. This makes it easy to migrate CPU Custom Operators to Neuron, and implement new beta operators, all without any intimate knowledge of the NeuronCore hardware.

---

**Note:** Neuron Custom C++ Operators feature is currently supported on NeuronCore-v2 architecture only, which is found in Trainium (Trn1) and second-generation Inferentia (Inf2) chips.

---

*This document is relevant for: Inf2, Trn1, Trn2*



### 7.3.1 API Reference Guide

*This document is relevant for: Inf2, Trn1, Trn2*

#### Custom Operators API Reference Guide [Beta]

This page provides the documentation for the C++ API available to creators of Neuron custom C++ operators (see *Neuron Custom C++ Operators [Beta]*).

##### Table of contents

- *Tensor Library*
- *Tensor Accessors*
- *Streaming Accessors*
- *TCM Accessor*
- *Writing Directly to Output Tensor*
- *Using multiple GPSIMD cores*
- *printf()*
- *Library Limitations*

#### Tensor Library

The tensor library used for Neuron custom C++ operators is based upon the PyTorch ATen tensor library. This includes the core Tensor class as well as select operations defined below. Users need to include the `<torch/torch.h>` header to access the tensor library. A small example of using the tensor library looks as follows.

```
#include <torch/torch.h>
...
torch::Tensor a = torch::zeros({32, 32, 3}, torch::kFloat);
```

#### Tensor Factory Functions

The tensor factory functions provide different means for creating new tensors.

They each take in a `size` argument that specifies the size of each dimension of the tensor created (with the exception of `eye`, which takes in two `int64`'s and creates a strictly 2-dimensional identity matrix.)

`c10::TensorOptions` allows the specification of optional properties for the tensor being created. Currently, only the `dtype` property has an effect on tensor construction, and it must be specified. Other properties, such as `layout` may be supported in the future. The example above shows a common way to use factory functions.

The following dtypes are supported:

- `torch::kFloat`
- `torch::kBFloat16`
- `torch::kHalf`

- torch::kInt
- torch::kChar
- torch::kLong
- torch::kShort
- torch::kByte

torch::Tensor **empty**(torch::IntArrayRef size, c10::TensorOptions options)

Creates a tensor filled with uninitialized data, with the specified size and options. Slightly faster than other factory functions since it skips writing data to the tensor.

torch::Tensor **full**(torch::IntArrayRef size, const Scalar &fill\_value, c10::TensorOptions options)

Creates a tensor filled with the specified `fill_value`, with the specified size and options.

torch::Tensor **zeros**(torch::IntArrayRef size, c10::TensorOptions options)

Creates a tensor filled with zeros, with the specified size and options.

torch::Tensor **ones**(torch::IntArrayRef size, c10::TensorOptions options)

Creates a tensor filled with ones, with the specified size and options.

torch::Tensor **eye**(int64\_t n, int64\_t m, c10::TensorOptions options)

Creates a 2-D tensor with ones on the diagonal and zeros elsewhere.

## Tensor Operation Functions

The tensor library provides commonly used operations defined below. The tensor operation functions do not support broadcasting; the shape of the operands must match if applicable.

The library provides two styles of functions for each tensor operation. For functions ending with `_out`, a tensor with the proper size must be provided to which the output is written. This is illustrated in the example below.

```
torch::exp_out(t_out, t_in);
```

Alternatively, for functions that do not end in `_out`, a new tensor that contains the results of the operation is allocated and returned as seen in the example below.

```
torch::Tensor t_out = torch::exp(t_in);
```

|  |
|--|
| <b>Warning:</b> Only operations that are documented below are supported. |
|--|

torch::Tensor &**abs\_out**(torch::Tensor &result, torch::Tensor &self)

torch::Tensor **abs**(torch::Tensor &self)

Computes the absolute value of each element in `self`.

torch::Tensor &**ceil\_out**(torch::Tensor &result, torch::Tensor &self)

torch::Tensor **ceil**(torch::Tensor &self)

Computes the ceiling of the elements of `self`, the smallest integer greater than or equal to each element.

torch::Tensor &**floor\_out**(torch::Tensor &result, torch::Tensor &self)

`torch::Tensor floor(torch::Tensor &self)`  
 Computes the floor of the elements of `self`, the largest integer less than or equal to each element.

`torch::Tensor &sin_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor sin(torch::Tensor &self)`  
 Computes the sine value of the elements of `self`.

`torch::Tensor &cos_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor cos(torch::Tensor &self)`  
 Computes the cosine value of the elements of `self`.

`torch::Tensor &tan_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor tan(torch::Tensor &self)`  
 Computes the tangent value of the elements of `self`.

`torch::Tensor &log_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor log(torch::Tensor &self)`  
 Computes the natural logarithm of the elements of `self`.

`torch::Tensor &log2_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor log2(torch::Tensor &self)`  
 Computes the base-2 logarithm of the elements of `self`.

`torch::Tensor &log10_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor log10(torch::Tensor &self)`  
 Computes the base-10 logarithm of the elements of `self`.

`torch::Tensor &exp_out(torch::Tensor &result, torch::Tensor &self)`

`torch::Tensor exp(torch::Tensor &self)`  
 Computes the exponential of the elements of `self`.

`torch::Tensor &pow_out(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &exponent)`

`torch::Tensor &pow_out(torch::Tensor &result, const torch::Scalar &self, const torch::Tensor &exponent)`

`torch::Tensor &pow_out(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &exponent)`

`torch::Tensor pow(const torch::Tensor &self, const torch::Scalar &exponent)`

`torch::Tensor pow(const torch::Scalar &self, const torch::Tensor &exponent)`

`torch::Tensor pow(const torch::Tensor &self, const torch::Tensor &exponent)`  
 Takes the power of each element in `self` with `exponent`.

`torch::Tensor &clamp_out(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &minval, const torch::Scalar &maxval)`

`torch::Tensor clamp(const torch::Tensor &self, const torch::Scalar &minval, const torch::Scalar &maxval)`  
 Clamps all elements in `self` into the range `[minval, maxval]`.

`torch::Tensor &add_out(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other, const torch::Scalar &alpha = 1)`

torch::Tensor &**add\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other, const torch::Scalar &alpha = 1)

torch::Tensor **add**(const torch::Tensor &self, const torch::Scalar &other, const torch::Scalar &alpha = 1)

torch::Tensor **add**(const torch::Tensor &self, const torch::Tensor &other, const torch::Scalar &alpha = 1)

Adds *other*, scaled by *alpha*, to *input*,

$$out = self + alpha \times other.$$

torch::Tensor &**sub\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other, const torch::Scalar &alpha = 1)

torch::Tensor &**sub\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other, const torch::Scalar &alpha = 1)

torch::Tensor **sub**(const torch::Tensor &self, const torch::Tensor &other, const torch::Scalar &alpha = 1)

torch::Tensor **sub**(const torch::Tensor &self, const torch::Scalar &other, const torch::Scalar &alpha = 1)

Subtracts *other*, scaled by *alpha*, to *input*,

$$out = self - alpha \times other.$$

torch::Tensor &**mul\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor &**mul\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor **mul**(const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor **mul**(const torch::Tensor &self, const torch::Tensor &other)

Multiplies *self* by *other*.

torch::Tensor &**div\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor &**div\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor **div**(const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor **div**(const torch::Tensor &self, const torch::Tensor &other)

Divides *self* by *other*.

---

**Note:** For tensor-tensor bitwise operations, all the bitwise operations are elementwise between two tensors. For scalar-tensor bitwise operations, the scalar is casted to the datatype of the tensor before computing the bitwise operation.

---

torch::Tensor &**bitwise\_and\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor &**bitwise\_and\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor &**bitwise\_and\_out**(torch::Tensor &result, const torch::Scalar &self, const torch::Tensor &other)

torch::Tensor **bitwise\_and**(const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor **bitwise\_and**(const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor **bitwise\_and**(const torch::Scalar &self, const torch::Tensor &other)

Computes the bitwise AND of **self** and **other**. The input tensors must be of integral types.

torch::Tensor &**bitwise\_or\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor &**bitwise\_or\_out**(torch::Tensor &result, const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor &**bitwise\_or\_out**(torch::Tensor &result, const torch::Scalar &self, const torch::Tensor &other)

torch::Tensor **bitwise\_or**(const torch::Tensor &self, const torch::Tensor &other)

torch::Tensor **bitwise\_or**(const torch::Tensor &self, const torch::Scalar &other)

torch::Tensor **bitwise\_or**(const torch::Scalar &self, const torch::Tensor &other)

Computes the bitwise OR of **self** and **other**. The input tensors must be of integral types.

torch::Tensor &**bitwise\_not\_out**(torch::Tensor &result, const torch::Tensor &self)

torch::Tensor **bitwise\_not**(torch::Tensor &result, const torch::Tensor &self)

Computes the bitwise NOT of **self**. The input tensor must be of integral types.

## Class torch::Tensor

### Constructors

Users should not call the Tensor constructor directly but instead use one of the Tensor factory functions.

### Member Functions

template<typename T, size\_t N>

TensorAccessor<T, N, true> **accessor**() const &

Return a **TensorAccessor** for element-wise random access of a Tensor's elements. Scalar type and dimension template parameters must be specified. This const-qualified overload returns a read-only **TensorAccessor**, preventing the user from writing to Tensor elements. See the Tensor Accessors section below for more details.

template<typename T, size\_t N>

TensorAccessor<T, N, false> **accessor**() &

Return a **TensorAccessor** for element-wise random access of a Tensor's elements. Scalar type and dimension template parameters must be specified. This non-const-qualified overload returns a **TensorAccessor** that can be used to both read and write to Tensor elements. See the Tensor Accessors section below for more details.

template<typename T>

TensorReadStreamAccessor<T> **read\_stream\_accessor**() const &

Opens a streaming accessor for read on a tensor. Template parameter T is the scalar type of the tensor data. See Streaming Accessors section below for more details.

template<typename T>

TensorWriteStreamAccessor<T> **write\_stream\_accessor**() &

Opens a streaming accessor for write on a tensor. Template parameter T is the scalar type of the tensor data. See Streaming Accessors section below for more details.

CoherencyEnforcer::Policy **get\_accessor\_coherence\_policy()** const

Get the Tensor accessor coherence policy. See Coherence section below for more details.

void **set\_accessor\_coherence\_policy**(CoherencyEnforcer::Policy policy) const

Set the Tensor accessor coherence policy. See Coherence section below for more details.

TensorTcmAccessor<true> **tcm\_accessor()** const &

Opens a TCM accessor on a tensor. This const-qualified overload returns a read-only `TensorTcmAccessor`, preventing the user from writing to Tensor elements. See TCM Accessor section below for more details.

TensorTcmAccessor<false> **tcm\_accessor()** &

Opens a TCM accessor on a tensor. This non-const-qualified overload returns a `TensorTcmAccessor` that can be used to both read and write to Tensor elements. See TCM Accessor section below for more details.

torch::Tensor &**fill\_**(const torch::Scalar &value) const

Fill a tensor with the specified value.

## Tensor Operators

Tensor &**operator**=(const Tensor &x) &

Tensor &**operator**=(Tensor &&x) &

Assignment operators

## Tensor Accessors

The standard tensor accessor provides element-wise random access to `Tensor` elements. They can be created by calling `Tensor::accessor()`. It can be used similarly to the Pytorch ATen version (see [https://pytorch.org/cppdocs/notes/tensor\\_basics.html#cpu-accessors](https://pytorch.org/cppdocs/notes/tensor_basics.html#cpu-accessors)). However, it is not as fast as other methods of accessing a `Tensor`, such as the streaming accessor or TCM accessor.

**Warning:** The standard tensor accessors can only be used in single core mode. Using standard tensor accessors in multicore mode is undefined behaviour and is going to cause race condition, yielding incorrect result.

## Example Usage

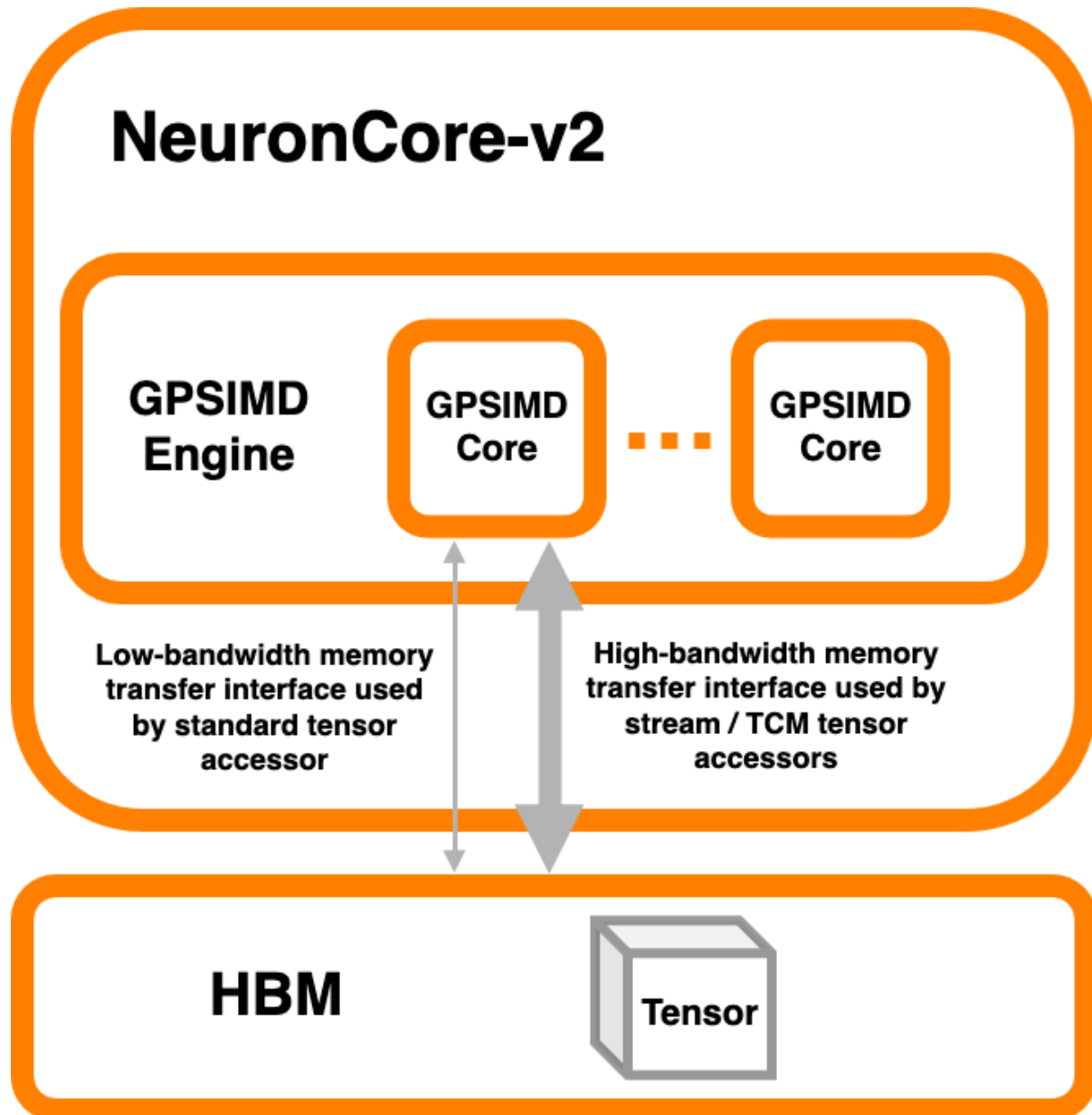
Element-wise add of two 1D tensors using `TensorAccessor`.

```
torch::Tensor tensor_add_compute(const torch::Tensor& t1, const torch::Tensor& t2) {
    size_t num_elem = t1.numel();
    assert(t1.sizes() == t2.sizes());
    torch::Tensor t_out = torch::empty({num_elem}, torch::kFloat);

    auto t1_acc = t1.accessor<float, 1>();
    auto t2_acc = t2.accessor<float, 1>();
    auto t_out_acc = t_out.accessor<float, 1>();
    for (size_t i = 0; i < num_elem; i++) {
        t_out_acc[i] = t1_acc[i] + t2_acc[i];
    }
    return t_out;
}
```

## Memory Architecture

Tensor data is stored in HBM. The various types of accessors enable users to access tensor data from their custom C++ operator code running on the GPSIMD engine.



## Streaming Accessors

Streaming accessors provide the user the ability to access `Tensor` elements in sequential order, faster than the standard tensor accessor. There are two stream accessor classes, one for reading and one for writing. Users should not construct stream accessors directly, but should get them from a `Tensor` using `Tensor::read_stream_accessor` and `Tensor::write_stream_accessor()`.

An active stream accessor is defined as a stream accessor that has been instantiated and not yet closed (via the `close()` method or by going out-of-scope).

The user is responsible for managing stream accessors concurrently accessing the same `Tensor`. For safest usage, no stream accessor should be active while there is an active `TensorWriteStreamAccessor` on the same `Tensor`. The user may either have multiple `TensorReadStreamAccessors` active on the same `Tensor`, or only have a single `TensorWriteStreamAccessor` active on that `Tensor`. Stream accessors should not be used concurrently with standard tensor accessors on the same `Tensor`.

An unlimited number of active stream accessors (in total, across all `Tensors`) are functionally supported, but only up to 4 active stream accessors will be performant. Additional stream accessors beyond the 4th will have performance similar to that of a standard tensor accessor.

**Warning:** Streaming Accessors can only be used in single core mode. Using streaming accessors in multicore mode is undefined behaviour and is going to cause race condition, yielding incorrect result.

## Example Usage

Element-wise add of two tensors using `TensorReadStreamAccessor` and `TensorWriteStreamAccessor`.

```
torch::Tensor tensor_add_compute(const torch::Tensor& t1, const torch::Tensor& t2) {
    assert(t1.sizes() == t2.sizes());
    torch::Tensor t_out = torch::empty(t1.sizes(), torch::kFloat);

    auto t1_rd_stm_acc = t1.read_stream_accessor<float>();
    auto t2_rd_stm_acc = t2.read_stream_accessor<float>();
    auto t_out_wr_stm_acc = t_out.write_stream_accessor<float>();
    for (int i = 0; i < t1.numel(); i++) {
        auto sum = t1_rd_stm_acc.read() + t2_rd_stm_acc.read();
        t_out_wr_stm_acc.write(sum);
    }
    return t_out;
}
```

## Class `torch::TensorWriteStreamAccessor`

`template<typename T> class TensorReadStreamAccessor`

The class template parameter `T` is the scalar type of the tensor data.



## Member Functions

### T read()

Reads from next element in the stream. User is responsible for knowing when to stop reading from `TensorReadStreamAccessor`. Reading past the end of the stream or on a closed stream results in undefined behaviour.

### int close()

Closes stream. Do not read from the stream after calling `close()`.

## Class `torch::TensorWriteStreamAccessor`

### `template<typename T> class torch::TensorWriteStreamAccessor`

The class template parameter `T` is the scalar type of the tensor data.

## Member Functions

### void write(T value)

Writes to next element in the stream. Written value is not guaranteed to be written back to the Tensor's memory until the `TensorWriteStreamAccessor` goes out of scope, or the user explicitly calls `close()`. User is responsible for knowing when to stop writing to a stream accessor. Writing past the end of the stream or on a closed stream results in undefined behaviour.

### int close()

Closes stream. Flushes write data to the Tensor's memory. Do not write to the stream after calling `close()`.

## Coherence

Stream accessors cache Tensor data in GPUSIMD tightly-coupled memory (TCM), but do not ensure their caches remain coherent. When exactly they read from or write back to HBM is opaque to the user (except for `close()` which forces a write back).

The safest way to use them is to ensure that no stream accessor is active (instantiated and not yet closed) while there is an active write stream accessor on the same Tensor. The user should either have multiple read stream accessors active on the same Tensor, or only have a single write stream accessor active on that Tensor.

The standard tensor accessors read/write HBM directly. Therefore, tensor accessors can safely concurrently access the same Tensor, but it is safest not to use them concurrently with stream accessors since HBM isn't guaranteed to be coherent with the stream accessor caches.

These coarse-grained guidelines are best practices, but it is possible to ignore them with careful usage of the accessors (making sure elements are read before they are written to, elements written to are written back before being read again, etc).

The coherence policy of a Tensor determines what to do when there is potentially incoherent access by an accessor of that Tensor. It can either cause an error, or allow it but print a warning, or do nothing. In the case of the latter two options, it is the user's responsibility to ensure they carefully use accessors coherently. Coherence policy for Tensors is `torch::CoherencyEnforcer::Policy::COHERENT` by default, but can be changed using `Tensor::set_accessor_coherence_policy()`.

```
// class torch::CoherencyEnforcer
enum Policy {
    // Enforce a resource is acquired in a way that guarantees coherence
    // Causes an error if it encounters potentially incoherent access
    COHERENT,

    // Allows potentially incoherent access, but will print a warning
    INCOHERENT_VERBOSE,

    // Allows potentially incoherent access, no error or warnings
    INCOHERENT_QUIET
};
```

## TCM Accessor

TCM accessors provide the fastest read and write performance. TCM accessors allow the user to manually manage copying data between larger, but slower-access HBM to faster GP-SIMD tightly-coupled memory (TCM). It may be beneficial to see the diagram under [Memory Architecture](#). Create a `TensorTcmAccessor` from a `Tensor` by calling `Tensor::tcm_accessor()`. Users can allocate and free TCM memory using `tcm_malloc()` and `tcm_free()`. Users have access to a 16KB pool of TCM memory. Note the streaming accessors also allocate from this pool (4KB each). TCM accessors do not do any coherence checks.

---

**Note:** See [Neuron Custom C++ Operators Performance Optimization](#) for a tutorial on how to use TCM accessors.

---

## Example Usage

Element-wise negate of a tensor using `TensorTcmAccessor`.

```
torch::Tensor tensor_negate_compute(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = torch::empty(t_in.sizes(), torch::kFloat);

    static constexpr size_t buffer_size = 1024;
    float *tcm_buffer = (float *)torch::neuron::tcm_malloc(sizeof(float) * buffer_size);

    if (tcm_buffer != nullptr) {
        // tcm_malloc allocated successfully, use TensorTcmAccessor
        auto t_in_tcm_acc = t_in.tcm_accessor();
        auto t_out_tcm_acc = t_out.tcm_accessor();
        for (size_t i = 0; i < num_elem; i += buffer_size) {
            size_t remaining_elem = num_elem - i;
            size_t copy_size = (remaining_elem > buffer_size) ? buffer_size : remaining_
↪elem;

            t_in_tcm_acc.tensor_to_tcm<float>(tcm_buffer, i, copy_size);
            for (size_t j = 0; j < copy_size; j++) {
                tcm_buffer[j] *= -1;
            }
            t_out_tcm_acc.tcm_to_tensor<float>(tcm_buffer, i, copy_size);
        }
    }
}
```

(continues on next page)

(continued from previous page)

```

    }

    torch::neuron::tcm_free(tcm_buffer);
} else {
    // Handle not enough memory...
}

return t_out;
}

```

## TCM Management Functions

void `*torch::neuron::tcm_malloc`(size\_t nbytes)

Allocate nbytes bytes of memory from TCM and return pointer to this memory. Upon failure, returns null.

void `torch::neuron::tcm_free`(void \*ptr)

Free memory that was allocated by `tcm_malloc()`. Undefined behaviour if `ptr` was not returned from a previous call to `tcm_malloc()`.

## Class `torch::TensorTcmAccessor`

`template<bool read_only> class torch::TensorTcmAccessor`

The `read_only` template parameter controls whether or not you can write to the accessor's Tensor. A `const` Tensor will return a read-only `TensorTcmAccessor` from `Tensor::tcm_accessor()`.

## Member Functions

`template<typename T>`

void `tensor_to_tcm`(T \*tcm\_ptr, size\_t tensor\_offset, size\_t num\_elem)

Copy `num_elem` elements from the accessor's Tensor starting at the index `tensor_offset` to a TCM buffer starting at `tcm_ptr`. Tensor indexing is performed as if the tensor was flattened. Template parameter `T` is the scalar type of the tensor data. The TCM buffer's size should be at least `sizeof(T) * num_elem` bytes.

`template<typename T>`

void `tcm_to_tensor`(T \*tcm\_ptr, size\_t tensor\_offset, size\_t num\_elem)

Copy `num_elem` elements from a TCM buffer starting at `tcm_ptr` to the accessor's Tensor starting at the index `tensor_offset`. Tensor indexing is performed as if the tensor was flattened. The TCM buffer's size should be at least `sizeof(T) * num_elem` bytes.

## Writing Directly to Output Tensor

`torch::Tensor get_dst_tensor()`

Returns a reference to the Custom C++ operator output tensor (return value). If this method is called, it is assumed that data will be written to this output tensor, and the tensor returned from the C++ operator will be ignored. Using this method will improve performance by avoiding additional copying of the return value. See example below for function usage.

```
// Example of write to get_dst_tensor()
torch::Tensor example_kernel(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = get_dst_tensor();
    auto t_out_tcm_acc = t_out.tcm_accessor();

    float *tcm_buffer = (float *)torch::neuron::tcm_malloc(sizeof(float) * buffer_
↪size);

    // Populate tcm_buffer with results
    ...
    // Write to t_out through tcm_accessor
    t_out_acc.tcm_to_tensor<float>(tcm_buffer, offset, copy_size);

    ...
}
```

## Using multiple GPSIMD cores

**Note:** See *Neuron Custom C++ Operators Performance Optimization* for a tutorial on how to use multiple GPSIMD cores to execute the Custom C++ Operator

By default, Custom C++ operators target a single core of the GPSIMD-Engine. Performance of Custom C++ operators can be improved by targeting multiple cores. To enable usage of multiple GPSIMD cores, `multicore=True` should be passed to `custom_op.load()`.

```
custom_op.load(
    name=name,
    compute_srcs=compute_srcs,
    shape_srcs=shape_srcs,
    build_directory=os.getcwd(),
    multicore=True
)
```

Each GPSIMD core executes the same kernel function. The user can control the execution on each core by conditioning the Custom C++ operator logic on the core id (obtained via `get_cpu_id()` API). This is illustrated in the example below.

**Warning:** In multicore mode, tensors can only be accessed through TCM accessors. Using regular tensor accessors and streaming accessors are going to yield incorrect result.

The following functions are defined in `neuron/neuron-utils.hpp`

**uint32\_t get\_cpu\_id()**

Return the id of the core that the Custom C++ operator is executing on, id is in range [0, get\_cpu\_count())

**uint32\_t get\_cpu\_count()**

Return the total number of available GPSIMD cores.

```
torch::Tensor example_kernel(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = get_dst_tensor();

    uint32_t cpu_id = get_cpu_id();
    uint32_t cpu_count = get_cpu_count();

    uint32_t partition = num_elem / cpu_count;

    float *tcm_buffer = (float *)torch::neuron::tcm_malloc(sizeof(float) * buffer_size);
    // Populate tcm_buffer with desired results
    ...

    // Write to t_out with a offset computed from cpu_id and cpu_count
    t_out_tcm_acc.tcm_to_tensor<float>(tcm_buffer, partition*cpu_id, copy_size);

    ...
}
```

**Return Value Handling**

When using multiple GPSIMD cores, the `get_dst_tensor()` API must be used to write the return value of the Custom C++ operators. Data not written to the tensor reference returned by `get_dst_tensor()`, or not invoking `get_dst_tensor()` will result in undefined behavior. The user is responsible for writing the appropriate portion of the output reference tensor from a given GPSIMD core. Since there is no synchronization between GPSIMD cores, it is advised that each GPSIMD core writes to a mutually exclusive partition of the output reference tensor.

**printf()**

Custom C++ operators support the use of C++'s `printf()` to send information to the host's terminal. Using `printf()` is the recommended approach to functional debug. With it, the programmer can check the value of inputs, outputs, intermediate values, and control flow within their operator.

**Usage**

To use `printf()` within a Custom C++ operator, the programmer must set the following environment variables before running their model in order to receive the messages printed by their operator:

Table 7.9: Environment Variables

| Name  | Description                             | Type    | Value to Enable printf                                    | Default Value  |
|---|---|---------|---|--|
| NEURON_RUNTIME_LOG_LEVEL                        | Verbose level                           | String  | at least INFO   | See ( <i>NeuronX Runtime Configuration</i> ) for more options.   |
| NEURON_SIZE_GPSIMD_STDOUT_QUEUE_BUFFER_IN_BYTES | Size of the stdio queue buffer in bytes | Integer | Power of two that is equal to or less than 131072 (128KB) | Recommend setting a value of 131072 to maximize the size of printf's buffer. Setting a value of 0 disables printf. |

Within a Custom C++ operator, `printf()` can be used as normal from within a C++ program. For more information, consult a reference such as (<https://cplusplus.com/reference/cstdio/printf/>)

### Example

```
#include <torch/torch.h>
#include <stdio.h> // Contains printf()

torch::Tensor tensor_negate_compute(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = torch::zeros({num_elem}, torch::kFloat);

    auto t_in_acc = t_in.accessor<float, 1>();
    auto t_out_acc = t_out.accessor<float, 1>();
    for (size_t i = 0; i < num_elem; i++) {
        float tmp = -1 * t_in_acc[i];
        printf("Assigning element %d to a value of %f\n", i, tmp);
        t_out_acc[i] = tmp;
    }
    return t_out;
}
```

Print statements then appear on the host's terminal with a header message prepended:

```
2023-Jan-26 00:25:02.0183 4057:4131 INFO TDRV:pool_stdio_queue_consume_all_entries
↳ Printing stdout from GPSIMD:
Assigning element 0 to a value of -1.000000
Assigning element 1 to a value of -2.000000
Assigning element 2 to a value of -3.000000
Assigning element 3 to a value of -4.000000
Assigning element 4 to a value of -5.000000
Assigning element 5 to a value of -6.000000
Assigning element 6 to a value of -7.000000
Assigning element 7 to a value of -8.000000
```

## Limitations

- Performance: using `printf()` significantly degrades the operator's performance.
  - The programmer can disable it by unsetting `NEURON_RT_GPSIMD_STDOUT_QUEUE_SIZE_BYTES` or setting it to 0.
    - \* We recommend that you disable `printf()` if you are running the model in a performance-sensitive context.
  - To maximize performance, remove calls to `printf()` from within the operator.
    - \* Even if `printf()` is disabled, calling the function incurs overhead.
- Buffer size: output from `printf()` is buffered during model execution and read by the Neuron runtime after execution.
  - The model can still execute successfully if you overflow the buffer.
  - Overflowing the buffer causes the oldest data in the buffer to be overwritten.
- Print statements are processed and printed to the host's terminal at the end of model execution, not in real time.
- `printf()` is only supported in single core mode, or on GPSIMD core 0 only when using multiple GPSIMD cores.

## Library Limitations

- Tensors passed into and returned from CustomOp functions can either have up to 8 dimensions where the maximum size of each dimension is 65535, or up to 4 dimensions where the maximum size of each dimension is 4294967295.
- When using multiple GPSIMD cores, only `TensorTcmAccessor` is supported. Usage of other accessors results in undefined behaviour.
- Each model can only have one CustomOp library, and the library can have 10 functions registered. For more information on function registration in PyTorch, see *Implementing an operator in C++* in the [Neuron Custom C++ Operators Developer Guide \[Beta\]](#).
  - However, models using `torch.sort` cannot have any CustomOps.

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

*This document is relevant for:* Inf2, Trn1, Trn2

### 7.3.2 Developer Guide

*This document is relevant for:* Inf2, Trn1, Trn2

## Neuron Custom C++ Operators Developer Guide [Beta]

This document gives an overview of the Neuron Custom C++ Operator feature and APIs . Currently, CustomOp support is limited to the PyTorch framework.

Please refer to the following documents for further information regarding Neuron Custom C++ Operators:

- [Neuron Custom C++ Operators in MLP Training](#)
- [Neuron Custom C++ Operators Performance Optimization](#)
- [Custom Operators API Reference Guide \[Beta\]](#)

### Table of contents

- [Setup & Installation](#)
- [Implementing an operator in C++](#)
- [Building and executing operators](#)
- [Performance Guidance](#)
- [Functional Debug](#)

## Setup & Installation

**Note:** The name of `aws-neuronx-gpsimd-customop` has been changed to `aws-neuronx-gpsimd-customop-lib` as of the neuron 2.10 release.

We provide tooling and library packages (RPM and DEB) that can be installed on TRN1 and INF2 instances:

```
aws-neuronx-gpsimd-tools-0.3
aws-neuronx-gpsimd-customop-lib-0.3
```

For AL2023 only, the following packages need be installed as dependencies:

```
sudo yum install libnsl
sudo yum install libxcrypt-compat
```

On AL2 and AL2023, they can be installed with the following commands:

```
sudo yum remove python3-devel -y
sudo yum remove aws-neuronx-gpsimd-tools-0.* -y
sudo yum remove aws-neuronx-gpsimd-customop-lib-0.* -y

sudo yum install python3-devel -y
sudo yum install aws-neuronx-gpsimd-tools-0.* -y
sudo yum install aws-neuronx-gpsimd-customop-lib-0.* -y
```

On Ubuntu, they can be installed with the following commands:

```
sudo apt-get remove python3-dev -y
sudo apt-get remove aws-neuronx-gpsimd-tools=0.* -y
```

(continues on next page)



(continued from previous page)

```

sudo apt-get remove aws-neuronx-gpsimd-customop-lib=0.* -y

sudo apt-get install python3-dev -y
sudo apt-get install aws-neuronx-gpsimd-tools=0.* -y
sudo apt-get install aws-neuronx-gpsimd-customop-lib=0.* -y

```

## Implementing an operator in C++

Custom operators require a function that defines the custom computation. We define this as the **kernel function**. Neuron Custom C++ Operators also contain a **shape function** separate from the normal compute code. This *shape function* defines the shapes of output tensors for a given set of inputs to the operator. This is needed because PyTorch Neuron (torch-neuronx) is based on the PyTorch/XLA software package and uses a Just-In-Time (JIT) compilation strategy. At runtime the operators in the model will be compiled into a binary to be executed on the NeuronCore. During compilation the shapes of the input and output tensors to operators are computed. The **shape function** is executed on the host, whereas the **kernel function** is executed on the NeuronCore.

### Kernel Function

The kernel function contains the C++ implementation of the CustomOp, as shown in the example below. By including torch.h in the source, the developer has access to a NeuronCore-ported subset of the torch C++ api (<https://pytorch.org/cppdocs/>). The port contains everything required for CustomOp development and model integration, specifically Tensor and Scalar classes in c10, and a subset of aTen operators.

```

#include <stdint.h>
#include <stdlib.h>
#include <torch/torch.h>

torch::Tensor tensor_negate_compute(const torch::Tensor& t_in) {
    size_t num_elem = t_in.numel();
    torch::Tensor t_out = torch::zeros({num_elem}, torch::kFloat);

    auto t_in_acc = t_in.accessor<float, 1>();
    auto t_out_acc = t_out.accessor<float, 1>();
    for (size_t i = 0; i < num_elem; i++) {
        t_out_acc[i] = -1 * t_in_acc[i];
    }
    return t_out;
}

```

The kernel function is the main computational code for the operator. We support a subset of the input types usable by regular PyTorch Custom Operators: torch::Tensor, torch::Scalar, double, and int64\_t. However we do not support std::vector or std::tuple of these types at this time. Note that similar to regular PyTorch Custom Operators, only double and not float, and only int64\_t and not other integral types such as int, short or long are supported. The return value must be a torch::Tensor.

**Warning:** Tensors passed into and returned from CustomOp functions can either have up to 8 dimensions where the maximum size of each dimension is 65535, or up to 4 dimensions where the maximum size of each dimension is 4294967295.

The body of the kernel function may exercise C/C++ libraries, `torch::Tensor` classes, and select `aten` operators, as is customary for Torch programming. For high performance, feature offerings provide faster memory access, via new Tensor Accessor classes and stack management compiler flags. Additionally, higher performance can be obtained by parallelizing execution of the kernel over multiple GPUSIMD cores. See the [Custom Operators API Reference Guide \[Beta\]](#) for more details.

Finally, because the kernel is specially compiled for and run by the NeuronCore target, its tooling, libraries, and environment differ from the host pytorch installation. For example, while the host may run Pytorch 1.13 and a C++17 compatible compiler in a linux environment, the NeuronCore may run a port of Pytorch 1.12 (c10) and LLVM's libc++ C++14 version 10.0.1 without linux. Developers must develop for the compiler, torch version, and environment of their targeted NeuronCore. See the [Custom Operators API Reference Guide \[Beta\]](#) for more details.

## Shape Function

The shape function has the same function signature as the kernel function, but does not perform any computations. Rather, it only defines the shape of the output tensor but not the actual values.

```
#include <stdint.h>
#include <stdlib.h>
#include <torch/torch.h>

torch::Tensor tensor_negate_shape(torch::Tensor t1) {
    size_t num_elem = t1.numel();
    torch::Tensor t_out = torch::zeros({num_elem}, torch::kFloat);

    return t_out;
}
```

The body of the shape function may exercise C/C++ libraries or `torch::Tensor` classes. The body may not access the data of input tensors since these are XLA Tensors and do not have any data storage allocated yet. However, any of the functions that access shape information such as `numel` (to get the number of elements) may be used.

## Building and executing operators

Once you have the kernel and shape functions for your operators you can build them into a library to use them from PyTorch in your model. Just like regular PyTorch Custom Operators, Neuron Custom C++ Operators use a registration macro to associate the kernel and shape functions with the name of the operator that will be called from Python.

Similar to PyTorch, Neuron Custom C++ Operators are grouped into libraries defined within the `NEURON_LIBRARY(<lib_name>, m)` scope, where `lib_name` is the name of your library of custom operators. Within this scope, calls to `m.def(<op_name>, <shape_fcn>, <kernel_fcn>)` define each operator in your library. The `op_name` is the name to call the operator with in the model (i.e. `torch.ops.lib_name.op_name()`). The `shape_fcn` is a function pointer to the shape function to call during compilation. Finally the `kernel_fcn` is the name of the function to be executed on the NeuronCore at runtime.

```
#include <stdint.h>
#include <stdlib.h>
#include <torch/torch.h>
#include "torchneuron/register.h"

torch::Tensor tensor_negate_shape(torch::Tensor t1) {
    size_t num_elem = t1.numel();
    torch::Tensor t_out = torch::zeros({num_elem}, torch::kFloat);
```

(continues on next page)

(continued from previous page)

```

    return t_out;
}

NEURON_LIBRARY(my_ops, m) {
    m.def("tensor_negate", &tensor_negate_shape, "tensor_negate_compute");
}

```

Notice that the `NEURON_LIBRARY` macro is used in the same C++ file as the `shape` function. This is because the registration is loaded on the host.

**Warning:** Each model can only have one CustomOp library, and the library can have 10 functions registered. However, models using `torch.sort` cannot have any CustomOps.

The custom op library is built by calling the `load` API in Python like:

```

import torch_neuronx
from torch_neuronx.xla_impl import custom_op

custom_op.load(
    name='my_ops',
    compute_srcs=['kernel.cpp'],
    shape_srcs=['shape.cpp'],
    multicore=False
)

```

In the example above, `name` refers to the name of the library file to be created (i.e. `libmy_ops.so`) and the `compute_srcs` and `shape_srcs` are lists of files to be compiled. After the `load` API completes, the library will have been compiled and loaded into the current PyTorch process.

**Warning:** The library file name should not be “builtin” as it is a reserved keyword.

CustomOp also supports multicore execution mode. If you want the library to run in multicore mode, pass the flag `multicore=True` into the `load` API. Notice that the execution mode is specified at the library level, so all the functions in the library run in the same mode. For more details of multicore CustomOp, please refer to *Using multiple GPUSIMD cores* section in *Custom Operators API Reference Guide [Beta]*.

Similar to PyTorch, the Neuron custom op will be available at `torch.ops.<lib_name>.<op_name>` where `lib_name` is defined in the `NEURON_LIBRARY` macro, and `op_name` is defined in the call to `m.def`.

```

import torch

out_tensor = torch.ops.my_ops.tensor_negate(in_tensor)

```

## Loading a previously built library

The library can also be built ahead of time or in a separate process and loaded later. In the load API, specify the `build_directory` argument and the library will be written to that location on disk.

```
import torch_neuronx
from torch_neuronx.xla_impl import custom_op

custom_op.load(
    name='my_ops',
    compute_srcs=['kernel.cpp'],
    shape_srcs=['shape.cpp'],
    build_directory=os.getcwd(),
)
```

Then, later, this library can be loaded by calling the `load_library` API and using the ops in the exact same way.

```
import torch
import torch_neuronx
from torch_neuronx.xla_impl import custom_op

custom_op.load_library('/home/user/libmy_ops.so')

out_tensor = torch.ops.my_ops.tensor_negate(in_tensor)
```

Note: The `load_library` API does not need to be called in the same process where the library is built with the `load` API. Similar to regular PyTorch Custom Operators, Neuron Custom C++ Operators are built and loaded at the same time when the `load` API is called.

## Performance Guidance

When possible, it is recommended that operators supported by the designated framework with supported compilation onto Neuron devices are used. These operators have been highly optimized for the Neuron architecture. However, for other scenarios where Custom C++ operators are the required solution, the following recommendations can be followed to improve performance:

- Use the provided memory management accessors (streaming and tcm accessor). Both of these accessors improve data fetch overhead. See the [Custom Operators API Reference Guide \[Beta\]](#) for more information.
- You can optionally specify the estimated amount of stack space (in bytes) used in your Custom C++ operator via the `extra_cflags` argument in the call to `custom_op.load()`. For instance, if you anticipate your operator using ~20KB of stack space, include the argument `extra_cflags=['-DSTACK_SIZE=20000']` in the call to `custom_op.load()`. **This is necessary only if you anticipate the stack to grow beyond ~8KB.** This flag is used to decide whether to place the stack in faster local memory, which significantly improves performance, or if we will need to place the stack in larger NeuronCore memory with longer access latency. If you do not specify this flag, or the estimate you provide is small enough (less than ~8KB), the stack will go in local memory. Note, when placed in local memory, the stack space will not be restricted by your estimate, but if your stack grows beyond ~8KB, there's a risk of a stack overflow, and you will be notified with an error message from GPSIMD should such a case occur. If you do specify a stack size, the maximum supported stack size is 400KB.
- Use multiple GPSIMD cores when possible to parallelize (and hence improve performance) of Custom C++ operator, refer to *Using multiple GPSIMD cores* section in [Custom Operators API Reference Guide \[Beta\]](#) for more information.

## Functional Debug

Custom C++ operators support the use of the C++ language's `printf()`. For functional debug, the recommended approach is using `printf()` to print input, intermediate, and final values. Consult the [Custom Operators API Reference Guide \[Beta\]](#) for more information.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 7.3.3 Tutorials

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 7.3.4 Misc (Neuron Custom C++ Operators)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## Neuron Custom C++ Tools Release Notes

---

**Note:** Neuron Custom C++ Operators feature is currently supported on NeuronCore-v2 architecture only, which is found in Trainium (Trn1) and second-generation Inferentia (Inf2) chips.

---

### aws-neuronx-gpsimd-tools [0.13]

Date: 12/12/2024

- Neuron Custom C++ Operators feature is currently supported on NeuronCore-v2 architecture only, which is found in Trainium (Trn1) and second-generation Inferentia (Inf2) chips.

### aws-neuronx-gpsimd-tools [0.1]

Date: 02/08/2023

- First release of aws-neuronx-gpsimd-tools. This release provides the required tools to support the building of Neuron Custom C++ operators.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### Neuron Custom C++ Library Release Notes

---

**Note:** Neuron Custom C++ Operators feature is currently supported on NeuronCore-v2 architecture only, which is found in Trainium (Trn1) and second-generation Inferentia (Inf2) chips.

---

#### aws-neuronx-gpsimd-customop-lib [0.13]

Date: 12/12/2024

- Neuron Custom C++ Operators feature is currently supported on NeuronCore-v2 architecture only, which is found in Trainium (Trn1) and second-generation Inferentia (Inf2) chips.

#### aws-neuronx-gpsimd-customop-lib [0.3]

Date: 04/28/2023

- Add initial support for using Multiple GPSIMD Cores for Custom C++ Operators
- Package name was changed to aws-neuronx-gpsimd-customop-lib

#### aws-neuronx-gpsimd-customop [0.1]

Date: 02/08/2023

- First release of aws-neuronx-gpsimd-customop. This release provides tensor library support required for building Neuron Custom C++ operators.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### API Reference Guide

- *Custom Operators API Reference Guide [Beta]*

### Developer Guide

- *Neuron Custom C++ Operators Developer Guide [Beta]*

### Tutorials

- *Neuron Custom C++ Operators in MLP Training*
- *Neuron Custom C++ Operators Performance Optimization*

**Misc**

- *Neuron Custom C++ Tools Release Notes*
- *Neuron Custom C++ Library Release Notes*

*This document is relevant for:* **Inf2, Trn1, Trn2**





## LEARNING NEURON

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.1 Neuron Architecture

The Neuron Architecture provides insights into Neuron enabled instances system, software and chip capabilities. The Amazon EC2 Trn and Inf instance architecture provides an overview of the Amazon EC2 instances powered by AWS Inferentia and AWS Trainium chips (Neuron devices), and the corresponding system features like inbox and network connectivity, memory hierarchy, and NeuronCore versions and capabilities. The Neuron model architecture fit provides insights to what is the best match between deep-learning model architectures and the NeuronCore version.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### 8.1.1 Instance and UltraServer Architecture

For a detailed description of Trn Instances:

- Amazon EC2 Trn2 Architecture
- Amazon EC2 Trn1 Architecture

For a detailed description of Inf Instances:

- Amazon EC2 Inf2 Architecture
- Amazon EC2 Inf1 Architecture

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

#### 8.1.2 Amazon EC2 AI Chips Architecture

Amazon EC2 AI Chips (Neuron Devices) are the accelerated machine learning chips (e.g. Inferentia or Trainium) that enable Trn and Inf instance.

For a detailed description of current Neuron chips:

- trainium2-arch
- trainium-arch
- inferentia2-arch
- inferentia-arch

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.1.3 AWS NeuronCore Architecture

NeuronCores are fully-independent heterogenous compute-units that power Tranium, Tranium2, Inferentia, and Inferentia2 chips. For a detailed description of current generation NeuronCore (NeuronCore-v3) hardware engines, see:

- [neuroncores-v3-arch](#)

For more information about previous generation NeuronCores, see:

- [neuroncores-v2-arch](#)
- [neuroncores-v1-arch](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

- [Instance and UltraServer Architecture](#)
- [Amazon EC2 AI Chips Architecture](#)
- [AWS NeuronCore Architecture](#)

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 8.2 Neuron Features

Neuron features provide insights into Neuron capabilities that enable high-performance and improve usability of developing and deploying deep learning acceleration on top of Inferentia and Trainium based instances.

|                         |                                 |                    |                     |
|-------------------------|---------------------------------|--------------------|---------------------|
| Data types              | Neuron rounding modes           | Neuron batching    | NeuronCore pipeline |
| Neuron persistent cache | Neuron collective communication | Logical NeuronCore |                     |
| configuration           | Neuron custom C++ operators     |                    |                     |

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.2.1 Data Types

#### Table of contents

- [Introduction](#)
- [NeuronCore v1 Data Types](#)
  - [Neuron Data-Types](#)
- [NeuronCore v2 Data Types](#)
  - [Model Type Conversion](#)
  - [NeuronCore v2 Rounding Modes](#)

## Introduction

Inferentia and Trainium NeuronDevices include different NeuronCore versions, which support different data-types. This section describes what data-types are supported in each NeuronCore version, for details about NeuronCore versions see `neuron_hw_arch`.

## NeuronCore v1 Data Types

### Neuron Data-Types

Neuron enables developers to choose from multiple data-types. The supported data-types are FP32, FP16, and BF16. Developers can train their models on their platform of choice (e.g. EC2 P3 instances), and then easily move their trained models to EC2 Inf1 for execution.

### FP16/BF16 models

Models natively trained in FP16/BF16 will be executed in their trained data-types. This is a straightforward migration from the training platform to Inf1.

### FP32 models

Neuron SDK supports **automatic model conversion** from FP32 to BF16 by default. This capability allows developers to train their models using FP32 format for the highest accuracy, and achieve performance benefits without having to worry about low-precision training (e.g. no need for loss-scaling during training). ML models are typically robust to FP32 to BF16 conversion, with minimal to no impact on accuracy. The conversion accuracy is model dependent; therefore, users are encouraged to benchmark the accuracy of the auto-converted model against the original FP32 trained model.

When the compiler is supplied with an unmodified FP32 model input it will automatically compile the model to run as BF16 on Inferentia. During inference the FP32 input data will be auto-converted internally by Inferentia to BF16 and the output will be converted back to FP32 data-type. For explicit FP16 inferencing, either use an FP16 trained model, or use an external tool (like AMP) to make the explicit conversions.

## NeuronCore v2 Data Types

The NeuronCore v2 supports the following data types:

- 32 and 16-bit Floating Point (FP32 / FP16)
- TensorFloat-32 (TF32)
- Brain Floating Point (BFloat16)
- 8-bit Floating point with configurable range and precision (cFP8)
- Unsigned 8-bit integer (UINT8)

---

**Note:** Neuron Compiler support for cFP8 and UINT8 is planned for a future Neuron SDK release. For INT8, see [Neuron Compiler: Enable Neuron INT8 support](#) for details.

---

The layout for these is as follows:

## Model Type Conversion

The Neuron SDK supports automatic model conversion from FP32 to BF16 by default. This capability allows developers to train their models using FP32 format for the highest accuracy, and then achieve run-time performance benefits without having to worry about low-precision training (e.g. no need for loss-scaling during training). ML models are typically robust to FP32 to BF16 conversion, with minimal to no impact on accuracy. Since conversion accuracy is model dependent, users are encouraged to benchmark the accuracy of the auto-converted model against the original FP32 trained model.

See *Mixed Precision and Performance-accuracy Tuning for Training* for more details on supported data types and their properties.

The Neuron compiler offers the `--auto-cast` and `--auto-cast-type` options to specify automatic casting of FP32 tensors to other data types to address performance and accuracy tradeoffs. See the *Neuron Compiler CLI Reference Guide* for a description of these options.

## NeuronCore v2 Rounding Modes

Because floating point values are represented by a finite number of bits, they cannot represent all real numbers accurately. Floating point calculations that exceed their defined data type size are rounded. The NeuronCore v2 performs a Round-to-Nearest (RNE) algorithm with ties to Even by default. It also provides a new Stochastic Rounding mode. When Stochastic Rounding is enabled, the hardware will round the floating point value up or down using a proportional probability. This could lead to improved model convergence. Use the environment variable `NEURON_RT_STOCHASTIC_ROUNDING_EN` to select a rounding mode.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

### 8.2.2 Neuron Rounding Modes

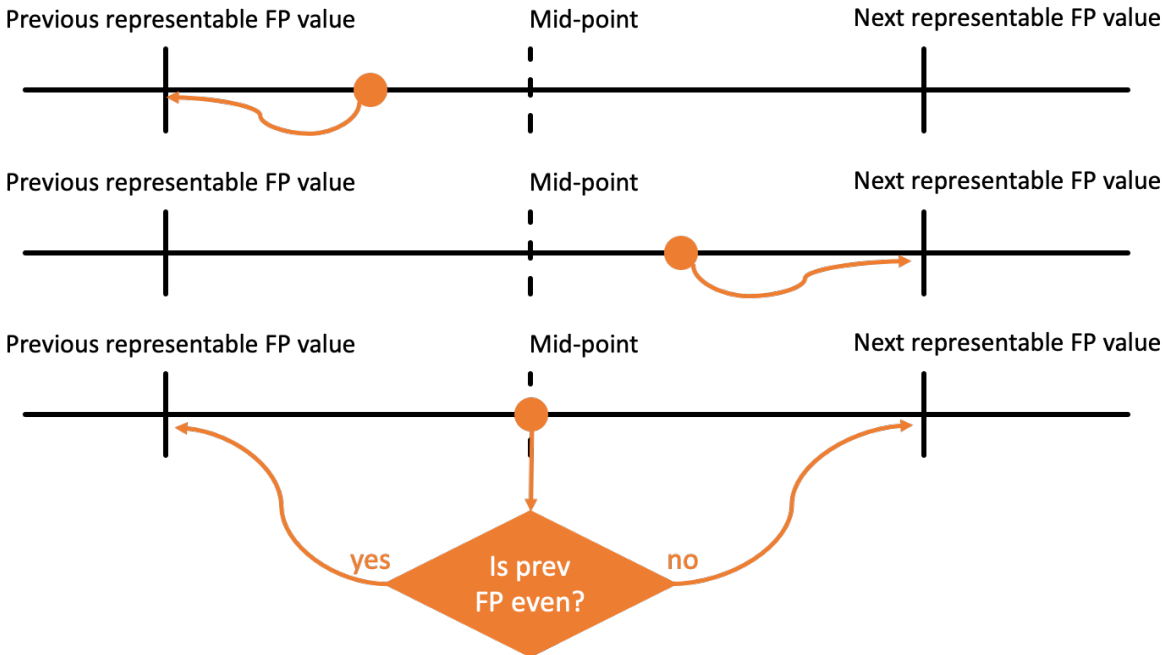
#### Table of contents

- *Round Nearest, ties to Even (RNE)*
- *Stochastic Rounding (SR)*
- *Quick Tests*

#### Round Nearest, ties to Even (RNE)

When the exact result of a floating point operation cannot be exactly represented as a floating point value, it must be rounded. The IEEE 754-2008 standard defines the default rounding mode to be ‘Round Nearest, ties to Even’ (RNE for short). Under this scheme, numbers are rounded to the nearest representable value, and in case of a ‘tie’ (i.e. the number is exactly between the two nearest representable values) numbers will be rounded to the nearest even number.

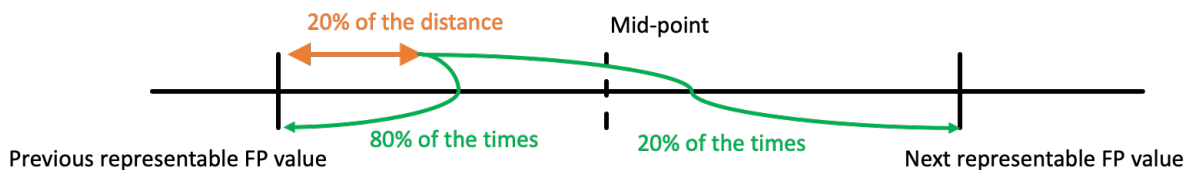
All NeuronCore generations support the RNE rounding scheme, which is the most commonly used rounding scheme for Machine Learning workloads. Below is an illustration of the RNE rounding scheme:



### Stochastic Rounding (SR)

One downside of the RNE rounding scheme (and other rounding schemes described in the IEEE 754-2008 standard), is that when adding floating point values of significantly different magnitudes, rounding can squash small values and prevent them from accumulating over time.

To improve this, starting from the second generation of the NeuronCore (NeuronCore-v2), customers can choose between the RNE rounding scheme described above, and a second rounding scheme called 'Stochastic Rounding' (SR for short). Stochastic rounding prevents the computation precision-loss described above, by performing the rounding operations in a probabilistic manner, according to the relative distance from the two nearest representable values, as illustrated below:



By performing the rounding in a probabilistic manner, this scheme allows for small increments to accumulate over time, even when added to numbers of significantly higher magnitude, which leads to more precise results when performing large floating point computations (as done for machine learning).

## Quick Tests

As an example, we examine the code-snippet below:

```
import torch
import torch_xla
import torch_xla.core.xla_model as xm
device = xm.xla_device()

a = torch.tensor(1024.0).half().to(device)

for i in range(2048) :
    a = (a + 0.5)
    xm.mark_step()

print(a)
```

This code shows that rounding can significantly impact the calculation's precision over time. To use standard RNE rounding, use the environment variable `NEURON_RT_STOCHASTIC_ROUNDING_EN=0`. To enable stochastic rounding, use the environment variable `NEURON_RT_STOCHASTIC_ROUNDING_EN=1`.

NOTE: Stochastic rounding mode is enabled by default in PyTorch-Neuron when `XLA_USE_BF16=1`.

The first test continues to show 1024 due to RNE rounding after each addition, and the second test shows result that is mostly in line with expectation.

```
$ NEURON_RT_STOCHASTIC_ROUNDING_EN=0 python3 rounding_mode_test.py

tensor(1024., device='xla:1', dtype=torch.float16)

$ NEURON_RT_STOCHASTIC_ROUNDING_EN=1 python3 rounding_mode_test.py

tensor(2056., device='xla:1', dtype=torch.float16)
```

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.2.3 Neuron Batching

Batching refers to the process of grouping multiple samples together, and processing them as a group (i.e. passing them together through the neural network). Batching is typically used as an optimization for improving throughput at the expense of higher latency (and potentially higher memory footprint). Batching considerations are slightly different between inference and training workloads, and we thus cover them separately below.

#### Table of contents

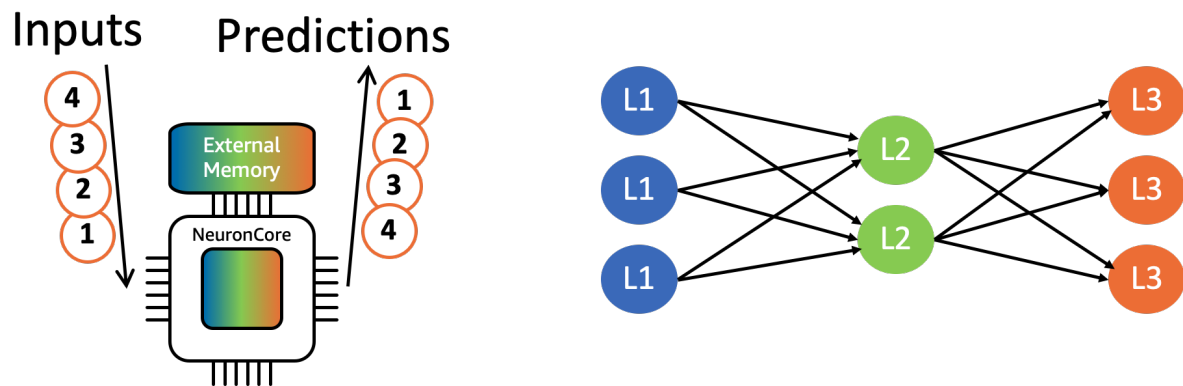
- *Batching in inference workloads*
  - *What is batched inference?*
  - *What are the benefits of batched Inference?*
  - *How to determine the optimal batch-size for inference workloads?*
  - *How to set the batch-size?*

- *Batching in training workloads*
  - *Batch-size naming*
  - *How to determine the optimal batch-size for training workloads?*

## Batching in inference workloads

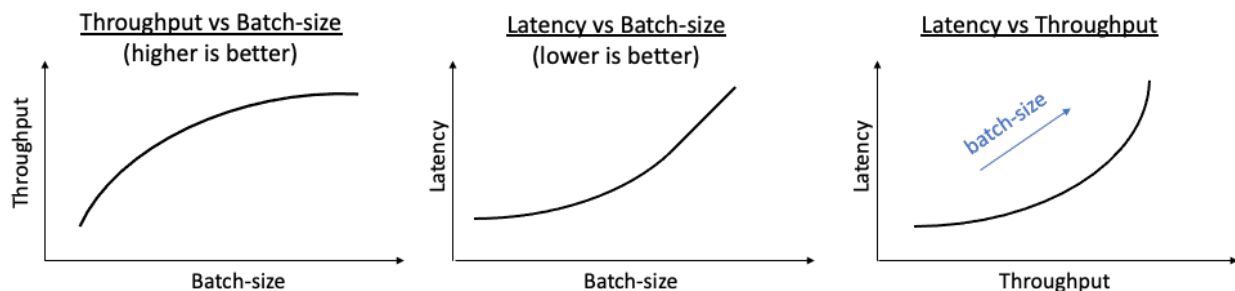
### What is batched inference?

The concept of batched inference is conceptually illustrated below, with a single NeuronCore performing batched computation of a 3 layer neural network with a batch-size of 4. The NeuronCore reads the parameters for a certain layer from the external memory, and then performs the corresponding computations for all 4 inference-requests, before reading the next set of parameters (thus, performing more compute for every parameter read from memory).

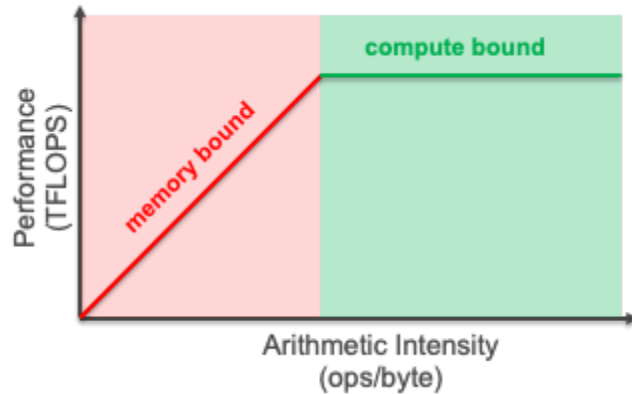


### What are the benefits of batched Inference?

For inference, batching is typically used as a trade-off knob between throughput and latency: higher batch-size typically leads to better hardware utilization and thus higher throughput, but at the same time batching requires to perform more computation until getting the first results, and hence leads to higher latency.



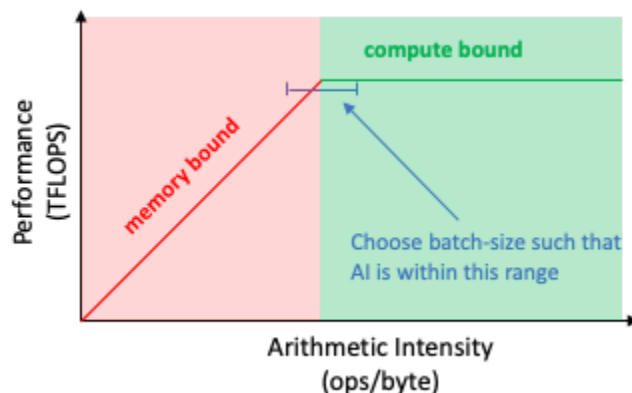
To understand why batching tends to improve throughput (up to a certain max value), it is useful to consider an intuitive visual performance-model called 'the roofline model', which provides with a theoretical bound on the system's performance:



The X-axis indicates the arithmetic intensity (AI) of the workload, which is the ratio between the number of operations and the number of bytes read-from/written-to memory. The Y-axis indicates the theoretical extractable performance. For small(large) AI values, the workload is expected to be memory(compute) bound. For inference workloads, AI is often approximated by dividing the model's number of operations by its memory footprint ( $\#params \times dtype\_size$ ). To a first order approximate, the AI value is linearly dependent on the batch-size, which means that the workloads performance (throughput) is expected to increase with the batch-size. To understand this more intuitively, for a larger batch size, Neuron can better amortize the cost of reading parameters from the external memory, and thus improve the overall hardware efficiency. It should be noted that while the roofline model can be very useful, it is not perfectly accurate (e.g. it doesn't take into account spill/fills from/to on-chip SRAM memories), and thus users are encouraged to use it as a tool for **estimating** the optimal batch-size for their workloads.

### How to determine the optimal batch-size for inference workloads?

The optimal batch size is dependent on the application-level requirements: some applications require strict latency guarantees (in which case, check out the [NeuronCore Pipeline](#) technology), while other applications strictly aim to maximize throughput. We thus encourage our users to try out multiple batch-sizes, and compare performance between them. A good starting for batch-size exploration can be identified using the roofline model: we can choose a batch-size that achieves an Arithmetic Intensity which is at the edge of the compute bound region. By doing that, we aim to achieve max throughput with a minimal batch-size, and thus minimal impact to latency.



This can be expressed via the following equation:  $batch\_size(Inference) = \lceil 0.5 \times (\frac{NeuronDevice\ PeakFLOPs}{NeuronDevice\ MemBW}) / (\frac{model\ FLOPs}{\#model\_dense\_params \times dtype\_size}) \rceil$  (for NeuronDevice PeakFLOPs and MemBW, see the [trainium-arch](#), [inferentia-arch](#) and [inferentia2-arch](#) pages).

For example, a BF16 BERT-Large model, with a sequence length of 128, will have the following approximated batch sizes:



| Model                   | Neuron-Device | Peak TFLOPS (BF16) | MemBW (GB/sec) | Model GFLOPs | Model Dense Params (Millions) | Data-type size (BF16) | Approximated optimal batch-size |
|-------------------------|---------------|--------------------|----------------|--------------|-------------------------------|-----------------------|---------------------------------|
| BERT-Large (SeqLen=128) | Inferentia    | 64                 | 50             | 77.3         | 302                           | 2                     | 6                               |
| BERT-Large (SeqLen=128) | Trainium      | 210                | 820            | 77.3         | 302                           | 2                     | 2                               |
| ResNet-50               | Inferentia    | 64                 | 50             | 7.8          | 25                            | 2                     | 5                               |
| ResNet-50               | Trainium      | 210                | 820            | 7.8          | 25                            | 2                     | 1                               |

We recommend to evaluate multiple batch sizes and compare the performance between them, in order to determine the optimal latency/throughput deployment-point.

### How to set the batch-size?

The Neuron compiler takes a model and its sample input, as inputs for the compilation process. For example, the code snippet below will compile a model with a batch-size of 4:

```
import torch
import torch_neuron
from torchvision import models

# Load the model and set it to evaluation mode
model = models.resnet50(pretrained=True)
model.eval()

# Compile with an example input of batch size 4
image = torch.rand([4, 3, 224, 224])

model_neuron = torch.neuron.trace(model, image, dynamic_batch_size=True)

# Execute with a batch of 12 images
batch = torch.rand([12, 3, 224, 224])
results = model_neuron(batch)
```

For ahead-of-time compiled inference graphs (i.e. Inf1), dynamic batching can be used (as shown in the above code snippet) to process a larger client-side inference batch-size, and allow the framework to automatically break up the user-batch (12 in our case) into smaller batch sizes, to match the compiled batch-size (4 in our case). This technique increases the achievable throughput by hiding the framework-to-neuron overhead, and amortizing it over a larger batch size.

See also:

- *Dynamic Batching* in torch-neuronx
- *Special Flags* in tensorflow-neuronx.

## Batching in training workloads

Unlike inference workloads, training is inherently an offline process, and thus doesn't have latency requirements. This means that training is almost always batched to some degree.

### Batch-size naming

For distributed processing, defining the batch size depends on the observation level. There are multiple terms you should be aware of when running a distributed training job, especially global batch size (GBS) and micro-batch. Knowing the batch size in advance is crucial for precompiling the computational graph and for setting the hyperparameters.

#### **micro-batch size**

Smallest unit of the number of samples getting processed in a single step in the accelerator. For very large models, it is frequently chosen to be 1.

#### **gradient accumulation**

Process of iterating over a micro-batch multiple times and summing up the gradients before an optimizer update. This can happen in a dedicated loop for gradient accumulation or as part of multiple iterations of samples in pipeline parallelism. See *Developer guide for Pipeline Parallelism* for more details on pipeline parallelism.

#### **data-parallel size (or DP degree)**

Number of model replicas that process different portions of data in parallel. Each replica maintains a complete copy of the model while processing unique data chunks, after which their gradients are synchronized for the optimizer update. See *Neuron Glossary* for more details.

#### **global batch-size**

Number of total samples used for an update of the optimizer. This includes all the respective gradients that get added up from data-parallel processing or gradient accumulation.  $\text{global\_batch\_size} = \text{micro\_batch\_size} * \text{data\_parallel\_size} * \text{gradient\_accumulation\_steps}$

#### **mini-batch or replica-batch size**

Number of samples that contribute to a gradient within one data-parallel rank. A mini-batch gradient is obtained by aggregating multiple micro-batch gradients within or without a pipeline (aka. gradient accumulation).  $\text{mini\_batch\_size} = \text{micro\_batch\_size} * \text{gradient\_accumulation\_steps}$

#### **worker batch**

The portion of mini-batch samples processed by a worker. The idea behind a worker batch is that one worker (node) might have a subset of the dp-degrees and we care about how much data gets tackled by this worker.

## How to determine the optimal batch-size for training workloads?

Determining the optimal batch-size for training workloads can be a non-trivial task. In most cases, we'd want to choose the largest batch-size that we can get away with.

The most dominant factor for determining the optimal batch-size in training workloads is memory footprint: training workloads have higher memory footprint compared to inference, as they require saving more tensors aside from the model parameters, such as gradients, intermediate activations (passed between forward-pass and backward-pass), and optimizer-state. If the batch-size is increased beyond a certain point, one can run out of device memory (indicated by an 'Out of device memory' error, typically abbreviated as OOM).

To estimate the memory footprint of a model, we look at the different contributors:

1. Weights and gradients:

1. typically 2B each, thus 4B per parameter
2. Optimizer state:
  1. typically 4B - 12B per parameter
3. Intermediate activations:
  1. sum of all tensor sizes for forward pass
  2. for example, for a transformer neural network, this is roughly  $16 \times x \times \text{<num\_layers>} \times x \times x = 100\text{MB} \times$

For training workloads, determining the optimal batch size can be a little more tricky, due to two reasons:

1. *Higher memory footprint:* Training workloads have higher memory footprint compared to inference, as they require saving more tensors aside from the model parameters, such as gradients, intermediate-state and optimizer-state. If the batch-size is increased too much, one can run out of device memory (indicated by an 'Out of memory' error, typically abbreviated as OOM).
2. *Arithmetic intensity estimation:* Arithmetic intensity is harder to estimate in training workloads, compared to inference workloads, as the majority of the external memory access are due to reads/writes of intermediate activation state (rather than parameters), which requires lower level familiarity with the model to estimate correctly.

A good first order approximate for the optimal batch-size in a training workload, is the largest one that can fit in the device's memory (i.e. won't lead to OOM error).  $\text{batch-size(Training)} = 0.6 \times (\text{<TP-Rank>} \times \text{<PP-Rank>} \times \text{``<NeuronCore MemoryCapacity>} / \text{``(<\#model-dense-params>} \times \text{``<model-state-bytes-per-parameter>)})$

Note TP-rank stands for Tensor-Parallelism rank, i.e. how many NeuronCores participate in a single Tensor-Parallelism group. Similarly, PP-rank stands for Pipeline-Parallelism rank, i.e. how many NeuronCores participate in a single Pipeline-Parallelism group.

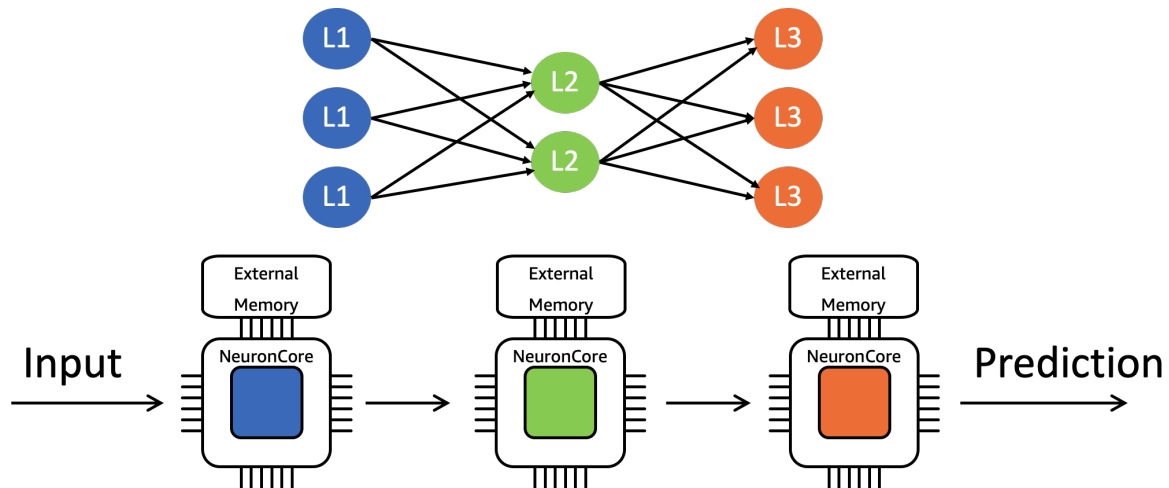
For example, for BERT-Large Ph1 training, with a model-state of 4B per parameter (2B weights, 2B parameters), and TP-rank = PP-rank = 1, the approximated optimal per-NeuronCore training batch-size would be:  $\text{batch-size(Training/Trainium)} = 0.6 \times (1 \times 1 \times 16\text{e}+9) / (300\text{e}+6 \times 4) = 8$

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1

## 8.2.4 NeuronCore Pipeline

The Neuron software feature referred to as a NeuronCore Pipeline refers to the process of sharding a compute-graph across multiple NeuronCores, caching the model parameters in each core's on-chip memory (cache), and then streaming inference requests across the cores in a pipelined manner. Based on the number of NeuronCores selected, the model might get seamlessly sharded across up-to 16 Inferentia devices (i.e. 64 NeuronCores). This enables users to optimize for both throughput and latency, as it enables the NeuronCores to process neural-networks with locally cached data and avoid the cost of accessing external memory.



One benefit to this approach is that NeuronCore Pipeline can typically hit maximal hardware efficiency without the need for batching (e.g. BERT, ResNet50).

For maximal performance, users should choose an instance-size that can cache the entire model by using sufficient NeuronCores. Inf1 instance types have different number of Inferentia devices, each of which has 4 NeuronCores, as shown here <https://aws.amazon.com/ec2/instance-types/inf1/>

To enable the NeuronCore Pipeline optimization, the compiler should be invoked with the following flags: `--neuroncore-pipeline-cores N`. The number of NeuronCores is typically chosen to be the minimal number that can fit the entire model, which is currently done through a trial-and-error process (compiling to different number of cores and looking for compilation success/failure message). This process will be automated in the future. A simple formula to help define the number of NeuronCores that may be an appropriate choice is

```
neuroncore-pipeline-cores = 4 * round( number-of-weights-in-model / (2 * 10^7) )
```

This allocates a set of NeuronCores based on the size of the given model's weights and normalizes to multiples of 4 so it uses full Inferentias.

The code snippet below shows how to compile a model with NeuronCore Pipeline for 16 NeuronCores (instance size inf1.6xlarge).

```
import numpy as np
import tensorflow.neuron as tfn

example_input = np.zeros([1,224,224,3], dtype='float16')
tfn.saved_model.compile("rn50_fp16",
                        "rn50_fp16_compiled/1",
                        model_feed_dict={'input_1:0' : example_input },
                        compiler_args = ['--neuroncore-pipeline-cores', '16'])
```

*This document is relevant for:* Inf1

*This document is relevant for:* Inf2, Trn1, Trn2

## 8.2.5 Neuron Collective Communication

### Table of contents

- *Introduction*
- *trn1.32xlarge topology*
- *trn1.2xlarge topology*
- *inf2.48xlarge topology*
- *Inf2 other instance sizes topologies*

### Introduction

Collective Communications is an integral component of distributed ML training. Multiple training nodes exchange information during ML training via Collective Communication operators such as all-reduce. Neuron provides hardware support for the execution of Collective Communication with the Neuron SDK responsible for the hardware configuration and for the execution orchestration. Neuron provides the following Collective Communication operators:

- all-reduce
- all-gather
- reduce-scatter

Neuron also provides the following peer to peer operators:

- send
- receive

Support for additional Collective Communication operators might be added in future releases. Neuron devices are connected via NeuronLinks within a single instance and EFA links between instances. All NeuronLinks transfer the data directly between Neuron device and between Neuron devices and EFA devices bypassing the host to achieve high bandwidth and low latency.

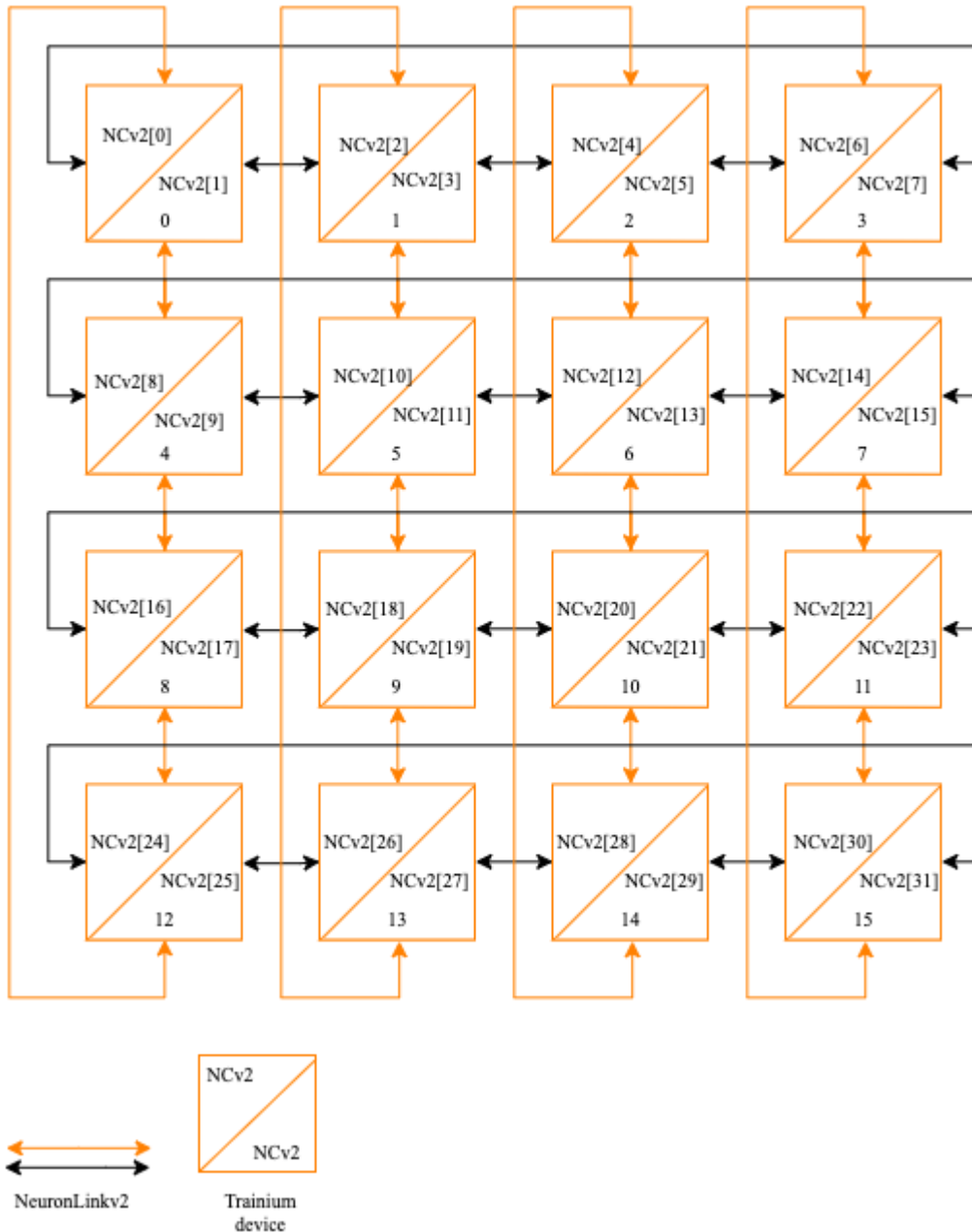
Collective Communication support on Neuron requires installation of 3 separate packages:

- `aws-neuronx-runtime-lib` - supports execution on Neuron, not specific to Collective Communication and is always required
- `aws-neuronx-collectives` - supports Collective Communication execution on a single instance and on multiple instances.
- `efa_installer` - low level libraries and drivers to support Collective Communication execution over EFA, required to support Collective Communication on multiple instances.

ML models need to be compiled by the Neuron compiler before they can be executed on Neuron devices. The result of the compilation is a binary object containing computational instruction and data movement instructions. Any Collective Communication operators encountered during compilation are converted to the place holder instructions to be filled by the runtime/collectives libraries during load and execution. This approach allows Neuron compiler to be unaware of the specific physical topology connecting Neuron devices. Once a compiled mode is placed on Neuron devices the runtime/collectives libraries generate the appropriate data movement instructions based on the placement. For example, a different set of instructions is generated when the next rank is connected via NeuronLinks or via EFA. Neuron executes Collective Communication operators using dedicated hardware that is not shared with computational resources. That allows Neuron to execute compute and communication in parallel. For example Neuron can all-reduce gradients of one

layer while the gradients for another layer are computed. Overlapping compute and communication can result in lower latency and higher performance.

### trn1.32xlarge topology



### Trn1.32xl 2D torus topology

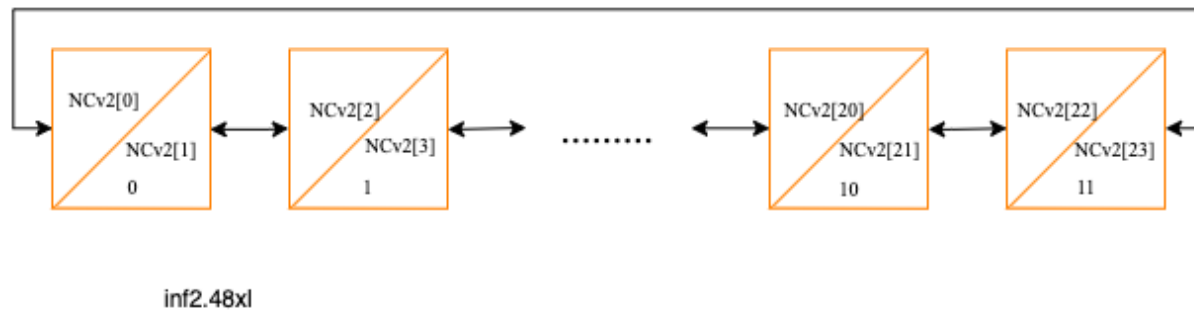
On a single trn1.32xlarge instance Neuron devices are connected in a 2D torus topology supporting Collective Communication operators in sets of 2, 8 and 32 ranks. Other set sizes might be supported in future releases. A single instance topology can be further extended across multiple instances using EFA NeuronLinks.

For example an 8x4 topology on a single instance, such as 8 rank tensor parallel and 4 ranks data parallel can be extended across multiple instances creating a large tensor/data parallel training cluster.

### trn1.2xlarge topology

Trn1.2xlarge instance type contains a single Neuron device with two NeuronCores. This instance type supports only 2 ranks Collective Communication operators. EFA is not available on trn1.2xlarge and the ranks cannot be extended beyond a single instance.

### inf2.48xlarge topology



### inf2.48xlarge topology

On inf2.48xlarge instance Neuron devices are connected in a ring via NeuronLink. Any **even** number of ranks for Collective Communication operators is supported provided that the ranks occupy consecutive Neuron devices. However, when using any number of ranks other than 24 (full instance) full performance of the ring is not utilized.

### Inf2 other instance sizes topologies



### inf2 other instance sizes topologies

On other inf2 instance sizes Neuron devices are connected bi-directionally. Any **even** number of ranks for Collective Communication operators is supported provided that the ranks occupy consecutive Neuron devices. Collective Communication performance is similar to the performance on inf2.48xlarge when fewer than 24 ranks are used.

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 8.2.6 Logical NeuronCore configuration

Logical NeuronCore configuration (LNC) is a set of compiler and runtime settings for instances powered by AWS Trainium2 that determines the number of NeuronCores exposed to your machine learning (ML) applications. LNC configuration works by combining the compute and memory resources of multiple physical NeuronCores into a single logical NeuronCore. You can configure these settings to reduce the number of worker process needed for training and deployment of large-scale models.

### Concepts

- *Logical NeuronCores*
- *Compiler and runtime settings*
- *Logical NeuronCore configurations*

## Logical NeuronCores

A logical NeuronCore is a grouping of physical NeuronCores that the Neuron Compiler, Neuron Runtime, Neuron Tools, and Frameworks handle as a single unified NeuronCore. Every Trainium2 device contains eight physical NeuronCore-v3.

## Compiler and runtime settings

LNC configuration is controlled with the following runtime and compiler settings:

### Neuron Runtime

The `NEURON_LOGICAL_NC_CONFIG` runtime environment variable controls how many physical NeuronCores are grouped to make up a logical NeuronCore.

### Neuron compiler flags

The `--logical-nc-config` or `-lnc` command-line options control the degree of model sharding the compiler performs on an input graph. You must compile your Models to use the LNC configuration set by the Neuron Runtime environment variable. AWS Neuron currently doesn't support setting the compiler flag to a different LNC configuration than the Neuron Runtime environment variable.

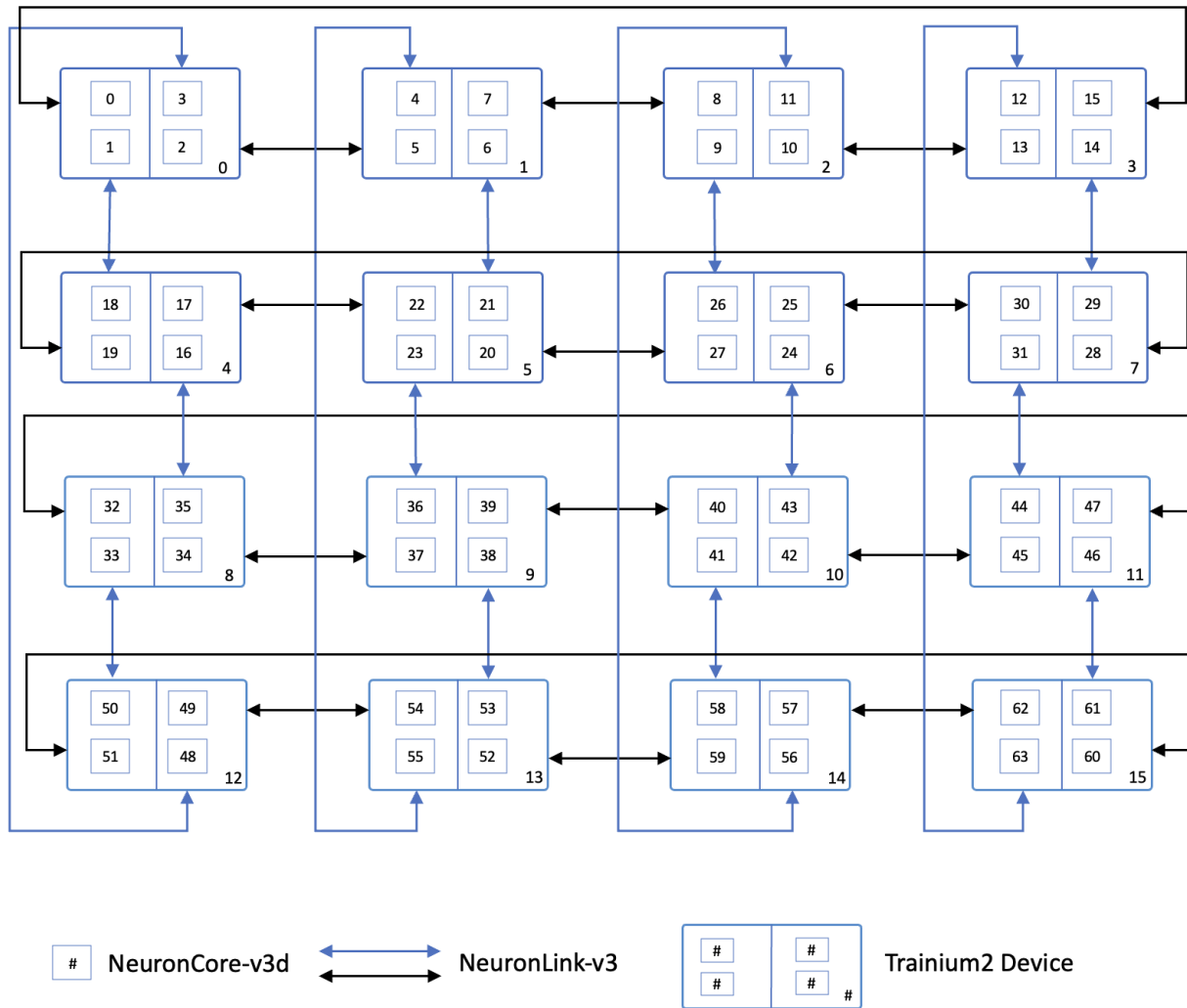
## Logical NeuronCore configurations

AWS Neuron supports the following Logical NeuronCore configurations:

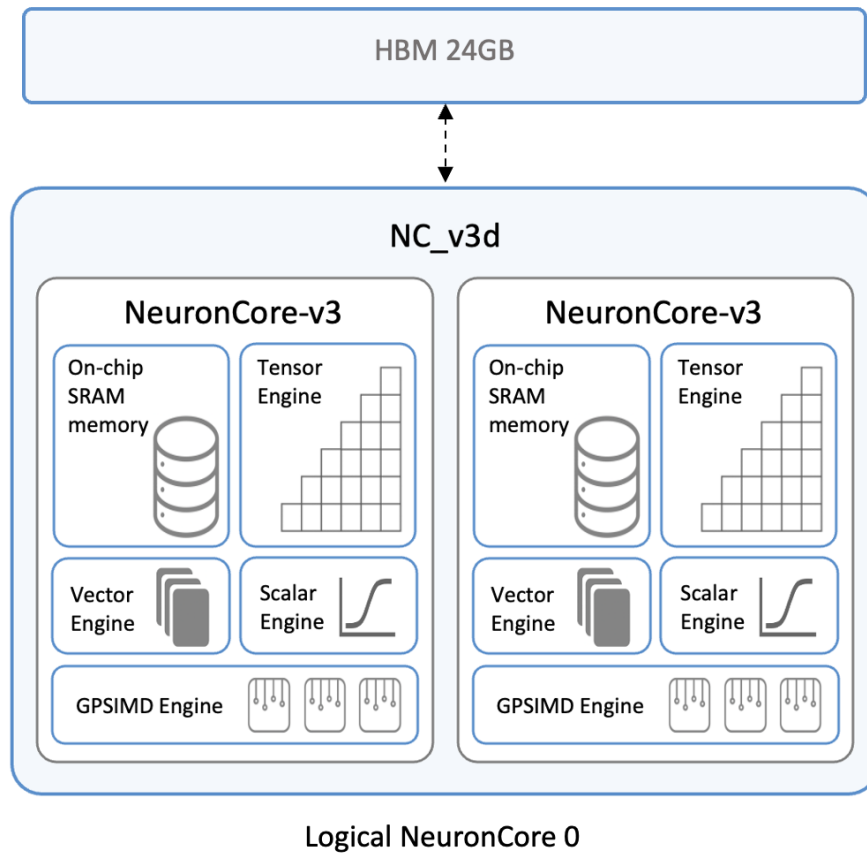
### LNC = 2

A Logical NeuronCore configuration (LNC) of two is the default setting on Trainium2 devices. It combines two physical NeuronCore-v3 into a logical NeuronCore with the software id `NC_v3d`. When you set Logical NeuronCore configuration to two, it directs Trainium2 devices to expose four `NC_v3d` to your machine learning applications. On this setting, a `Trn2.48xlarge` instance presents 64 available NeuronCores. The following high-level diagram shows a `Trn2.48xlarge` instance, connected in a 2D torus topology, with the Logical NeuronCore configuration set to two.





Trainium2 devices contain four 24GB HBM banks. Each bank is shared by two physical NeuronCore-v3. When LNC=2, the two physical NeuronCores share a single address space. Workers on each of the two physical NeuronCores can access tensors and perform local collective operations without accessing the network. The following diagram shows how a logical NeuronCore is presented to the software under this configuration.



To set the Logical NeuronCore configuration to two, use the following runtime and compiler flag combination:

**Runtime environment variable:**

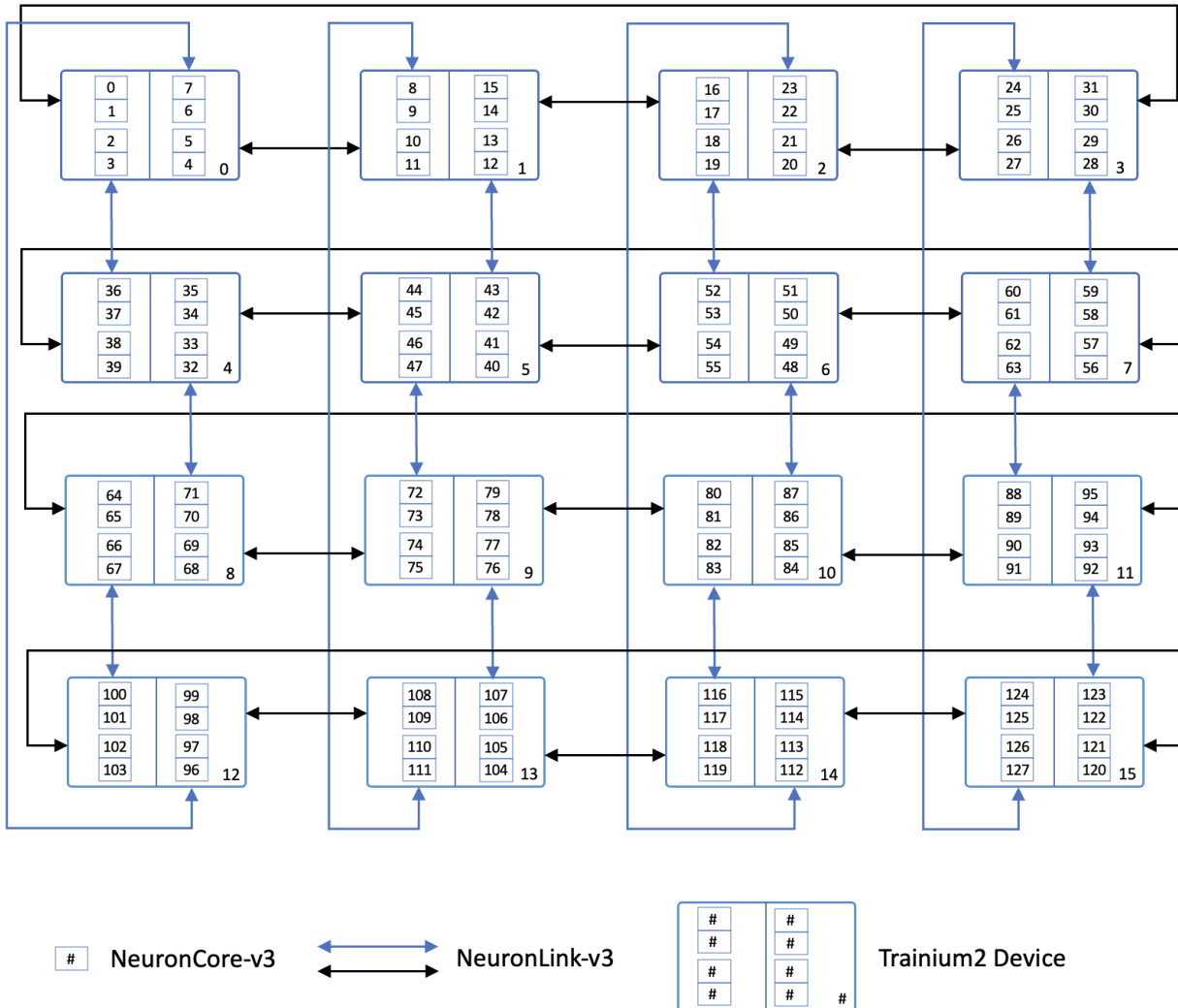
`NEURON_LOGICAL_NC_CONFIG = 2`

**Compiler flag:**

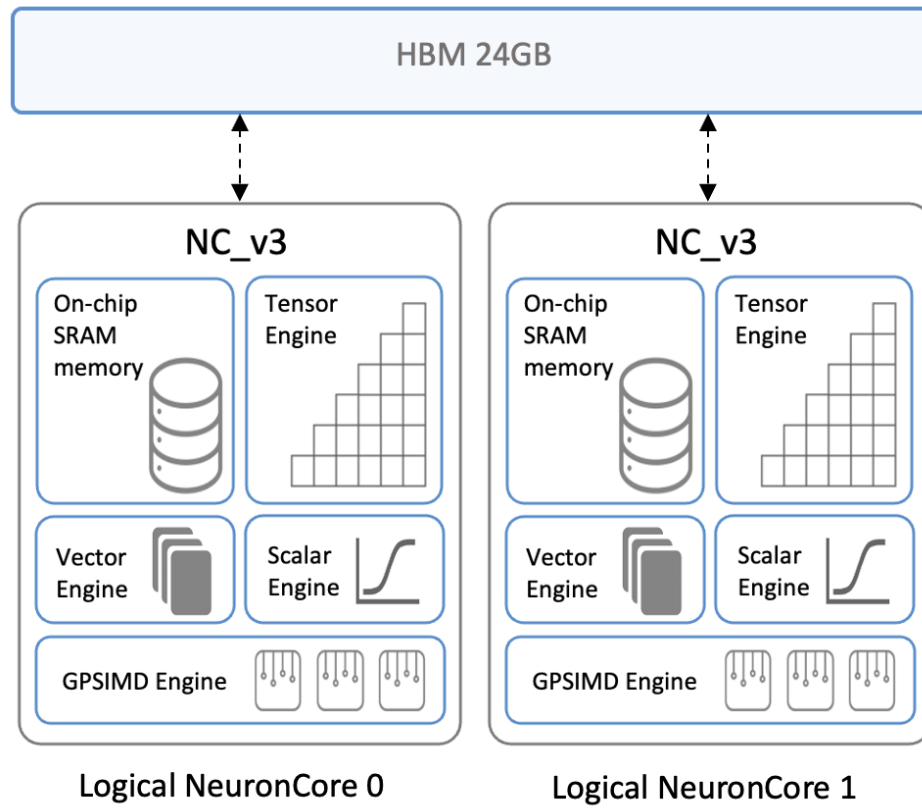
`-lnc = 2`

**LNC = 1**

When you set the Logical NeuronCore configuration to one, it assigns each physical NeuronCore-v3 to a single logical NeuronCore with the software id NC\_V3. This directs Trainium2 devices to expose eight NC\_v3 to your machine learning applications. On this setting, a Trn2.48xlarge instance presents 128 available NeuronCores. The following high-level diagram shows a Trn2.48xlarge instance, connected in a 2D torus topology, with the Logical NeuronCore configuration set to one.



Trainium2 devices contain four 24GB HBM banks. Each bank is shared by two physical NeuronCore-v3. When the Logical NeuronCore configuration is set to one, both physical NeuronCores have access to the entire 24GB HBM bank. The following diagram shows how logical NeuronCores are presented to the software under this configuration.



To set the Logical NeuronCore configuration to one, use the following runtime and compiler flag combination:

**Runtime environment variable:**

`NEURON_LOGICAL_NC_CONFIG = 1`

**Compiler flag:**

`-lnc = 1`

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf2, Trn1, Trn2*

## 8.2.7 Neuron Custom C++ Operators

Neuron Custom C++ Operators enable developers to write C++ Custom Operators (“CustomOps”) that run on NeuronCores. This enables developers to extend operator support beyond what is officially supported by Neuron.

Developers can use standard PyTorch custom operators programming interfaces to leverage Neuron Custom C++ Operators feature. This makes it easy to migrate CPU Custom Operators to Neuron, and implement new beta operators, all without any intimate knowledge of the NeuronCore hardware.

For more details see [Neuron Custom C++ Operators \[Beta\]](#)

*This document is relevant for: Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 8.3 Neuron Application Notes

### Neuron 2.x

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.3.1 Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability (GA)

Neuron release 2.3 is the first release of Neuron 2.x that enables GA of the new EC2 Trn1 instances. Neuron release 2.3 extends the latest release of Neuron 1.x (Neuron 1.19.2), adding support for Deep Learning training on the AWS Trainium chips.

Starting with Neuron release 2.3, developers can run Deep Learning training workloads on Trn1 instances, saving training costs by up to 50% over equivalent GPU-based EC2 instances, while achieving the highest training performance in the AWS cloud for popular NLP models. Neuron 2.x introduces new capabilities and major architectural updates to support training neural-networks with the Trn1 instances.

In addition, starting with this release, Neuron introduces new packages, renames several packages, and updates Neuron installation and update instructions. This release also ends support for Neuron Runtime 1.x.

## More about the release

|                                   |   |
|-----------------------------------|---|
| What's New                        | <ul style="list-style-type: none"> <li>• <a href="#">rn2.3.0_new</a></li> <li>• <a href="#">neuron-packages-changes</a></li> <li>• <a href="#">announce-aws-neuron-github-org</a></li> <li>• <a href="#">announce-neuron-rtd-eol</a></li> </ul> |
| Get started with Neuron           | <ul style="list-style-type: none"> <li>• <a href="#">torch_quick_start</a></li> <li>• <a href="#">Neuron Quick Links</a></li> </ul>   |
| Tested workloads and known issues | <ul style="list-style-type: none"> <li>• <a href="#">rn2.3.0_tested</a></li> <li>• <a href="#">rn2.3.0-known-issues</a></li> </ul>  |
| Frequently Asked Questions (FAQ)  | <ul style="list-style-type: none"> <li>• <a href="#">neuron2-intro-faq</a></li> <li>• <a href="#">neuron-training-faq</a></li> </ul>  |
| Troubleshooting                   | <ul style="list-style-type: none"> <li>• <a href="#">PyTorch Neuron Troubleshooting on Trn1</a></li> <li>• <a href="#">Neuron Runtime Troubleshooting on Inf1 and Trn1</a></li> </ul>   |
| Neuron architecture and features  | <ul style="list-style-type: none"> <li>• <a href="#">Neuron Architecture</a></li> <li>• <a href="#">Neuron Features</a></li> </ul>  |
| Neuron Components release notes   | <ul style="list-style-type: none"> <li>• <a href="#">components-rn</a></li> </ul>   |

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## Neuron Runtime library

*This document is relevant for: Inf1*

### 8.3.2 Introducing Neuron Runtime 2.x (libnrt.so)

#### Table of contents

- [What are we changing?](#)
- [Why are we making this change?](#)
- [How will this change affect the Neuron SDK?](#)
  - [Neuron Driver](#)

- *Neuron Runtime*
- *Neuron framework extensions*
- *TensorFlow model server*
- *Neuron tools*
- *How will this change affect me?*
  - *Neuron installation and upgrade*
  - *Migrate your application to Neuron Runtime 2.x (libnrt.so)*
- *Troubleshooting*
  - *Application fails to start*
  - *Application fails to start although I installed latest aws-neuron-dkms*
  - *Application unexpected behavior when upgrading to release Neuron 1.16.0 or newer*
  - *Application unexpected behavior when downgrading to releases before Neuron 1.6.0 (from Neuron 1.16.0 or newer)*
  - *Neuron Core is in use*
- *Frequently Asked Questions (FAQ)*
  - *Do I need to recompile my model to run it with Neuron Runtime 2.x (libnrt.so)?*
  - *Do I need to change my application launch command?*
  - *Can libnrt.so and neuron-rtd co-exist in the same environment?*
  - *Are there Neuron framework versions that will not support Neuron Runtime 2.x (libnrt.so)?*

## What are we changing?

Starting with the *Neuron 1.16.0* release, *Neuron Runtime 1.x* (*neuron-rtd*) is entering maintenance mode and is being replaced by *Neuron Runtime 2.x*, a shared library named (*libnrt.so*). For more information on Runtime 1.x see *maintenance\_rtd*.

Upgrading to *libnrt.so* simplifies the Neuron installation and upgrade process, introduces new capabilities for allocating NeuronCores to applications, streamlines container creation, and deprecates tools that are no longer needed.

This document describes the capabilities of *Neuron Runtime 2.x* in detail, provides information needed for successful installation and upgrade, and provides information needed for successful upgrade of Neuron applications using *Neuron Runtime 1.x* (included in releases before *Neuron 1.16.0*) to *Neuron Runtime 2.x* (included in releases *Neuron 1.16.0* or newer).

## Why are we making this change?

Before *Neuron 1.16.0*, Neuron Runtime was delivered as a daemon (*neuron-rtd*), and communicated with Neuron framework extensions through a gRPC interface. *neuron-rtd* was packaged as an rpm or debian package (*aws-neuron-runtime*) and required a separate installation step.

Starting with *Neuron 1.16.0*, *Neuron Runtime 2.x* is delivered as a shared library (*libnrt.so*) and is directly linked to Neuron framework extensions. *libnrt.so* is packaged and installed as part of the Neuron framework extensions (e.g. TensorFlow Neuron, PyTorch Neuron or MXNet Neuron), and does not require a separate installation step. Installing Neuron Runtime as part of the Neuron framework extensions simplifies installation and improves the user experience.

In addition, since `libnrt.so` is directly linked to the Neuron framework extensions, faster communication between the Neuron Runtime and Neuron Frameworks is enabled by eliminating the gRPC interface overhead.

For more information see [How will this change affect the Neuron SDK?](#) and [Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#).

### How will this change affect the Neuron SDK?

#### Neuron Driver

Use the latest Neuron Driver. For successful installation and upgrade to *Neuron 1.16.0* or newer, you must install or upgrade to Neuron Driver (`aws-neuron-dkms`) version *2.1.5.0* or newer. Neuron applications using *Neuron 1.16.0* will fail if they do not detect *Neuron Driver version 2.1.5.0* or newer. For installation and upgrade instructions see `install-guide-index`.

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

To see details of Neuron component versions please see `neuron-release-content`.

---

#### Important:

**For successful installation or update to Neuron 1.16.0 and newer from previous releases:**

- Stop Neuron Runtime 1.x daemon (`neuron-rtd`) by running: `sudo systemctl stop neuron-rtd`
  - Uninstall `neuron-rtd` by running: `sudo apt remove aws-neuron-runtime` or `sudo yum remove aws-neuron-runtime`
  - Install or upgrade to the latest Neuron Driver (`aws-neuron-dkms`) by following the `install-guide-index` instructions.
  - Starting with Neuron version 2.3, `aws-neuron-dkms` the package name has been changed to `aws-neuronx-dkms`, see [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)
- 

#### Neuron Runtime

- **Installation** Starting from *Neuron 1.16.0*, Neuron releases will no longer include the `aws-neuron-runtime` packages and Neuron Runtime will be part of the Neuron framework extension of choice (TensorFlow Neuron, PyTorch Neuron or MXNet Neuron). Installing any Neuron framework package will install the Neuron Runtime library (`libnrt.so`).
  - For installation and upgrade instructions see `install-guide-index`.
- **Configuring Neuron Runtime**

Before *Neuron 1.16.0*, *Neuron Runtime 1.x* was configured in configuration files (e.g. `/opt/aws/neuron/config/neuron-rtd.config`). Starting from *Neuron 1.16.0*, *Neuron Runtime 2.x* can be configured through environment variables. See [NeuronX Runtime Configuration](#) for details.
- **Starting and Stopping Neuron Runtime**

Before introducing `libnrt.so`, `neuron-rtd` ran as a daemon that communicated through a gRPC interface. Whenever `neuron-rtd` took ownership of a Neuron device, it continued owning that device until it was stopped. This created the need to stop `neuron-rtd` in certain cases. With the introduction of



`libnrt.so`, *Neuron Runtime* as it runs inside the context of the application. With *Neuron Runtime 2.x*, the act of starting and stopping a Neuron application causes `libnrt.so` to automatically claim or release ownership of the required Neuron devices.

- **NeuronCore Groups (NCG) deprecation**

Before the introduction of *Neuron Runtime 2.x*, NeuronCore Group (NCG) was used to define an execution group of one or more NeuronCores where models could be loaded and executed. It also provided separation between processes.

With the introduction of *Neuron Runtime 2.x*, strict separation of NeuronCores into groups is no longer necessary and NeuronCore Groups (NCG) has been deprecated. See [eol-ncg](#) for more information.

- **Running multiple *Neuron Runtimes***

Before the introduction of `libnrt.so`, it was necessary to run multiple `neuron-rtd` daemons to allocate Neuron devices for each `neuron-rtd`, using configuration files. After the introduction of `libnrt.so`, it will no longer necessary to run multiple `neuron-rtd` daemons to allocate Neuron devices to a specific Neuron application. With `libnrt.so` NeuronCores (A Neuron device includes multiple NeuronCores) are allocated to a particular application by using `NEURON_RT_VISIBLE_CORES` or `NEURON_RT_NUM_CORES` environment variables, for example:

```
NEURON_RT_VISIBLE_CORES=0-3 myapp1.py
NEURON_RT_VISIBLE_CORES=4-11 myapp2.py
```

Or

```
NEURON_RT_NUM_CORES=3 myapp1.py &
NEURON_RT_NUM_CORES=4 myapp2.py &
```

See [NeuronX Runtime Configuration](#) for details.

- **Logging**

Similar to *Neuron Runtime 1.x*, *Neuron Runtime 2.x* logs into syslog (verbose logging). To make debugging easier, *Neuron Runtime 2.x* also logs into the console (error-only logging). Refer to [NeuronX Runtime Configuration](#) to see how to increase or decrease logging verbosity.

- **Multi-process access to NeuronCores**

With the introduction of `libnrt.so`, it is no longer possible to load models from multiple processes on the same NeuronCore. A NeuronCore can only be accessed from the same process. Instead you can load models on a specific NeuronCore, using multiple threads from the same process.

---

**Note:** For optimal performance of multi-model execution, each NeuronCore executes a single model.

---

- **Neuron Runtime architecture**

*Neuron Runtime 2.x* is delivered as a shared library (`libnrt.so`) and is directly linked to Neuron framework extensions. `libnrt.so` is packaged and installed as part of Neuron framework extensions (e.g. TensorFlow Neuron, PyTorch Neuron, or MXNet Neuron), and does not require a separate installation step. Installing Neuron Runtime as part of the Neuron framework extensions simplifies installation and improves the user experience. In addition, since `libnrt.so` is directly linked to Neuron framework extensions, it enables faster communication between Neuron Runtime and Neuron Frameworks by eliminating gRPC interface overhead.

### Neuron framework extensions

Starting from *Neuron 1.16.0*, Neuron framework extensions (TensorFlow Neuron, PyTorch Neuron, or MXNet Neuron) are packaged together with `libnrt.so`. It is required to install the `aws-neuron-dkms` Driver version 2.1.5.0 or newer for proper operation. The `neuron-rtd` daemon that was installed in previous releases no longer works starting with Neuron 1.16.0.

To see details of Neuron component versions see `neuron-release-content`.

### TensorFlow model server

Starting from *Neuron 1.16.0*, the TensorFlow Neuron model server is packaged together with `libnrt.so` and expects `aws-neuron-dkms` version 2.1.5.0 or newer for proper operation.

---

**Note:** The TensorFlow Neuron model server included in *Neuron 1.16.0* runs from the directory in which it was installed and will not run properly if copied to a different location, due to its dependency on `libnrt.so`.

---

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

### Neuron tools

- `neuron-cli` - Starting from *Neuron 1.16.0*, `neuron-cli` enters maintenance mode. See `maintenance_neuron-cli` for more information.
- `neuron-top` - Starting from *Neuron 1.16.0*, `neuron-top` has a new user interface. See [Neuron Top User Guide](#) for more information.
- `neuron-monitor` - `neuron-monitor` was updated to support Neuron Runtime 2.x (`libnrt.so`)
  - See [Neuron Monitor User Guide](#) for an updated user guide of `neuron-monitor`.
  - See `neuron-monitor-upg` for a list of changes between *Neuron Monitor 2.x* and *Neuron Monitor 1.0*
  - See `neuron-monitor-bwc` for instructions for using *Neuron Monitor 2.x* with *Neuron Runtime 1.x* (`neuron-rtd`).

### How will this change affect me?

### Neuron installation and upgrade

As explained in “[How will this change affect the Neuron SDK?](#)”, starting from *Neuron 1.16.0*, `libnrt.so` requires the latest Neuron Driver (`aws-neuron-dkms`). In addition, it is no longer necessary to install `aws-neuron-runtime`. To install Neuron or to upgrade to latest Neuron version, follow the installation and upgrade instructions below:

- **PyTorch Neuron**
  - `install-neuron-pytorch`.
  - `update-neuron-pytorch`.
- **TensorFlow Neuron**

- install-neuron-tensorflow.
- update-neuron-tensorflow.
- **MXNet Neuron**
  - install-neuron-mxnet.
  - update-neuron-mxnet.

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

## Migrate your application to Neuron Runtime 2.x (libnrt.so)

For a successful migration from previous releases of your application to *Neuron 1.16.0* or newer, make sure you perform the following:

### 1. Prerequisite

Read “[How will this change affect the Neuron SDK?](#)”.

### 2. Make sure you are not using *Neuron Runtime 1.x* (`aws-neuron-runtime`)

- Remove any code that installs `aws-neuron-runtime` from any CI/CD scripts.
- Stop `neuron-rtd` by running `sudo systemctl stop neuron-rtd`
- Uninstall `neuron-rtd` by running `sudo apt remove aws-neuron-runtime` or `sudo yum remove aws-neuron-runtime`

### 3. Upgrade to your Neuron Framework of choice:

- update-neuron-pytorch.
- update-neuron-tensorflow.
- update-neuron-mxnet.

### 4. If you have code that starts and/or stops `neuron-rtd`

Remove any code that starts or stops `neuron-rtd` from any CI/CD scripts.

### 5. Application running multiple `neuron-rtd`

If your application runs multiple processes and requires running multiple `neuron-rtd` daemons:

- Remove the code that runs multiple `neuron-rtd` daemons.
- Instead of allocating Neuron devices to `neuron-rtd` through configuration files, use `NEURON_RT_VISIBLE_CORES` or `NEURON_RT_NUM_CORES` environment variables to allocate NeuronCores. See [NeuronX Runtime Configuration](#) for details.

If your application uses `NEURONCORE_GROUP_SIZES`, see the next item.

---

**Note:** `NEURON_RT_VISIBLE_CORES` and `NEURON_RT_NUM_CORES` environment variables enable you to allocate NeuronCores to an application. Allocating NeuronCores improves application granularity, because Neuron devices include multiple NeuronCores.

---

### 6. Application running multiple processes using `NEURONCORE_GROUP_SIZES`

- Consider using `NEURON_RT_VISIBLE_CORES` or `NEURON_RT_NUM_CORES` environment variables instead of `NEURONCORE_GROUP_SIZES`, which is being deprecated.

See *NeuronX Runtime Configuration* for details.

- If you are using TensorFlow Neuron (tensorflow-neuron (TF2.x)) and you are replacing NEURONCORE\_GROUP\_SIZES=AxB which enables auto multicore replication, see the new API *TensorFlow 2.x (tensorflow-neuron) Auto Multicore Replication (Beta)* for usage and documentation.
- The behavior of your application will remain the same as before if you do not set NEURON\_RT\_VISIBLE\_CORES and do not set NEURON\_RT\_NUM\_CORES.
- If you are considering migrating to NEURON\_RT\_VISIBLE\_CORES or NEURON\_RT\_NUM\_CORES:
  - NEURON\_RT\_VISIBLE\_CORES takes precedence over NEURON\_RT\_NUM\_CORES.
  - If you are migrating to NEURON\_RT\_VISIBLE\_CORES:
    - \* For TensorFlow applications or PyTorch applications make sure that NEURONCORE\_GROUP\_SIZES is unset, or that NEURONCORE\_GROUP\_SIZES allocates the same or smaller number of NeuronCores as allocated by NEURON\_RT\_VISIBLE\_CORES.
    - \* For MXNet applications, setting NEURONCORE\_GROUP\_SIZES and NEURON\_RT\_VISIBLE\_CORES environment variables at the same time is not supported. Use NEURON\_RT\_VISIBLE\_CORES only.
    - \* See *NeuronX Runtime Configuration* for more details on how to use NEURON\_RT\_VISIBLE\_CORES.
  - If you are migrating to NEURON\_RT\_NUM\_CORES:
    - \* Make sure that NEURONCORE\_GROUP\_SIZES is unset.
    - \* See *NeuronX Runtime Configuration* for more details on how to use NEURON\_RT\_NUM\_CORES.

#### 7. Application running multiple processes accessing the same NeuronCore

If your application accesses the same NeuronCore from multiple processes, this is no longer possible with `libnrt.so`. Instead, modify your application to access the same NeuronCore from multiple threads.

---

**Note:** Optimal performance of multi-model execution is achieved when each NeuronCore executes a single model.

---

#### 8. Neuron Tools

- If you are using Neuron Monitor, see `neuron-monitor-upg` for details.
- If you are using `neuron-cli` remove any call to `neuron-cli`. For more information, see `maintenance_neuron-cli`.

#### 9. Containers

If your application is running within a container, and it previously executed `neuron-rtd` within the container, you need to re-build your container, so it will not include or install `aws-neuron-runtime`. See `neuron-containers` and `containers-migration-to-runtime2` for details.

## Troubleshooting

### Application fails to start

#### Description

Starting with the *Neuron 1.16.0* release, Neuron Runtime (`libnrt.so`) requires *Neuron Driver 2.0* or greater (`aws-neuron-dkms`). Neuron Runtime requires the Neuron Driver (`aws-neuron-dkms` package) to access Neuron devices.

If `aws-neuron-dkms` is not installed, the application will fail with an error message on the console and syslog similar to the following:

```
NRT:nrt_init      Unable to determine Neuron Driver version. Please check aws-neuron-
↪dkms package is installed.
```

If an old `aws-neuron-dkms` is installed, the application will fail with an error message on the console and syslog similar to the following:

```
NRT:nrt_init      This runtime requires Neuron Driver version 2.0 or greater. Please
↪upgrade aws-neuron-dkms package.
```

#### Solution

Follow the installation steps in `install-guide-index` to install `aws-neuron-dkms`.

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

### Application fails to start although I installed latest `aws-neuron-dkms`

#### Description

Starting from the *Neuron 1.16.0* release, Neuron Runtime (`libnrt.so`) requires *Neuron Driver 2.0* or greater (`aws-neuron-dkms`). If an old `aws-neuron-dkms` is installed, the application will fail. You may try to install `aws-neuron-dkms` and still face application failure, because the `aws-neuron-dkms` installation failed as a result of `neuron-rtd` daemon that was still running.

#### Solution

- Stop `neuron-rtd` by running: `sudo systemctl stop neuron-rtd`
- Uninstall `neuron-rtd` by running: `sudo apt remove aws-neuron-runtime` or `sudo yum remove aws-neuron-runtime`
- Install `aws-neuron-dkms` by following steps in `install-guide-index`

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

## Application unexpected behavior when upgrading to release *Neuron 1.16.0* or newer

### Description

When upgrading to release *Neuron 1.16.0* or newer from previous releases, the OS may include two different versions of *Neuron Runtime*: the `libnrt.so` shared library and `neuron-rtd` daemon. This can happen if the user did not stop `neuron-rtd` daemon or did not make sure to uninstall the existing Neuron version before upgrade. In this case the user application may behave unexpectedly.

### Solution

If the OS includes two different versions of *Neuron Runtime*, `libnrt.so` shared library and `neuron-rtd` daemon:

- Before running applications that use `neuron-rtd`, restart `neuron-rtd` by calling `sudo systemctl restart neuron-rtd`.
- Before running applications linked with `libnrt.so`, stop `neuron-rtd` by calling `sudo systemctl stop neuron-rtd`.

## Application unexpected behavior when downgrading to releases before *Neuron 1.6.0* (from *Neuron 1.16.0* or newer)

### Description

When upgrading to release *Neuron 1.16.0* or newer from previous releases, and then downgrading back to releases before *Neuron 1.6.0*, the OS may include two different versions of *Neuron Runtime*: the `libnrt.so` shared library and `neuron-rtd` daemon. This can happen if the user did not make sure to uninstall the existing Neuron version before the upgrade or downgrade. In this case the user application may behave unexpectedly.

### Solution

If the OS include two different versions of *Neuron Runtime*, `libnrt.so` shared library and `neuron-rtd` daemon:

- Before running applications that use `neuron-rtd`, restart `neuron-rtd` by calling `sudo systemctl restart neuron-rtd`.
- Before running applications linked with `libnrt.so`, stop `neuron-rtd` by calling `sudo systemctl stop neuron-rtd`.

## Neuron Core is in use

### Description

A Neuron Core cannot be shared between two applications. If an application started using a Neuron Core all other applications trying to use the NeuronCore will fail during runtime initialization with the following message in the console and in syslog:

```
ERROR    NRT:nrt_allocate_neuron_cores          NeuronCore(s) not available -  
↪Requested:nc1-nc1 Available:0
```

## Solution

Terminate the the process using NeuronCore and then try launching the application.

## Frequently Asked Questions (FAQ)

### Do I need to recompile my model to run it with Neuron Runtime 2.x (libnrt.so)?

No.

### Do I need to change my application launch command?

No.

### Can libnrt.so and neuron-rtd co-exist in the same environment?

Although we recommend upgrading to the latest Neuron release, we understand that for a transition period you may continue using `neuron-rtd` for old releases. If you are using Neuron Framework (PyTorch, TensorFlow or MXNet) from releases before *Neuron 1.16.0*:

- Install the latest Neuron Driver (`aws-neuron-dkms`)

---

**Important:** Starting with Neuron version 2.3, the `aws-neuron-dkms` package name has been changed to `aws-neuronx-dkms`. See [Introducing the first release of Neuron 2.x enabling EC2 Trn1 General Availability \(GA\)](#)

---

- For development, we recommend using different environments for Neuron Framework (PyTorch, TensorFlow or MXNet) from releases before *Neuron 1.16.0* and for Neuron Framework (PyTorch, TensorFlow or MXNet) from *Neuron 1.16.0* and newer. If that is not possible, make sure to stop `neuron-rtd` before executing models using Neuron Framework (PyTorch, TensorFlow or MXNet) from *Neuron 1.16.0* and newer.
- For deployment, when you are ready to upgrade, upgrade to Neuron Framework (PyTorch, TensorFlow or MXNet) from *Neuron 1.16.0* and newer. See [Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#) for more information.

**Warning:** Executing models using Neuron Framework (PyTorch, TensorFlow or MXNet) from *Neuron 1.16.0* and newer in an environment where `neuron-rtd` is running may cause undefined behavior. Make sure to stop `neuron-rtd` before executing models using Neuron Framework (PyTorch, TensorFlow or MXNet) from *Neuron 1.16.0* and newer.

### Are there Neuron framework versions that will not support Neuron Runtime 2.x (libnrt.so)?

All supported PyTorch Neuron and TensorFlow framework extensions, in addition to Neuron MXnet 1.8.0 framework extensions support Neuron Runtime 2.x.

Neuron MxNet 1.5.1 does not support Neuron Runtime 2.x (`libnrt.so`) and has now entered maintenance mode. See `maintenance_mxnet_1_5` for details.

*This document is relevant for: Inf1*

## Performance (Inf1)

*This document is relevant for: Inf1*

### 8.3.3 Performance Tuning

---

**Important:** NeuronCore Groups (NCG) have been deprecated. See `eol-ncg` and *Migrate your application to Neuron Runtime 2.x ([libnrt.so](#))* for more details.

---

This guide is intended to provide the reader with an in-depth understanding of how to optimize neural network performance on Inferentia for both throughput and latency. For simplicity, the guide uses the TensorFlow and ResNet-50 models as teaching examples to show how to choose between different compile-time optimizations (e.g., Batching and NeuronCore Pipeline), as well as model-serving optimizations (e.g., multi-threading and dynamic-batching) to improve inference performance.

The following guides are considered to be prerequisites for this tutorial:

- `/src/examples/tensorflow/tensorflow_resnet50/resnet50.ipynb`
- `tensorflow-serving-neurocore-group`
- *Neuron Batching*
- *NeuronCore Pipeline*

#### Batching and pipelining (technical background)

Neuron provides developers with various performance optimization features.

Two of the most widely used features are batching and pipelining. Both techniques aim to keep the data close to the compute engines, but they achieve this data locality in different ways. In batching it is achieved by loading the data into an on-chip cache and reusing it multiple times for multiple different model-inputs, while in pipelining it is achieved by caching all model parameters into the on-chip cache across multiple NeuronCores and streaming the calculation across them.

As a general rule of thumb, batching is preferred for applications that aim to optimize throughput and cost at the expense of latency, while pipelining is preferred for applications with a high-throughput requirement under a strict latency budget.

#### Compiling for batching optimization

To enable batching optimization, the model must first be compiled for a target batch-size. This is done by specifying the batch size in the input tensor's batch dimension during compilation. Users are encouraged to evaluate multiple batch size, in order to determine the optimal latency/throughput deployment-point, which is application-dependent.

For example, the code snippet below enables batching on a ResNet50 model, with a batch-size of 5:

```
import numpy as np
import tensorflow.neuron as tfn

# To change the batch size, change the first dimension in example_input
batch_size = 5
example_input = np.zeros([batch_size,224,224,3], dtype='float16')
```

(continues on next page)



(continued from previous page)

```
tfn.saved_model.compile("rn50_fp16",
                        "rn50_fp16_compiled/1",
                        model_feed_dict={'input_1:0': example_input },
                        dynamic_batch_size=True)
```

**Note:** Depending on the size of the neural network, Neuron has a maximum batch size that works optimally on Inferentia. If an unsupported batch size is used, an internal compiler error message will be displayed. A simple way to explore optimal batch size for your specific model is to increment the batch size from 1 upward, one at a time, and test application performance.

### Compiling for pipeline optimization

In NeuronCore Pipeline mode, Neuron stores the model parameters in Inferentia's local cache and streams inference requests across the available NeuronCores, as specified by the `--neuroncore-pipeline-cores` compiler argument. For example, to compile the model to fit a pipeline size of four Inferentia devices (16 NeuronCores) available in the `inf1.6xlarge` instance size:

```
import numpy as np
import tensorflow.neuron as tfn

compiler_args = ['--neuroncore-pipeline-cores', '16']
example_input = np.zeros([1,224,224,3], dtype='float16')
tfn.saved_model.compile("rn50_fp16",
                        "rn50_fp16_compiled/1",
                        model_feed_dict={'input_1:0': example_input },
                        compiler_args=compiler_args)
```

The minimum number of NeuronCores needed to run a compiled model can be found using the Neuron Check Model tool. See [Neuron Check Model](#).

### Model-serving inference optimizations

To fully realize the maximum throughput of the compiled model (for either batching and pipelining), users need to launch multiple host CPU threads to feed inputs into the Neuron pipeline. The number of threads needs to be larger than the specified maximum number of NeuronCores.

Additionally, dynamic batching can be used to process a larger client-side inference batch-size and the framework automatically breaks up the user-batch into smaller batch sizes, to match the compiled batch-size. This technique increases the achievable throughput by hiding the framework-to-neuron overhead, and amortizing it over a larger batch size. To use dynamic batching, set the argument `--dynamic_batch_size=True` during compilation and send a larger inference batch size (user inference batch size) that is equal to a multiple of the compiled batch size.

Both methods can be applied together if this improves performance. However, multi-threading is always needed as a first step to achieve high throughput. You need to experiment to find optimal settings for your application.

By default the framework sets the number of outstanding inference requests to the total number of NeuronCores plus three. This can be changed by setting the `NEURON_MAX_NUM_INFERS` environment variable. For example, if the compiled model includes CPU partitions (e.g., if the Neuron compiler decides that some operations are more efficient to execute on CPU), the number of threads needs to be increased to account for the additional compute performed on the

CPU. Note that the available instance host memory size needs to be taken into consideration to prevent out-of-memory errors. As above, you need to experiment in order to find the optimal settings for your application.

---

**Note:** By default the framework allocates a NeuronCore Group size to match the size of the compiled model. The size of the model is the number of NeuronCores limit passed to compiler during compilation (`--neuroncore-pipeline-cores` option). For more information see the `tensorflow-serving-neurocore-group`.

---

## Other considerations

### Mixed Precision

You can find more information about performance and accuracy trade offs in *Mixed precision and performance-accuracy tuning (neuron-cc)*.

### Operator support

The Neuron Compiler maintains an evolving list of supported operators for each framework: *Neuron Supported operators*

AWS Neuron handles unsupported operators by partitioning the graph into subgraphs and executing them on different targets (e.g., NeuronCore partition, CPU partition). If the entire model can run on Inferentia (i.e., all operators are supported), then it will be compiled into a single subgraph, which will be executed by a NeuronCore Group.

### Debug

You can examine the post-compiled model to view the compilation results using the Neuron plugin for TensorBoard. See *Visualize graphs executed on Neuron*.

### ResNet-50 optimization example

For an example demonstrating the concepts described here, see `/src/examples/tensorflow/keras_resnet50/keras_resnet50.ipynb`

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## 8.3.4 Parallel Execution using NEURON\_RT\_NUM\_CORES

---

**Important:** `NEURONCORE_GROUP_SIZES` will no longer be supported starting with the Neuron 1.19.0 release. If your application uses `NEURONCORE_GROUP_SIZES` see *Migrate your application to Neuron Runtime 2.x (libnrt.so)* and `eol-ncgs-env_2` for more details.

---

## Introduction

Inf1 instances are available with a different number of Inferentia chips. Each Inferentia chip consists of 4 NeuronCores and an Inf1 instance includes 4 to 64 NeuronCores, depending on the size of the instance. This guide shows you how to load one or more compiled models into different consecutive groups of NeuronCores using your framework of choice.

## Data Parallel Execution

In PyTorch and TensorFlow, the same compiled model can run in parallel on an Inf1 instance by loading it multiple times, up to the total number of NeuronCores specified in `NEURON_RT_NUM_CORES` or `NEURON_RT_VISIBLE_CORES`. For more information about `NEURON_RT_NUM_CORES` and `NEURON_RT_VISIBLE_CORES`, refer to *Neuron Runtime Configuration*.

## Running multiple models using single process

To run multiple models using a single process, set the environment variable `NEURON_RT_NUM_CORES` with a list of the number of cores in each group, separated by commas.

You can set the `NEURON_RT_NUM_CORES` environment variable at runtime:

```
#!/bin/bash
NEURON_RT_NUM_CORES=13 python your_neuron_application.py
```

Or from within the Python process running your models (NOTE: You can only set it once in the same process at the beginning of the script):

```
#!/usr/bin/env python
import os

# Set Environment
os.environ['NEURON_RT_NUM_CORES']='13'

# Load models and run inferences ...
```

The following examples allow you to load 4 models into 4 groups of NeuronCores within one process. For example, if there are 4 models A, B, C, D compiled to 2, 4, 3, and 4 NeuronCores respectively, directly load the models A, B, C, D in sequence within your TensorFlow or PyTorch Neuron process. This example requires an `inf1.6xlarge` instance with 16 NeuronCores, as the total number of NeuronCores within the NeuronCore Groups is 13.

In MXNet, mapping from models to NeuronCores is controlled by context `mx.neuron(neuron_core_index)` where `neuron_core_index` is the NeuronCore index at the start of the group. In the example above, map model A to `mx.neuron(0)` context, model B to `mx.neuron(2)` context, model C to `mx.neuron(6)` context and model D to `mx.neuron(9)` context. For further details, refer to *Flexible Execution Group (FlexEG) in Neuron-MXNet*.

For PyTorch

See *Data Parallel Inference on Torch Neuron* for more details.

For Tensorflow

```
# Set Environment
os.environ['NEURON_RT_NUM_CORES']='13'

# Load models (TF2)
```

(continues on next page)

(continued from previous page)

```

model0 = tf.keras.models.load_model(model0_file) # loaded into the first group of NC0-NC1
model1 = tf.keras.models.load_model(model1_file) # loaded into the second group of NC2-
↳ NC5
model2 = tf.keras.models.load_model(model1_file) # loaded into the third group of NC6-NC8
model3 = tf.keras.models.load_model(model1_file) # loaded into the fourth group of NC9-
↳ NC12

# run inference by simply calling the loaded model
results0 = model0(inputs0)
results1 = model1(inputs1)
results2 = model2(inputs2)
results3 = model3(inputs3)

```

For MXNet 2.x:

```

# Set Environment
os.environ['NEURON_RT_NUM_CORES']='13'

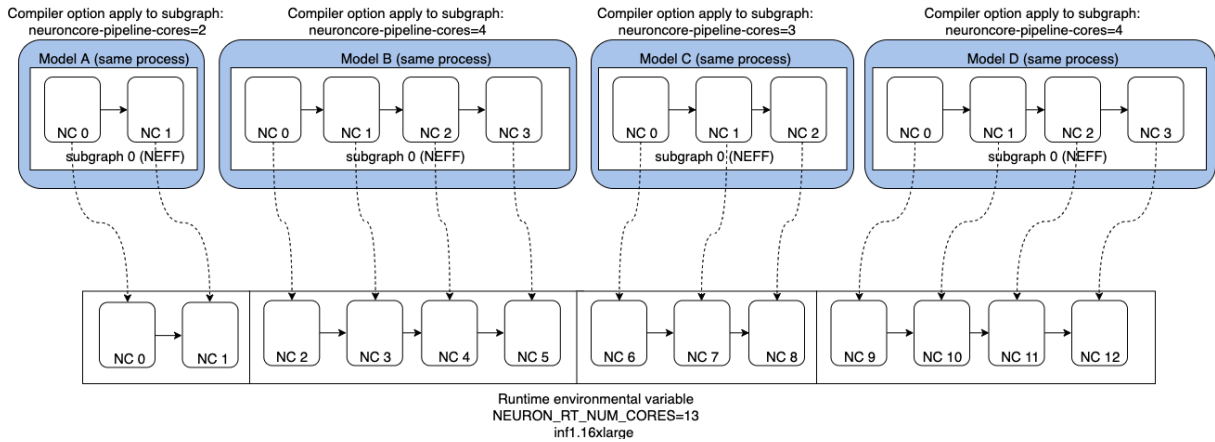
# Load models (MXNet)
# loaded into the first group of NC0-NC1
sym, args, aux = mx.model.load_checkpoint(mx_model0_file, 0)
model0 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')
# loaded into the second group of NC2-NC5
sym, args, aux = mx.model.load_checkpoint(mx_model1_file, 0)
model1 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')
# loaded into the third group of NC6-NC8
sym, args, aux = mx.model.load_checkpoint(mx_model2_file, 0)
model2 = sym.bind(ctx=mx.neuron(6), args=args, aux_states=aux, grad_req='null')
# loaded into the fourth group of NC9-NC12
sym, args, aux = mx.model.load_checkpoint(mx_model3_file, 0)
model3 = sym.bind(ctx=mx.neuron(9), args=args, aux_states=aux, grad_req='null')

# run inference by simply calling the loaded model
results0 = model0.forward(data=inputs0)
results1 = model1.forward(data=inputs1)
results2 = model2.forward(data=inputs2)
results3 = model3.forward(data=inputs3)

```

You can identify the NeuronCores used by each application with the `neuron-top` command line tool. For more information about the neuron-top user interface, see [Neuron Top User Guide](#).

```
$ neuron-top
```

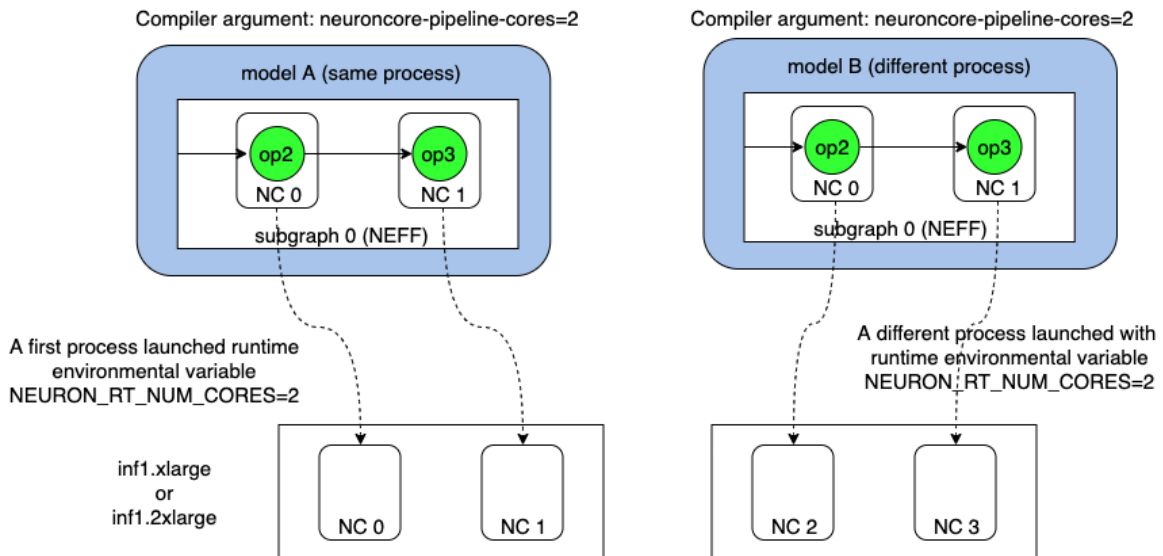


## Running multiple models using multiple processes

You can also run multiple models in parallel processes, when you set `NEURON_RT_NUM_CORES` per process:

```
$ NEURON_RT_NUM_CORES=2 python your_1st_neuron_application.py
$ NEURON_RT_NUM_CORES=2 python your_2nd_neuron_application.py
```

The first process automatically selects a first set of 2 unused NeuronCores for its new group. The second process automatically selects a new set of 2 unused NeuronCores for its new group.



## Running multiple models on the same NeuronCore group

You can load more than one model in a NeuronCore group within one process. Neuron runtime handles switching from one model to the next model within the NeuronCore group, when the next model is run within the application. In TensorFlow or PyTorch, simply load the additional models after the initial number of models have been loaded, to fill the NeuronCore groups associated with the process.

For PyTorch:

```
# Set Environment
os.environ['NEURON_RT_NUM_CORES']='2'

# Load models (PT)
model0 = torch.jit.load(model0_file) # loaded into the first group of NC0-NC1
model1 = torch.jit.load(model1_file) # loaded into the first group of NC0-NC1

# run inference by simply calling the loaded model
results0 = model0(inputs0)
results1 = model1(inputs1)
```

For TensorFlow 2.x:

```
# Set Environment
os.environ['NEURON_RT_NUM_CORES']='2'

# Load models (TF2)
model0 = tf.keras.models.load_model(model0_file) # loaded into the first group of NC0-NC1
model1 = tf.keras.models.load_model(model1_file) # loaded into the first group of NC0-NC1

# run inference by simply calling the loaded model
results0 = model0(inputs0)
results1 = model1(inputs1)
```

In MXNet, use context `mx.neuron(neuron_core_index)` and use the same NeuronCore start index for the additional models.

```
# Set Environment
os.environ['NEURON_RT_NUM_CORES']='2'

# Load models (MXNet)
# loaded into the first group of NC0-NC1
sym, args, aux = mx.model.load_checkpoint(mx_model0_file, 0)
model0 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')
# loaded into the first group of NC0-NC1
sym, args, aux = mx.model.load_checkpoint(mx_model1_file, 0)
model1 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')

# run inference by simply calling the loaded model
results0 = model0.forward(data=inputs0)
results1 = model1.forward(data=inputs1)
```

The total `NEURON_RT_NUM_CORES` across all processes cannot exceed the number of NeuronCores available on the instance. For example, on an `inf1.xlarge` with default configurations where the total number of NeuronCores visible to TensorFlow-Neuron is 4, you can launch one process with `NEURON_RT_NUM_CORES=2` (pipelined) and another process with `NEURON_RT_NUM_CORES=2` (data-parallel).

Examples using NEURON\_RT\_NUM\_CORES include:

- PyTorch example
- MXNet example

### Auto Model Replication in TensorFlow Neuron (tensorflow-neuron) (Beta)

Refer to the following API documentation to see how to perform automatic replication on multiple cores. Note auto-replication will only work on models compiled with pipeline size 1: via `--neuroncore-pipeline-cores=1`. If automatic replication is not enabled, the model will default to replicate on up to 4 cores.

Python API (TF 2.x only):

*TensorFlow 2.x (tensorflow-neuron) Auto Multicore Replication (Beta)*

CLI API (TF 1.x and TF 2.x):

*TensorFlow Neuron TF2.x (tensorflow-neuronx TF2.x) Auto Multicore Replication CLI (Beta)*

### Auto Model Replication (Being Deprecated)

The Auto Model Replication feature in TensorFlow-Neuron enables you to load the model once and the data parallel replication will occur automatically. This reduces framework memory usage, as the same model is not loaded multiple times. This feature is beta and available in TensorFlow-Neuron only.

To enable Auto Model Replication, set NEURONCORE\_GROUP\_SIZES to Nx1, where N is the desired replication count (the number of NeuronCore groups, each group has size 1). For example, NEURONCORE\_GROUP\_SIZES=8x1 would automatically replicate the single-NeuronCore model 8 times.

```
os.environ['NEURONCORE_GROUP_SIZES'] = '4x1'
```

or

```
NEURONCORE_GROUP_SIZES=4x1 python3 application.py
```

When NEURONCORE\_GROUP\_SIZES is not set, the default is 4x1, where a single-NeuronCore model is replicated 4 times on any size of inf1 machine.

This feature is only available for models compiled with neuroncore-pipeline-cores set to 1 (default).

You will still need to use threads in the scaffolding code, to feed the loaded replicated model instance, to achieve high throughput.

Example of auto model replication: `/src/examples/tensorflow/openpose_demo/openpose.ipynb`

## FAQ

### Can I mix data parallel and NeuronCore Pipelines?

Yes. You can compile the model using the neuroncore-pipeline-cores option. This tells the compiler to set compilation to the specified number of cores for *NeuronCore Pipeline*. The Neuron Compiler returns a NEFF that fits within this limit. See the *Neuron compiler CLI Reference Guide (neuron-cc)* for instructions on how to use this option.

For example, on an inf1.2xlarge, you can load two model instances, each compiled with neuroncore-pipeline-cores set to 2, so they can run in parallel. The model instances can be loaded from different saved models or from the same saved model.

## Can I have a mix of multiple models in one Neuroncore group and single model in another one Neuroncore group?

Currently, you can do this in MXNet, by setting up two Neuroncore groups, then loading, for example, multiple models in one NCG, using context `mx.neuron(0)`, and loading a single model in the second NCG, using context `mx.neuron(2)`. You can also load a single model in the first NCG and multiple models in the second NCG. For example:

```
# Set Environment
os.environ['NEURON_RT_NUM_CORES']='6'

# Load models (MXNet)
# loaded into the first group of NC0-NC1
sym, args, aux = mx.model.load_checkpoint(mx_model0_file, 0)
model0 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')
# loaded into the second group of NC2-NC5
sym, args, aux = mx.model.load_checkpoint(mx_model1_file, 0)
model1 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')
# loaded into the second group of NC2-NC5
sym, args, aux = mx.model.load_checkpoint(mx_model2_file, 0)
model2 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')
# loaded into the second group of NC2-NC5
sym, args, aux = mx.model.load_checkpoint(mx_model3_file, 0)
model3 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')

# run inference by simply calling the loaded model
results0 = model0.forward(data=inputs0)
results1 = model1.forward(data=inputs1)
results2 = model2.forward(data=inputs2)
results3 = model3.forward(data=inputs3)
```

Loading multiple models in one NCG and a single model in another NCG is currently not supported in TensorFlow and PyTorch.

*This document is relevant for:* Inf1

## PyTorch Neuron (torch-neuron)

*This document is relevant for:* Inf1

### 8.3.5 Running R-CNNs on Inf1

This application note demonstrates how to compile and run [Detectron2](#)-based R-CNNs on Inf1. It also provides guidance on how to use profiling to improve performance of R-CNN models on Inf1.

#### Table of contents

- [R-CNN Model Overview](#)
  - [R-CNN Limitations and Considerations on Inferentia \(NeuronCore-v1\)](#)
- [Requirements](#)
- [Installation](#)



- *Compiling an R-CNN for Inference*
  - *Create a Detectron2 R-CNN Model*
  - *Profile the Model*
  - *Compiling the ResNet backbone to Inference*
- *Optimize the R-CNN model*
  - *Compiling the RPN*
  - *Fusing the Backbone and RPN Head*
  - *Compiling the RoI Heads*
- *End-to-end Compilation and Inference*
  - *Benchmarking*
  - *Other improvements*
    - \* *For latency sensitive applications:*
    - \* *For throughput sensitive applications:*

## R-CNN Model Overview

Region-based CNN (R-CNN) models are commonly used for object detection and image segmentation tasks. A typical R-CNN architecture consists of the following components:

- **Backbone:** The backbone extracts features from input images. In some models the backbone is a Feature Pyramid Network (FPN), which uses a top-down architecture with lateral connections to build an in-network feature pyramid from a single-scale input. The backbone is commonly a ResNet or Vision Transformer based network.
- **Region Proposal Network (RPN):** The RPN predicts region proposals with a wide range of scales and aspect ratios. RPNs are constructed using convolutional layers and anchor boxes, which that serve as references for multiple scales and aspect ratios.
- **Region of Interest (RoI):** The RoI component is used to resize the extracted features of varying size to the same size so that they can be consumed by a fully connected layer. RoI Align is typically used instead of RoI Pooling, because RoI Align provides better alignment.

The [Detectron2](#) library provides many popular PyTorch R-CNN implementations, including R-CNN, Fast R-CNN, Faster R-CNN, and Mask R-CNN. This application note focuses on the Detectron2 R-CNN models.

## R-CNN Limitations and Considerations on Inferentia (NeuronCore-v1)

R-CNN models may have limitations and considerations on Inferentia (NeuronCore-v1). See the Model Architecture Fit Guidelines for more information. These limitations are not applicable to NeuronCore-v2.

### Requirements

The process described in this application note is intended to be run on an `inf1.2xlarge`. In practice, R-CNN models can be run on any Inf1 instance size.

Verify that this Jupyter notebook is running the Python kernel environment that was set up according to the [PyTorch Installation Guide](#). Select the kernel from the “Kernel -> Change Kernel” option at the top of the Jupyter notebook page.

### Installation

This process requires the following pip packages:

- `torch==1.11.0`
- `torch-neuron`
- `neuron-cc`
- `opencv-python`
- `pycocotools`
- `torchvision==0.12.0`
- `detectron2==0.6`

The following section explains how to build `torchvision` from source and install the `Detectron2` package. It also reinstalls the Neuron packages, to ensure version compatibility.

The `torchvision roi_align_kernel.cpp` kernel is modified to use OMP threading for a multi-threaded inference on the CPU. This significantly improves the performance of RoI Align kernels on Inf1: OMP threading leads to a RoI Align latency reduction two to three times larger than the default `roi_align_kernel.cpp` kernel configuration.

```
# Install python3.7-dev for pycocotools (a Detectron2 dependency)
!sudo apt install python3.7-dev -y

# Install Neuron packages
!pip config set global.extra-index-url https://pip.repos.neuron.amazonaws.com
!pip uninstall -y torchvision
!pip install --force-reinstall torch-neuron==1.11.0.* neuron-cc[tensorflow] "protobuf==3.
↪20.1" ninja opencv-python

# Change cuda to 10.2 for Detectron2
!sudo rm /usr/local/cuda
!sudo ln -s /usr/local/cuda-10.2 /usr/local/cuda

# Install Torchvision 0.12.0 from source
!git clone -b release/0.12 https://github.com/pytorch/vision.git

# Update the RoI Align kernel to use OMP multithreading
with open('vision/torchvision/csrc/ops/cpu/roi_align_kernel.cpp', 'r') as file:
```

(continues on next page)

(continued from previous page)

```

content = file.read()

# Enable OMP Multithreading and set the number of threads to 4
old = "// #pragma omp parallel for num_threads(32)"
new = "#pragma omp parallel for num_threads(4)"
content = content.replace(old, new)

# Re-write the file
with open('vision/torchvision/csrc/ops/cpu/roi_align_kernel.cpp', 'w') as file:
    file.write(content)

# Build Torchvision with OMP threading
!cd vision && CFLAGS="-fopenmp" python setup.py bdist_wheel
%pip install vision/dist/*.whl

# Install Detectron2 release v0.6
!python -m pip install 'git+https://github.com/facebookresearch/detectron2.git@v0.6'

```

## Compiling an R-CNN for Inf1

By default, R-CNN models are not compilable on Inf1, because they cannot be traced with `torch.jit.trace`, which is a requisite for inference on Inf1. The following section demonstrates techniques for compiling a Detectron2 R-CNN model for inference on Inf1.

Specifically, this section explains how to create a standard Detectron2 R-CNN model, using a ResNet-101 backbone. It demonstrates how to use profiling to identify the most compute-intensive parts of the R-CNN that need to be compiled for accelerated inference on Inf1. It then explains how to manually extract and compile the ResNet backbone (the dominant compute component) and inject the compiled backbone back into the full model, for improved performance.

## Create a Detectron2 R-CNN Model

Create a Detectron2 R-CNN model using the `COCO-Detection/faster_rcnn_R_101_FPN_3x.yaml` pretrained weights and config file. Download a sample image from the COCO dataset and run an example inference.

```

from detectron2 import model_zoo
from detectron2.engine import DefaultPredictor
from detectron2.config import get_cfg

def get_model():

    # Configure the R-CNN model
    CONFIG_FILE = "COCO-Detection/faster_rcnn_R_101_FPN_3x.yaml"
    WEIGHTS_FILE = "COCO-Detection/faster_rcnn_R_101_FPN_3x.yaml"
    cfg = get_cfg()
    cfg.merge_from_file(model_zoo.get_config_file(CONFIG_FILE))
    cfg.MODEL.WEIGHTS = model_zoo.get_checkpoint_url(WEIGHTS_FILE)
    cfg.MODEL.ROI_HEADS.SCORE_THRESH_TEST = 0.5
    cfg.MODEL.DEVICE = 'cpu' # Send to CPU for Neuron Tracing

    # Create the R-CNN predictor wrapper

```

(continues on next page)

(continued from previous page)

```

predictor = DefaultPredictor(cfg)
return predictor

```

```

import os
import urllib.request

# Define a function to get a sample image
def get_image():
    filename = 'input.jpg'
    if not os.path.exists(filename):
        url = "http://images.cocodataset.org/val2017/000000439715.jpg"
        urllib.request.urlretrieve(url, filename)
    return filename

```

```

import time
import cv2

# Create an R-CNN model
predictor = get_model()

# Get a sample image from the COCO dataset
image_filename = get_image()
image = cv2.imread(image_filename)

# Run inference and print inference latency
start = time.time()
outputs = predictor(image)
print(f'Inference time: {(time.time() - start):0.3f} s')

```

## Profile the Model

Use the [PyTorch Profiler](#) to identify which operators contribute the most to the model's runtime on CPU. Ideally, you can compile these compute intensive operators onto Inf1 for accelerated inference.

```

import torch.autograd.profiler as profiler

with profiler.profile(record_shapes=True) as prof:
    with profiler.record_function("model_inference"):
        predictor(image)
print(prof.key_averages().table(sort_by="cpu_time_total", row_limit=30))

```

We see that convolution operators (`aten::convolution`) contribute the most to inference time. By compiling these convolution operators to Inf1, you can improve performance of the R-CNN model. Print the R-CNN model architecture to see which layers contain the `aten::convolution` operators:

```

print(predictor.model)

```

Note that the ResNet FPN backbone (`predictor.model.backbone L17-L162`) contains the majority of convolution operators in the model. The RPN (`predictor.model.proposal_generator L181-L533`) also contains several convolutions. Based on this, compile the ResNet backbone and RPN onto Inf1 to maximize performance.

## Compiling the ResNet backbone to Inf1

This section demonstrates how to compile the ResNet backbone to Inf1 and use it for inference.

Extract the backbone by accessing it with `predictor.model.backbone`. Compile the backbone using `strict=False`, because the backbone outputs a dictionary. Use a fixed input shape (800 x 800) for compilation, as all inputs will be resized to this shape during inference. This section also defines a basic preprocessing function (mostly derived from the Detectron2 R-CNN `DefaultPredictor` module L308-L318) that reshapes inputs to 800 x 800.

Create a `NeuronRCNN` wrapper to inject the compiled backbone back into the model by dynamically replacing the `predictor.model.backbone` attribute with the compiled model.

```
import torch
import torch_neuron

example = torch.rand([1, 3, 800, 800])

# Use `with torch.no_grad():` to avoid a jit tracing issue in the ResNet backbone
with torch.no_grad():
    neuron_backbone = torch_neuron.trace(predictor.model.backbone, example, strict=False)

backbone_filename = 'backbone.pt'
torch.jit.save(neuron_backbone, backbone_filename)
```

```
from detectron2.modeling.meta_arch.rcnn import GeneralizedRCNN
from torch.jit import ScriptModule

class NeuronRCNN(torch.nn.Module):
    """
    Creates a `NeuronRCNN` wrapper that injects the compiled backbone into
    the R-CNN model. It also stores the `size_divisibility` attribute from
    the original backbone.
    """

    def __init__(self, model: GeneralizedRCNN, neuron_backbone: ScriptModule) -> None:
        super().__init__()

        # Keep track of the backbone variables
        size_divisibility = model.backbone.size_divisibility

        # Load and inject the compiled backbone
        model.backbone = neuron_backbone

        # Set backbone variables
        setattr(model.backbone, 'size_divisibility', size_divisibility)

        self.model = model

    def forward(self, x):
        return self.model(x)
```

```
# Create the R-CNN with the compiled backbone
neuron_rcnn = NeuronRCNN(predictor.model, neuron_backbone)
```

(continues on next page)

(continued from previous page)

```

neuron_rcnn.eval()

# Print the R-CNN architecture to verify the backbone is now the
# `neuron_backbone` (shows up as `RecursiveScriptModule`)
print(neuron_rcnn)

def preprocess(original_image, predictor):
    """
    A basic preprocessing function that sets the input height=800 and
    input width=800. The function is derived from the preprocessing
    steps in the Detectron2 `DefaultPredictor` module.
    """

    height, width = original_image.shape[:2]
    resize_func = predictor.aug.get_transform(original_image)
    resize_func.new_h = 800 # Override height
    resize_func.new_w = 800 # Override width
    image = resize_func.apply_image(original_image)
    image = torch.as_tensor(image.astype("float32").transpose(2, 0, 1))
    inputs = {"image": image, "height": height, "width": width}
    return inputs

# Get a resized input using the sample image
inputs = preprocess(image, get_model())

# Run inference and print inference latency
start = time.time()
for _ in range(10):
    outputs = neuron_rcnn([inputs])[0]
print(f'Inference time: {(time.time() - start)/10:.3f} s')

with profiler.profile(record_shapes=True) as prof:
    with profiler.record_function("model_inference"):
        neuron_rcnn([inputs])
print(prof.key_averages().table(sort_by="cpu_time_total", row_limit=30))

```

By running the backbone on Inf1, the overall runtime is already significantly improved. The count and runtime of `aten::convolution` operators is also decreased. We now see a `neuron::forward_v2` operator that is the compiled backbone.

## Optimize the R-CNN model

### Compiling the RPN

Examine the profiling and note that there are still several `aten::convolution`, `aten::linear`, and `aten::addmm` operators that significantly contribute to the model's overall latency. By inspecting the model's architecture and code, we can determine that the majority of these operators are contained in the RPN module (`predictor.model.proposal_generator` L181-L533).

To improve the model's performance, extract the RPN Head and compile it on Inf1 to increase the number of operators running on Inf1. You need to compile the RPN Head, because the RPN Anchor Generator contains objects that are not

traceable with `torch.jit.trace`.

The RPN Head contains five layers that run inference on multiple resized inputs. To compile the RPN Head, create a list of tensors that contain the input (“features”) shapes used by RPN Head on each layer. These tensor shapes can be determined by printing the input shapes in the RPN Head forward function (`predictor.model.proposal_generator.rpn_head.forward`).

Create a new NeuronRCNN wrapper that injects both the compiled backbone and RPN Head into the R-CNN model.

```
import math

input_shape = [1, 3, 800, 800] # Overall input shape at inference time

# Create the list example of RPN inputs using the resizing logic from the RPN Head
features = list()
for i in [0, 1, 2, 3, 4]:
    ratio = 1 / (4 * 2**i)
    x_i_h = math.ceil(input_shape[2] * ratio)
    x_i_w = math.ceil(input_shape[3] * ratio)
    feature = torch.zeros(1, 256, x_i_h, x_i_w)
    features.append(feature)

# Extract and compile the RPN Head
neuron_rpn_head = torch_neuron.trace(predictor.model.proposal_generator.rpn_head,
    ↪[features])
rpn_head_filename = 'rpn_head.pt'
torch.jit.save(neuron_rpn_head, rpn_head_filename)

class NeuronRCNN(torch.nn.Module):
    """
    Creates a wrapper that injects the compiled backbone and RPN Head
    into the R-CNN model.
    """

    def __init__(self, model: GeneralizedRCNN, neuron_backbone: ScriptModule, neuron_rpn_
    ↪head: ScriptModule) -> None:
        super().__init__()

        # Keep track of the backbone variables
        size_divisibility = model.backbone.size_divisibility

        # Inject the compiled backbone
        model.backbone = neuron_backbone

        # Set backbone variables
        setattr(model.backbone, 'size_divisibility', size_divisibility)

        # Inject the compiled RPN Head
        model.proposal_generator.rpn_head = neuron_rpn_head

        self.model = model

    def forward(self, x):
        return self.model(x)
```

```
# Create the R-CNN with the compiled backbone and RPN Head
predictor = get_model()
neuron_rcnn = NeuronRCNN(predictor.model, neuron_backbone, neuron_rpn_head)
neuron_rcnn.eval()

# Print the R-CNN architecture to verify the compiled modules show up
print(neuron_rcnn)
```

```
# Run inference and print inference latency
start = time.time()
for _ in range(10):
    outputs = neuron_rcnn([inputs])[0]
print(f'Inference time: {(time.time() - start)/10}:0.3f s')
```

```
with profiler.profile(record_shapes=True) as prof:
    with profiler.record_function("model_inference"):
        neuron_rcnn([inputs])
print(prof.key_averages().table(sort_by="cpu_time_total", row_limit=30))
```

By running the compiled backbone and RPN Head on Inf1, overall runtime is improved. Once again, the number and runtime of `aten::convolution` operators is also decreased. There are now two `neuron::forward_v2` operators, which correspond to the compiled backbone and RPN Head.

## Fusing the Backbone and RPN Head

It is usually preferable to compile fewer independent models (“subgraphs”) on Inf1. Combining models and compiling them as a single subgraph enables the Neuron compiler to perform additional optimizations and reduces I/O data transfer between CPU and NeuronCores between each subgraph.

In this section, the ResNet backbone and RPN Head are “fused” into a single model to compile on Inf1. Create the `NeuronFusedBackboneRPNHead` wrapper as a compilable model that contains both the ResNet backbone (`predictor.model.backbone` L17-L162) and RPN Head (`predictor.model.proposal_generator` L181-L533). Output the features to be used downstream by the RoI Heads. Compile this `NeuronFusedBackboneRPNHead` wrapper as `neuron_backbone_rpn`, then create a separate `BackboneRPN` wrapper to inject the `neuron_backbone_rpn` in place of the original backbone and RPN Head. Copy the remainder of the RPN forward code (`predictor.model.proposal_generator.forward` L431-L480) to create a “fused” backbone + RPN module. Lastly, re-write the `NeuronRCNN` wrapper to use the fused backbone + RPN module. The `NeuronRCNN` wrapper also uses the `predictor.model` forward code to re-write the rest of the R-CNN model forward function.

```
class NeuronFusedBackboneRPNHead(torch.nn.Module):
    """
    Wrapper to compile the fused ResNet backbone and RPN Head.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.backbone = model.backbone
        self.rpn_head = model.proposal_generator.rpn_head
        self.in_features = model.proposal_generator.in_features

    def forward(self, x):
        features = self.backbone(x)
```

(continues on next page)



(continued from previous page)

```

features_ = [features[f] for f in self.in_features]
return self.rpn_head(features_), features

```

```

# Create the wrapper with the combined backbone and RPN Head
predictor = get_model()
backbone_rpn_wrapper = NeuronFusedBackboneRPNHead(predictor.model)
backbone_rpn_wrapper.eval()

# Compile the wrapper
example = torch.rand([1, 3, 800, 800])

with torch.no_grad():
    neuron_backbone_rpn_head = torch_neuron.trace(
        backbone_rpn_wrapper, example, strict=False)

backbone_rpn_filename = 'backbone_rpn.pt'
torch.jit.save(neuron_backbone_rpn_head, backbone_rpn_filename)

```

```

class BackboneRPN(torch.nn.Module):
    """
    Wrapper that uses the compiled `neuron_backbone_rpn` instead
    of the original backbone and RPN Head. We copy the remainder
    of the RPN `forward` code (`predictor.model.proposal_generator.forward`)
    to create a "fused" backbone + RPN module.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.backbone_rpn_head = NeuronFusedBackboneRPNHead(model)
        self._rpn = model.proposal_generator
        self.in_features = model.proposal_generator.in_features

    def forward(self, images):
        preds, features = self.backbone_rpn_head(images.tensor)
        features_ = [features[f] for f in self.in_features]
        pred_objectness_logits, pred_anchor_deltas = preds
        anchors = self._rpn.anchor_generator(features_)

        # Transpose the Hi*Wi*A dimension to the middle:
        pred_objectness_logits = [
            # (N, A, Hi, Wi) -> (N, Hi, Wi, A) -> (N, Hi*Wi*A, 1)
            score.permute(0, 2, 3, 1).flatten(1)
            for score in pred_objectness_logits
        ]
        pred_anchor_deltas = [
            # (N, A*B, Hi, Wi) -> (N, A, B, Hi, Wi) -> (N, Hi, Wi, A, B) -> (N, Hi*Wi*A, B)
            x.view(x.shape[0], -1, self._rpn.anchor_generator.box_dim,
                  x.shape[-2], x.shape[-1])
            .permute(0, 3, 4, 1, 2)
            .flatten(1, -2)

```

(continues on next page)

(continued from previous page)

```

        for x in pred_anchor_deltas
    ]

    proposals = self._rpn.predict_proposals(
        anchors, pred_objectness_logits, pred_anchor_deltas, images.image_sizes
    )
    return proposals, features

```

```

class NeuronRCNN(torch.nn.Module):
    """
    Wrapper that uses the fused backbone + RPN module and re-writes
    the rest of the R-CNN `model` `forward` function.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()

        # Use the fused Backbone + RPN
        self.backbone_rpn = BackboneRPN(model)

        self.roi_heads = model.roi_heads

        self.preprocess_image = model.preprocess_image
        self._postprocess = model._postprocess

    def forward(self, batched_inputs):
        images = self.preprocess_image(batched_inputs)
        proposals, features = self.backbone_rpn(images)
        results, _ = self.roi_heads(images, features, proposals, None)
        return self._postprocess(results, batched_inputs, images.image_sizes)

```

```

# Create the new NeuronRCNN wrapper with the combined backbone and RPN Head
predictor = get_model()
neuron_rcnn = NeuronRCNN(predictor.model)
neuron_rcnn.eval()

```

```

# Inject the Neuron compiled models
neuron_rcnn.backbone_rpn.backbone_rpn_head = neuron_backbone_rpn_head

# Print the R-CNN architecture to verify the compiled modules show up
print(neuron_rcnn)

```

```

# Run inference and print inference latency
start = time.time()
for _ in range(10):
    outputs = neuron_rcnn([inputs])[0]
print(f'Inference time: {(time.time() - start)/10:0.3f} s')

```

```

with profiler.profile(record_shapes=True) as prof:
    with profiler.record_function("model_inference"):
        neuron_rcnn([inputs])

```

(continues on next page)

(continued from previous page)

```
print(prof.key_averages().table(sort_by="cpu_time_total", row_limit=30))
```

By running the fused backbone + RPN Head on Inf1, overall runtime is improved even more. We now see a single `neuron::forward_v2` operator with a lower runtime than the previous combined runtime of the two separate `neuron::forward_v2` operators.

## Compiling the RoI Heads

This section describes how to extract and compile part of RoI Heads module (`predictor.model.roi_heads` L530-L778) which runs most of the remaining `aten::linear` and `aten::addmm` operators on Inf1. The entire RoI Heads module cannot be extracted, because it contains unsupported operators. So you need to create a `NeuronBoxHeadBoxPredictor` wrapper, extracts specific parts of the `roi_heads` for compilation. The example input for compilation is the shape of the input into the `self.roi_heads.box_head.forward` function. Write another wrapper, `ROIHead` that combines the compiled `roi_heads` into the rest of the RoI module. The `_forward_box` and `forward` functions are from the `predictor.model.roi_heads` module. Lastly, re-write the `NeuronRCNN` wrapper to use the optimized RoI Heads wrapper as well as the fused backbone + RPN module.

```
class NeuronBoxHeadBoxPredictor(torch.nn.Module):
    """
    Wrapper that extracts the RoI Box Head and Box Predictor
    for compilation.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.roi_heads = model.roi_heads

    def forward(self, box_features):
        box_features = self.roi_heads.box_head(box_features)
        predictions = self.roi_heads.box_predictor(box_features)
        return predictions
```

```
# Create the NeuronBoxHeadBoxPredictor wrapper
predictor = get_model()
box_head_predictor = NeuronBoxHeadBoxPredictor(predictor.model)
box_head_predictor.eval()

# Compile the wrapper
example = torch.rand([1000, 256, 7, 7])
neuron_box_head_predictor = torch_neuron.trace(box_head_predictor, example)

roi_head_filename = 'box_head_predictor.pt'
torch.jit.save(neuron_box_head_predictor, roi_head_filename)
```

```
class ROIHead(torch.nn.Module):
    """
    Wrapper that combines the compiled `roi_heads` into the
    rest of the RoI module. The `_forward_box` and `forward`
    functions are from the `predictor.model.roi_heads` module.
    """
```

(continues on next page)

(continued from previous page)

```

def __init__(self, model: GeneralizedRCNN) -> None:
    super().__init__()
    self.roi_heads = model.roi_heads
    self.neuron_box_head_predictor = NeuronBoxHeadBoxPredictor(model)

def _forward_box(self, features, proposals):
    features = [features[f] for f in self.roi_heads.box_in_features]
    box_features = self.roi_heads.box_pooler(
        features, [x.proposal_boxes for x in proposals])
    predictions = self.neuron_box_head_predictor(box_features)
    pred_instances, _ = self.roi_heads.box_predictor.inference(
        predictions, proposals)
    return pred_instances

def forward(self, images, features, proposals, targets=None):
    pred_instances = self._forward_box(features, proposals)
    pred_instances = self.roi_heads.forward_with_given_boxes(
        features, pred_instances)
    return pred_instances, {}

```

```

class NeuronRCNN(torch.nn.Module):
    """
    Wrapper that uses the fused backbone + RPN module and the optimized RoI
    Heads wrapper
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()

        # Create fused Backbone + RPN
        self.backbone_rpn = BackboneRPN(model)

        # Create Neuron RoI Head
        self.roi_heads = ROIHead(model)

        # Define pre and post-processing functions
        self.preprocess_image = model.preprocess_image
        self._postprocess = model._postprocess

    def forward(self, batched_inputs):
        images = self.preprocess_image(batched_inputs)
        proposals, features = self.backbone_rpn(images)
        results, _ = self.roi_heads(images, features, proposals, None)
        return self._postprocess(results, batched_inputs, images.image_sizes)

```

```

# Initialize an R-CNN on CPU
predictor = get_model()

# Create the Neuron R-CNN on CPU
neuron_rcnn = NeuronRCNN(predictor.model)
neuron_rcnn.eval()

```

(continues on next page)

(continued from previous page)

```
# Inject the Neuron compiled models
neuron_rcnn.backbone_rpn.backbone_rpn_head = neuron_backbone_rpn_head
neuron_rcnn.roi_heads.neuron_box_head_predictor = neuron_box_head_predictor
```

```
# Run inference and print inference latency
start = time.time()
for _ in range(10):
    outputs = neuron_rcnn([inputs])[0]
print(f'CPU Inference time: {((time.time() - start)/10):0.3f} s')
```

```
with profiler.profile(record_shapes=True) as prof:
    with profiler.record_function("model_inference"):
        neuron_rcnn([inputs])
print(prof.key_averages().table(sort_by="cpu_time_total", row_limit=30))
```

Although the overall latency did not change significantly, running more of the model on Inf1 instead of CPU frees up CPU resources when multiple models are running in parallel.

## End-to-end Compilation and Inference

This section provides standalone code that compiles and runs an optimized Detectron2 R-CNN on Inf1. Most of the code in this section is from the previous sections in this application note and is consolidated here for easy deployment. This section has the following main components:

- Preprocessing and compilation functions
- **Wrappers that extract the R-CNN ResNet backbone, RPN Head, and RoI Head for compilation on Inf1.**
- **A NeuronRCNN wrapper that creates an optimized end-to-end Detectron2 R-CNN model for inference on Inf1**
- **Benchmarking code that runs parallelized inference for optimized throughput on Inf1**

## Benchmarking

The benchmarking section explains how to load multiple optimized RCNN models and run them in parallel, to maximize throughput.

Use the beta NeuronCore placement API, `torch_neuron.experimental.neuron_cores_context()`, to ensure all compiled models in an optimized RCNN model are loaded onto the same NeuronCore. Note that the functionality and API of `torch_neuron.experimental.neuron_cores_context()` might change in future releases.

Define a simple benchmark function that loads four optimized RCNN models onto four separate NeuronCores, runs multithreaded inference, and calculates the corresponding latency and throughput. Benchmark various numbers of loaded models, to show the impact of parallelism.

Note that throughput increases (at the cost of latency) when more models are run in parallel on Inf1. Increasing the number of worker threads also improves throughput.

## Other improvements

There are many additional optimizations that can be applied to RCNN models on Inf1 depending on the application:

### For latency sensitive applications:

- Each of the five layers in the RPN head can be parallelized to decrease overall latency.
- The number of OMP Threads can be increased in the ROI Align kernel. Both of these optimizations improve latency, at the cost of decreasing throughput.

### For throughput sensitive applications:

- The input batch size can be increased to improve NeuronCore utilization.

```
import time
import os
import urllib.request
from typing import Any, Union, Callable

import cv2
import numpy as np
from concurrent.futures import ThreadPoolExecutor

import torch
import torch_neuron

from detectron2 import model_zoo
from detectron2.engine import DefaultPredictor
from detectron2.config import get_cfg
from detectron2.modeling.meta_arch.rcnn import GeneralizedRCNN

# -----
# Helper functions
# -----

def get_model():

    # Configure the R-CNN model
    CONFIG_FILE = "COCO-Detection/faster_rcnn_R_101_FPN_3x.yaml"
    WEIGHTS_FILE = "COCO-Detection/faster_rcnn_R_101_FPN_3x.yaml"
    cfg = get_cfg()
    cfg.merge_from_file(model_zoo.get_config_file(CONFIG_FILE))
    cfg.MODEL.WEIGHTS = model_zoo.get_checkpoint_url(WEIGHTS_FILE)
    cfg.MODEL.ROI_HEADS.SCORE_THRESH_TEST = 0.5
    cfg.MODEL.DEVICE = 'cpu' # Send to CPU for Neuron Tracing

    # Create the R-CNN predictor wrapper
    predictor = DefaultPredictor(cfg)
    return predictor
```

(continues on next page)

(continued from previous page)

```

def get_image():
    # Get a sample image
    filename = 'input.jpg'
    if not os.path.exists(filename):
        url = "http://images.cocodataset.org/val2017/000000439715.jpg"
        urllib.request.urlretrieve(url, filename)
    return filename

def preprocess(original_image, predictor):
    """
    A basic preprocessing function that sets the input height=800 and
    input width=800. The function is derived from the preprocessing
    steps in the Detectron2 `DefaultPredictor` module.
    """

    height, width = original_image.shape[:2]
    resize_func = predictor.aug.get_transform(original_image)
    resize_func.new_h = 800 # Override height
    resize_func.new_w = 800 # Override width
    image = resize_func.apply_image(original_image)
    image = torch.as_tensor(image.astype("float32").transpose(2, 0, 1))
    inputs = {"image": image, "height": height, "width": width}
    return inputs

# -----
# Neuron modules
# -----

class NeuronFusedBackboneRPNHead(torch.nn.Module):
    """
    Wrapper to compile the fused ResNet backbone and RPN Head.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.backbone = model.backbone
        self.rpn_head = model.proposal_generator.rpn_head
        self.in_features = model.proposal_generator.in_features

    def forward(self, x):
        features = self.backbone(x)
        features_ = [features[f] for f in self.in_features]
        return self.rpn_head(features_), features

class BackboneRPN(torch.nn.Module):
    """
    Wrapper that uses the compiled `neuron_backbone_rpn` instead

```

(continues on next page)

(continued from previous page)

```

of the original backbone and RPN Head. We copy the remainder
of the RPN `forward` code (`predictor.model.proposal_generator.forward`)
to create a "fused" backbone + RPN module.
"""

def __init__(self, model: GeneralizedRCNN) -> None:
    super().__init__()
    self.backbone_rpn_head = NeuronFusedBackboneRPNHead(model)
    self._rpn = model.proposal_generator
    self.in_features = model.proposal_generator.in_features

def forward(self, images):
    preds, features = self.backbone_rpn_head(images.tensor)
    features_ = [features[f] for f in self.in_features]
    pred_objectness_logits, pred_anchor_deltas = preds
    anchors = self._rpn.anchor_generator(features_)

    # Transpose the Hi*Wi*A dimension to the middle:
    pred_objectness_logits = [
        # (N, A, Hi, Wi) -> (N, Hi, Wi, A) -> (N, Hi*Wi*A)
        score.permute(0, 2, 3, 1).flatten(1)
        for score in pred_objectness_logits
    ]
    pred_anchor_deltas = [
        # (N, A*B, Hi, Wi) -> (N, A, B, Hi, Wi) -> (N, Hi, Wi, A, B) -> (N, Hi*Wi*A, B)
        x.view(x.shape[0], -1, self._rpn.anchor_generator.box_dim,
              x.shape[-2], x.shape[-1])
        .permute(0, 3, 4, 1, 2)
        .flatten(1, -2)
        for x in pred_anchor_deltas
    ]

    proposals = self._rpn.predict_proposals(
        anchors, pred_objectness_logits, pred_anchor_deltas, images.image_sizes
    )
    return proposals, features

class NeuronBoxHeadBoxPredictor(torch.nn.Module):
    """
    Wrapper that extracts the RoI Box Head and Box Predictor
    for compilation.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.roi_heads = model.roi_heads

    def forward(self, box_features):
        box_features = self.roi_heads.box_head(box_features)
        predictions = self.roi_heads.box_predictor(box_features)

```

(continues on next page)



(continued from previous page)

```

        return predictions

class ROIHead(torch.nn.Module):
    """
    Wrapper that combines the compiled `roi_heads` into the
    rest of the RoI module. The `_forward_box` and `forward`
    functions are from the `predictor.model.roi_heads` module.
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()
        self.roi_heads = model.roi_heads
        self.neuron_box_head_predictor = NeuronBoxHeadBoxPredictor(model)

    def _forward_box(self, features, proposals):
        features = [features[f] for f in self.roi_heads.box_in_features]
        box_features = self.roi_heads.box_pooler(
            features, [x.proposal_boxes for x in proposals])
        predictions = self.neuron_box_head_predictor(box_features)
        pred_instances, _ = self.roi_heads.box_predictor.inference(
            predictions, proposals)
        return pred_instances

    def forward(self, images, features, proposals, targets=None):
        pred_instances = self._forward_box(features, proposals)
        pred_instances = self.roi_heads.forward_with_given_boxes(
            features, pred_instances)
        return pred_instances, {}

class NeuronRCNN(torch.nn.Module):
    """
    Wrapper that uses the fused backbone + RPN module and the optimized RoI
    Heads wrapper
    """

    def __init__(self, model: GeneralizedRCNN) -> None:
        super().__init__()

        # Create fused Backbone + RPN
        self.backbone_rpn = BackboneRPN(model)

        # Create Neuron RoI Head
        self.roi_heads = ROIHead(model)

        # Define pre and post-processing functions
        self.preprocess_image = model.preprocess_image
        self._postprocess = model._postprocess

    def forward(self, batched_inputs):
        images = self.preprocess_image(batched_inputs)

```

(continues on next page)

(continued from previous page)

```

        proposals, features = self.backbone_rpn(images)
        results, _ = self.roi_heads(images, features, proposals, None)
        return self._postprocess(results, batched_inputs, images.image_sizes)

# -----
# Compilation functions
# -----

def compile(
    model: Union[Callable, torch.nn.Module],
    example_inputs: Any,
    filename: str,
    **kwargs
) -> torch.nn.Module:
    """
    Compiles the model for In1l if it doesn't already exist and saves it as the provided_
    ↪ filename.

    model: A module or function which defines a torch model or computation.
    example_inputs: An example set of inputs which will be passed to the
        `model` during compilation.
    filename: Name of the compiled model
    kwargs: Extra `torch_neuron.trace` kwargs
    """

    if not os.path.exists(filename):
        with torch.no_grad():
            compiled_model = torch_neuron.trace(model, example_inputs, **kwargs)
            torch.jit.save(compiled_model, filename)

# -----
# Benchmarking function
# -----

def benchmark(backbone_rpn_filename, roi_head_filename, inputs,
              n_models=4, batch_size=1, n_threads=4, iterations=200):
    """
    A simple benchmarking function that loads `n_models` optimized
    models onto separate NeuronCores, runs multithreaded inference,
    and calculates the corresponding latency and throughput.
    """

    # Load models
    models = list()
    for i in range(n_models):
        with torch_neuron.experimental.neuron_cores_context(i):
            # Create the RCNN with the fused backbone + RPN Head and compiled RoI Heads
            # Initialize an R-CNN on CPU
            predictor = get_model()

```

(continues on next page)

(continued from previous page)

```

# Create the Neuron R-CNN on CPU
neuron_rcnn = NeuronRCNN(predictor.model)
neuron_rcnn.eval()

# Inject the Neuron compiled models
neuron_rcnn.backbone_rpn.backbone_rpn_head = torch.jit.load(backbone_rpn_
↪filename)
neuron_rcnn.roi_heads.neuron_box_head_predictor = torch.jit.load(roi_head_
↪filename)

models.append(neuron_rcnn)

# Warmup
for _ in range(8):
    for model in models:
        model([inputs])

latencies = []

# Thread task
def task(i):
    start = time.time()
    models[i]([inputs])
    finish = time.time()
    latencies.append((finish - start) * 1000)

begin = time.time()
with ThreadPoolExecutor(max_workers=n_threads) as pool:
    for i in range(iterations):
        pool.submit(task, i % n_models)
end = time.time()

# Compute metrics
boundaries = [50, 95, 99]
names = [f'Latency P{i} (ms)' for i in boundaries]
percentiles = np.percentile(latencies, boundaries)
duration = end - begin

# Display metrics
results = {
    'Samples': iterations,
    'Batch Size': batch_size,
    'Models': n_models,
    'Threads': n_threads,
    'Duration (s)': end - begin,
    'Throughput (inf/s)': (batch_size * iterations) / duration,
    **dict(zip(names, percentiles)),
}

print('-' * 80)
pad = max(map(len, results))
for key, value in results.items():

```

(continues on next page)

(continued from previous page)

```

        if isinstance(value, float):
            print(f'{key + ":" :<{pad + 1}} {value:0.3f}')
        else:
            print(f'{key + ":" :<{pad + 1}} {value}')
    print()

if __name__ == "__main__":

    # Create and compile the combined backbone and RPN Head wrapper
    backbone_rpn_filename = 'backbone_rpn.pt'
    predictor = get_model()
    backbone_rpn_wrapper = NeuronFusedBackboneRPNHead(predictor.model)
    backbone_rpn_wrapper.eval()
    example = torch.rand([1, 3, 800, 800])
    compile(backbone_rpn_wrapper, example, backbone_rpn_filename, strict=False)

    # Create and compile the RoI Head wrapper
    roi_head_filename = 'box_head_predictor.pt'
    predictor = get_model()
    box_head_predictor = NeuronBoxHeadBoxPredictor(predictor.model)
    box_head_predictor.eval()
    example = torch.rand([1000, 256, 7, 7])
    compile(box_head_predictor, example, roi_head_filename)

    # Download a sample image from the COCO dataset and read it
    image_filename = get_image()
    image = cv2.imread(image_filename)
    inputs = preprocess(image, get_model())

    # Benchmark the Neuron R-CNN model for various numbers of loaded models
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=1, n_threads=1)
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=1, n_threads=2)
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=2, n_threads=2)
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=2, n_threads=4)
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=4, n_threads=4)
    benchmark(backbone_rpn_filename, roi_head_filename, inputs, n_models=4, n_threads=8)

```

*This document is relevant for:* Inf1

## PyTorch Neuronx (torch-neuronx)

This document is relevant for: Inf1, Inf2, Trn1, Trn2

### 8.3.6 Graph Partitioner on torch\_neuronx

#### Table of Contents

- *Introduction*
- *The Purpose of the Graph Partitioner*
- *How it Works*
  - *Determining Unsupported Operators*
  - *Customizability*
- *Examples*
  - *Default Usage*
  - *Specifying requirements*
  - *Specifying additional operators to partition*

#### Introduction

This guide introduces the graph partitioner for torch-neuronx. The following sections explain the purpose of the graph partitioner, how it works, and go over a few examples.

#### The Purpose of the Graph Partitioner

While neuronx-cc is very sophisticated and can compile most operators, there are some operator configurations that are not supported by the compiler. Usually in a model that contains unsupported operators, these are only a few operators while the supported parts of the model can benefit from the acceleration benefits that Neuron offers. With this in mind, we developed a graph partitioner that will partition out unsupported operators to be executed on CPU, while compiling and executing the supported operators on Neuron.

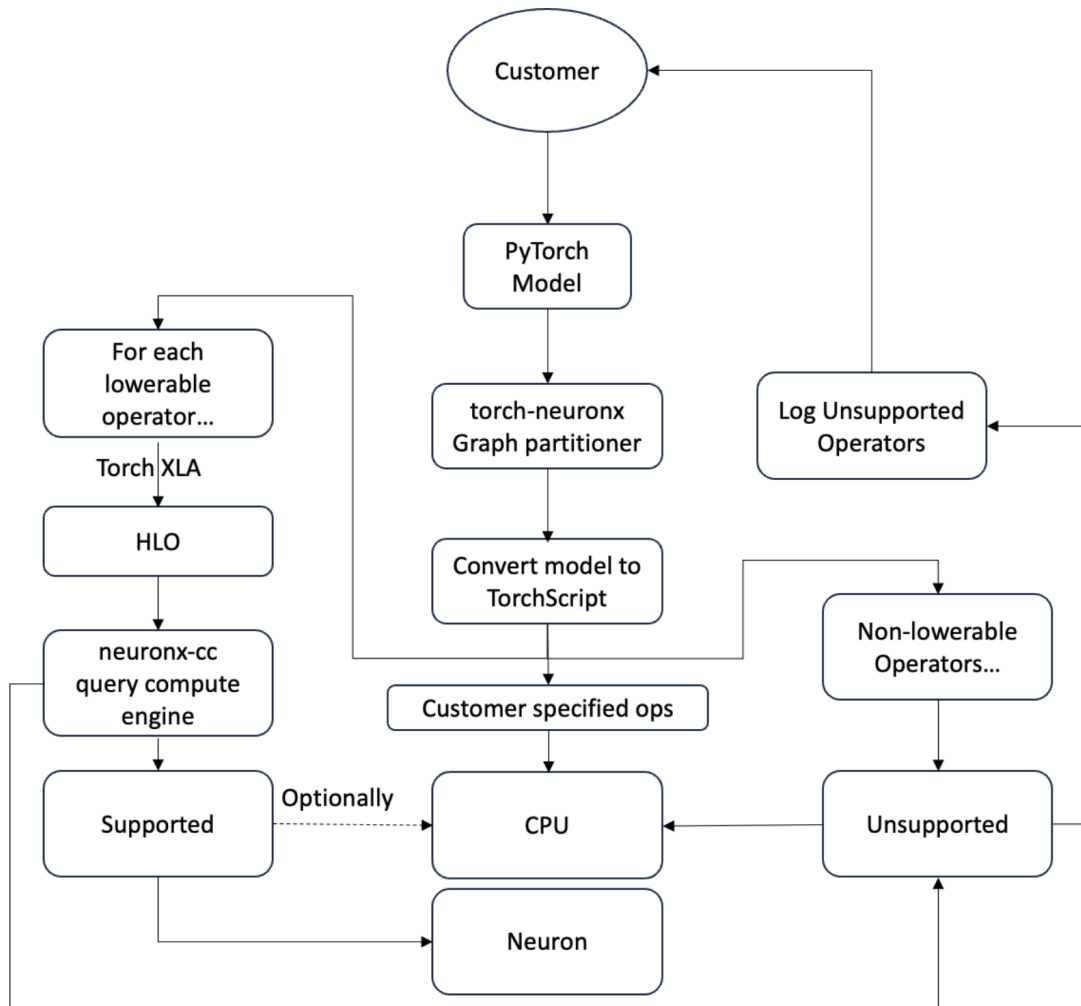
#### How it Works

##### Determining Unsupported Operators

Operator support is determined by the neuronx-cc compiler frontend. This is done because this gives us more flexibility than a static list. This is evident in cases where a specific operator configuration is supported but another configuration is not supported. For example, we support the square root operator, but do not support it with a C64 data type for example.

To check operator support, we use the `torch_neuronx.analyze()` API, which queries the compiler for device placement: Neuron or CPU, which gives the graph partitioner a base graph to start partitioning.

The below image shows the flow of the graph partitioner:



## Customizability

The graph partitioner has a wide range of customizability for a variety of situations. The customization options include:

1. **Minimum Operator Support:** Only partition the model if a minimum percentage of operators are supported.
2. **Minimum Subgraph Size:** The minimum number of operators in any given subgraph. This can be useful if having compute chokepoints with single operator subgraphs is not desired.
3. **Maximum Subgraph Count:** The maximum number of subgraphs. Too many subgraphs can fragment the computation graph causing performance degradation.
4. **Ops to Partition:** Additional operators to partition to CPU beyond the unsupported operators. This can be useful to suggest to the graph partitioner to partition to create a more balanced graph.

Furthermore, compiler flags/args can be passed into all Neuron subgraphs through the graph partitioner.

For the API Reference, visit [torch\\_neuronx.trace\(\)](#) and [torch\\_neuronx.PartitionerConfig](#)

**Note:** Dynamic batching has a case-by-case support with partitioned models, because it is highly dependent on how the final partition scheme looks like.

## Examples

The following sections provide example usages of the graph partitioner.

### Default Usage

The below model is a simple MLP model with sorted log softmax output. The sort operator, `torch.sort()` or `aten::sort`, is not supported by `neuronx-cc` at this time, so the graph partitioner will partition out the sort operator to CPU.

```
import torch
import torch_neuronx
import torch.nn as nn

import logging

# adjust logger level to see what the partitioner is doing
logger = logging.getLogger("Neuron")

class MLP(nn.Module):
    def __init__(
        self, input_size=28 * 28, output_size=10, layers=[4096, 2048]
    ):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)
        self.relu = nn.ReLU()

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = self.relu(f1)
        f2 = self.fc2(r1)
        r2 = self.relu(f2)
        f3 = self.fc3(r2)
        out = torch.log_softmax(f3, dim=1)
        sort_out, _ = torch.sort(out)
        return sort_out

n = MLP()
n.eval()

inputs = torch.rand(32, 784)

# Configure the graph partitioner with the default values
partitioner_config = torch_neuronx.PartitionerConfig()

# Trace a neural network with graph partitioner enabled
neuron_net = torch_neuronx.trace(n, inputs, partitioner_config=partitioner_config)

# Run inference on the partitioned model
output = neuron_net(inputs)
```

## Specifying requirements

This example is very similar to the previous example, but has two differences. The unsupported sort operator is sandwiched between the ReLU activation function after the first linear layer and the second linear layer. The second difference is that we are specifying a max subgraph count of 2.

```
import torch
import torch_neuronx
import torch.nn as nn

import logging

# adjust logger level to see what the partitioner is doing
logger = logging.getLogger("Neuron")

class MLP(nn.Module):
    def __init__(
        self, input_size=28 * 28, output_size=10, layers=[4096, 2048]
    ):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)
        self.relu = nn.ReLU()

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = self.relu(f1)
        sort_r1, _ = torch.sort(r1)
        f2 = self.fc2(sort_r1)
        r2 = self.relu(f2)
        f3 = self.fc3(r2)
        out = torch.log_softmax(f3, dim=1)
        return out

n = MLP()
n.eval()

inputs = torch.rand(32, 784)

# Configure the graph partitioner with the default values
partitioner_config = torch_neuronx.PartitionerConfig(max_subgraph_count=2)

# This trace will fail since the min_subgraph_size requirement can't be satisfied by the
# graph partitioner
neuron_net = torch_neuronx.trace(n, inputs, partitioner_config=partitioner_config)
```

Output:

```
ValueError: The partitioner has found 3 subgraphs which exceeds the specified max_
↳ subgraph count of 2.
```

This example fails because the sort operator placement generates 3 subgraphs, which is more than 2.



## Specifying additional operators to partition

This example shows a situation where we want to partition out the `log_softmax` operator despite it being supported. We also specify an 80% support percentage threshold.

```
import torch
import torch_neuronx
import torch.nn as nn

import logging

# adjust logger level to see what the partitioner is doing
logger = logging.getLogger("Neuron")
logger.setLevel(logging.INFO)

class MLP(nn.Module):
    def __init__(
        self, input_size=28 * 28, output_size=10, layers=[4096, 2048]
    ):
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(input_size, layers[0])
        self.fc2 = nn.Linear(layers[0], layers[1])
        self.fc3 = nn.Linear(layers[1], output_size)
        self.relu = nn.ReLU()

    def forward(self, x):
        f1 = self.fc1(x)
        r1 = self.relu(f1)
        f2 = self.fc2(r1)
        r2 = self.relu(f2)
        f3 = self.fc3(r2)
        out = torch.log_softmax(f3, dim=1)
        sort_out, _ = torch.sort(out)
        return sort_out

n = MLP()
n.eval()

inputs = torch.rand(32, 784)

# Configure the graph partitioner with the default values
partitioner_config = torch_neuronx.PartitionerConfig(min_operator_percentage_threshold=0.
    ↪8, ops_to_partition=set(["aten::log_softmax"]))

# This trace succeeds
neuron_net = torch_neuronx.trace(n, inputs, partitioner_config=partitioner_config)
```

Key Output logs:

```
...
Neuron: The following operations are currently supported:
Neuron: aten::linear
Neuron: aten::relu
```

(continues on next page)

(continued from previous page)

```

Neuron: aten::log_softmax
Neuron: The following operations are currently not supported:
Neuron: aten::sort, unsup.py(28): <stack_trace>
...
Neuron: 85.71% of arithmetic operations (6 of 7) are supported
Neuron: Num Partitions: 2

Neuron: Creating Partition #1 for device: Device.NEURON
Neuron: The following operators will be included in this partition:
Neuron: prim::GetAttr:9
Neuron: aten::linear:3
Neuron: aten::relu:2
...
Neuron: Creating Partition #2 for device: Device.CPU
Neuron: The following operators will be included in this partition:
Neuron: prim::Constant:4
Neuron: aten::sort:1
Neuron: aten::log_softmax:1

```

Notice that we still report that `aten::log_softmax` is still supported, but also report that `aten::log_softmax` is in Partition #2 which is for `Device.CPU`.

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## Transformers Neuron (transformers-neuronx)

*This document is relevant for: Inf2, Trn1, Trn2*

## 8.3.7 Generative LLM inference with Neuron

### Table of contents

- *Background*
- *Performance optimizations*
  - *KV-caching:*
  - *Model sharding:*
  - *Computation/communication overlap:*
  - *Compact data-types:*
  - *Bucketing:*
- *Model partitioning*
  - *How many NeuronCores do I need?*
  - *Which parallelism technique should I use?*
  - *What batch-size should I use?*

## Background

Large Language Models (LLMs) generate human-like text through a process known as generative inference. Fundamentally, given an input prompt, generative LLM inference generates text outputs, by iteratively predicting the next token in a sequence.

These models typically take a sequence of integers as input, which represent a sequence of tokens (words/subwords), and generate a prediction for the next token to be emitted. Below is a simple example that illustrates this in code:

```
# Vocabulary of tokens the model can parse. The position of each token in the
# vocabulary is used as the token_id (an integer representing that token)
vocab = ["having", "I", "fun", "am", "learning", ".", "Neuron"]

# input token_ids: list of integers that represent the input tokens in this
# case: "I", "am", "having", "fun"
input_token_ids = [1, 3, 0, 2]

# The LLM gets a vector of input token_ids, and generates a probability-distribution
# for what the output token_id should be (with a probability score for each token_id
# in the vocabulary)
output = LLM(input_token_ids)

# by taking argmax on the output, we effectively perform a 'greedy sampling' process,
# i.e. we choose the token_id with the highest probability. Other sampling techniques
# also exist, e.g. Top-K. By choosing a probabilistic sampling method we enable the model
# to generate different outputs when called multiple times with the same input.
next_token_id = np.argmax(output)

# map the token_id back into an output token
next_token = vocab[next_token_id]
```

To generate entire sentences, the application iteratively invokes the LLM to generate the next token's prediction, and at each iteration we append the predicted token back into the input:

```
def generate(input_token_ids, n_tokens_to_generate):
    for _ in range(n_tokens_to_generate): # decode loop
        output = LLM(input_token_ids) # model forward pass

        next_token_id = np.argmax(output) # greedy sampling

        if (next_token_id == EOS_TOK_ID)
            break # break if generated End Of Sentence (EOS)

        # append the prediction to the input, and continue to the next out_token
        input_token_ids.append(int(next_token_id))

    return input_token_ids[-n_tokens_to_generate :] # only return generated token_ids

input_token_ids = [1, 3] # "I" "am"
output_token_ids = generate(input_token_ids, 4) # output_token_ids = [0, 2, 4, 6]
output_tokens = [vocab[i] for i in output_token_ids] # "having" "fun" "learning" "Neuron"
```

This process, of predicting a future value (regression) and adding it back into the input (auto), is sometimes referred to as autoregression. For more details, Jay Mody's [GPT in 60 Lines of NumPy](#) is an excellent writeup on GPTs (Generative Pre-trained Transformers).

## Performance optimizations

The sheer size of state-of-the-art LLMs, as well as the sequential nature of text generation, poses multiple challenges for efficient generative LLM deployment.

First, the model is typically sharded across multiple devices, in order to fit the model in device memory. This creates communication overhead and complexity among devices. Secondly, certain deployments have strict application-level latency bounds, thus requiring substantial latency optimizations. This is especially challenging, due to the sequential nature of token-by-token generation. Finally, generating one token at a time often leads to poor device utilization, due to low arithmetic intensity, which can be improved via batching (see [What batch-size should I use?](#)).

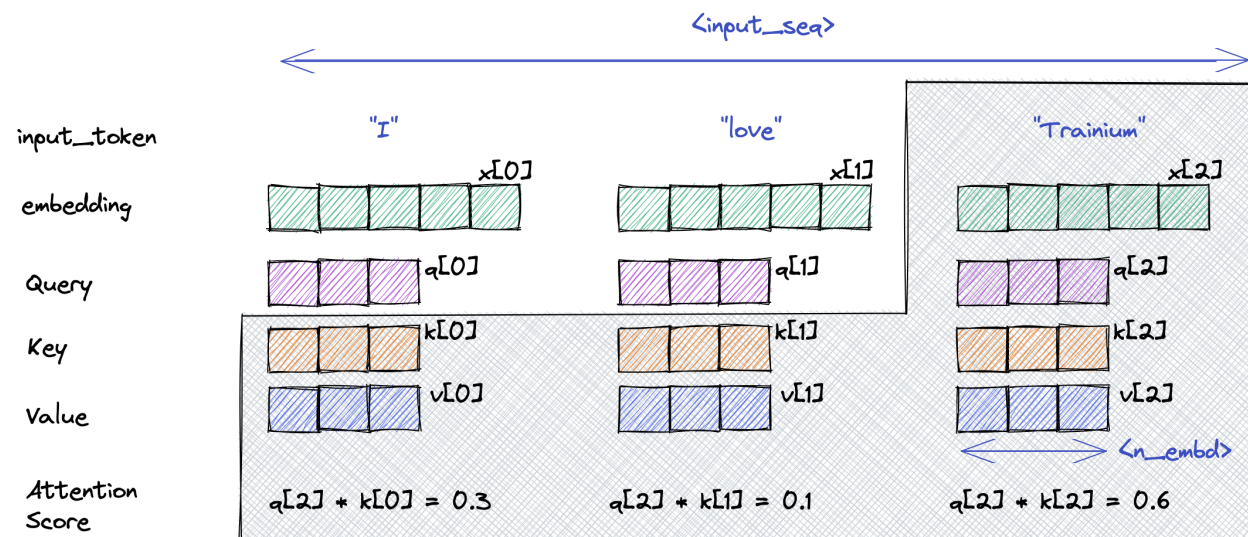
The Neuron SDK provides several built-in optimizations, allowing you to extract optimal performance when deploying LLM models, including:

### KV-caching:

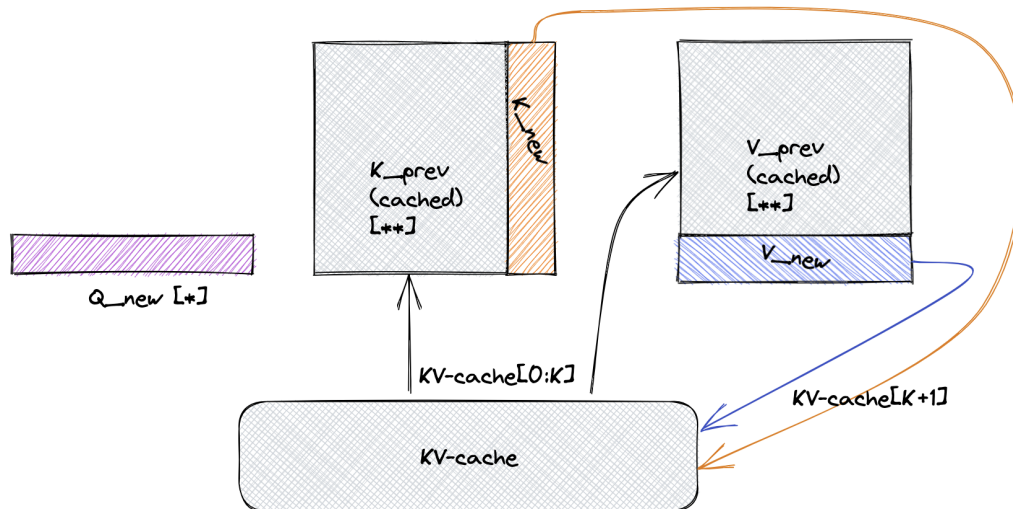
The `transformers-neuronx` library implements KV-cache optimization, which saves compute resources by reusing previously calculated SelfAttention key-value pairs, instead of recalculating them for each generated token.

To illustrate this concept, see the inner workings of the MaskedSelfAttention operator in the figure below.

At each token generation step, the Query vector of a single current token is multiplied by the Key vectors of all previous tokens in the sequence to create attention scores and these scores are further multiplied by the Value vectors of all previous tokens.



The core idea behind this optimization is that instead of re-computing the Key and Value vectors for all previous tokens at each token generation step, Neuron can perform only incremental computation for the current token and reuse previously computed Key/Value vectors from the KV-cache. The Key/Value vector of the current token is also appended to the KV-cache, for the next token generation step.



#### Notes:

- \* When processing token[K], we only need the  $K$ 'th row of  $Q$
- \*\* When processing token[K], we require the full  $K$  &  $V$  tensors, but we can mostly reuse the cached values (This enables skipping the computation of  $K$  &  $V$ )

Note that the first token in the output sequence is unique in two ways:

- No KV-cache is available at this point.
- Neuron needs to compute the entire KV-cache for `<input_len>` tokens (the input prompt), rather than one incremental KV-cache entry.

This means that first-token latency is typically higher than the following tokens.

### Model sharding:

Neuron enables you to shard the model across devices via Tensor Parallelism, Pipeline Parallelism (coming soon), or a combination of the two (coming soon).

Tensor Parallelism shards each layer across multiple devices, enabling you to achieve the optimal latency.

Pipeline Parallelism places different layers on different devices and creates a pipeline between them (as the name suggests) and is useful mainly when optimizing throughput and/or cost-per-inference.

To find the optimal Tensor/Pipeline parallelism configuration for your model, see the [Model partitioning](#) section.

### Computation/communication overlap:

The Neuron compiler automatically fuses Collective Communication primitives (e.g., AllReduce) with the following computation (e.g., GEMM) in the compute graph. This helps minimize any overhead caused by sharding the model across devices.

## Compact data-types:

Neuron supports INT8 and FP8 (coming soon), which can significantly reduce the model's memory bandwidth and capacity requirements. This is especially useful for Generative LLM inference, which is typically memory-bound. Therefore, using a compact data-type can improve the overall LLM inference performance with lower latency and higher throughput.

## Bucketing:

The transformers-neuronx library automatically uses bucketing to process the input prompt and output tokens. Bucketing makes it possible to handle variable sequence lengths, without requiring support for dynamic shapes. Using multiple progressively larger buckets helps minimize the portion of the KV-cache that needs to be read for each token.

## Model partitioning

### How many NeuronCores do I need?

Transformer models are typically defined via a hyper-parameter configuration, such as the following:

```
{
  "n_vocab": 50257, # number of tokens in our vocabulary
  "n_ctx": 2048, # maximum possible sequence length of the input
  "n_embd": 9216, # embedding dimension (determines the "width" of the network)
  "n_head": 72, # number of attention heads (n_embd must be divisible by n_head)
  "n_layer": 64 # number of layers (determines the "depth" of the network)
}
```

To determine the number of NeuronCores needed to fit the model, perform the following calculation:

```
weight_mem_footprint = 12 x <n_layer> x <n_embd>^2 x <dtype-size>
KV_cache_mem_footprint = <batch-size> x <n_layer> x <n_ctx> x <n_embd> x 2 x <dtype-size>
# <dtype-size> is 2 for BF16/FP16, or 1 for FP8/INT8

mem_footprint = weight_mem_footprint + KV_cache_mem_footprint
```

And from here, determining the number of NeuronCores is straightforward:

```
num_neuron_cores = ceil_to_closest_supported_size (mem_footprint / <NC-HBM-capacity>,
↪<instance-type>) # 16GiB per Inferentia2/Trainium1 NeuronCore
```

For example, when running OPT-66B on Inf2, with a batch-size of 16, the number of required NeuronCores can be computed as follows.

```
# OPT-66B example (BF16, Inf2)
# n_layer=64, n_ctx=2048, n_embd=9216, batch=16
weight_mem_footprint = 12 x 64 x 9216^2 x 2 = 121.5 GiB
KV_cache_mem_footprint = 16 x 64 x 2048 x 9216 x 2 x 2 = 72 GiB

mem_footprint = 121.5GiB + 72GiB = 193.5 GiB

num_neuron_cores = ceil_to_closest_supported_size (193.5GiB / 16GiB, Inf2)
                  = ceil_to_closest_supported_size (12.1) = 24
```

(continues on next page)

(continued from previous page)

```
## Currently, the Neuron runtime supports tensor-parallelism degrees 2, 4, 8, and 32 on Trn1
## and supports tensor-parallelism degrees 2, 4, 8, 12 and 24 on Inf2.
```

Use the *Neuron Calculator* to compute the number of cores needed for a custom hyper-parameter configuration.

### Which parallelism technique should I use?

Tensor parallelism improves latency, at the expense of increased intra-layer communication. Thus, as a general rule, it is recommended to use the smallest tensor parallelism degree that meets your latency requirement and then use pipeline/data parallelism from that point on.

If latency is not a major concern in your application (e.g., model evaluation) and the primary goal is to maximize throughput (i.e., minimize total cost per token), then it is most efficient to use pipeline parallelism and increase the batch-size as much as possible.

### What batch-size should I use?

Due to the serial token generation nature of generative LLM inference, this workload tends to be extremely memory bound. This means that throughput (and thus cost per inference) improves significantly by batching.

As a general rule, we recommend increasing the batch-size to the maximum amount that fits within the latency budget (up to batch=256. A larger batch-size typically does not help with performance.)

Note that the KV-cache grows linearly with the batch-size and can grow until it runs out of memory (typically referred to as OOM). If the latency budget allows, we recommend increasing the batch-size to the maximum value that does not result in OOM.

Users may also consider pipelining the model beyond what is necessary to fit model parameters / KV-cache on devices, in order to free up device-memory space and thus allow the batch-size to increase without causing OOM issues.

*This document is relevant for: Inf2, Trn1, Trn2*

## PyTorch 2.x

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 8.3.8 Introducing PyTorch 2.7 Support

#### Table of contents

- *What are we introducing?*
- *How is PyTorch NeuronX 2.7 different compared to PyTorch NeuronX 2.5?*
- *How can I install PyTorch NeuronX 2.7?*
- *Migrate your application to PyTorch 2.7*
  - *Migrating training scripts*
  - *Migrating inference scripts*



- *Troubleshooting and Known Issues*
  - *TypeError: AdamW.\_\_init\_\_() got an unexpected keyword argument 'decoupled\_weight\_decay'*
  - *Tensor split on second dimension of 2D array not working*
  - *Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`*
  - *Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.6 release, please downcast your model directly”*
  - *AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'xrt\_world\_size'*
  - *AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'get\_ordinal'*
  - *AttributeError: <module 'torch\_xla.core.xla\_model' ... does not have the attribute 'get\_local\_ordinal'*
  - *Socket Error: Socket failed to bind*
  - *AttributeError: module 'torch' has no attribute 'xla' Failure*
  - *Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API*
  - *ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023*
  - *FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path' Failure*
  - *Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile*
  - *Compiler assertion error when running Stable Diffusion training*
- *Frequently Asked Questions (FAQ)*
  - *Do I need to recompile my models with PyTorch 2.7?*
  - *Do I need to update my scripts for PyTorch 2.7?*
  - *What environment variables will be changed with PyTorch NeuronX 2.7 ?*
  - *What features will be missing with PyTorch NeuronX 2.7?*
  - *Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.7?*
  - *Can I still use PyTorch 2.6 version?*
  - *Can I still use PyTorch 2.5 version?*
  - *Can I still use PyTorch 2.1 version?*



## What are we introducing?

Starting with the [Neuron 2.24](#) release, customers will be able to upgrade to PyTorch NeuronX(torch-neuronx) supporting PyTorch 2.7.

[PyTorch Neuron \(torch-neuronx\) Setup](#) is updated to include installation instructions for PyTorch NeuronX 2.7 for Amazon Linux 2023 and Ubuntu 22. Note that PyTorch NeuronX 2.7 is supported on Python 3.9, 3.10, and 3.11.

Please review [migration guide](#) for possible changes to training scripts. No code changes are required for inference scripts.

## How is PyTorch NeuronX 2.7 different compared to PyTorch NeuronX 2.5?

PyTorch NeuronX 2.7 uses Torch-XLA v2.7 and PyTorch v2.7 which have C++11 ABI enabled by default.

Additionally, Torch-XLA v2.7 includes a fix for training performance issue <https://github.com/pytorch/xla/issues/9037>.

See [Torch-XLA 2.7 release](#) for a full list.

See [Migrate your application to PyTorch 2.7](#) for changes needed to use PyTorch NeuronX 2.7.

---

**Note:** GSPMD and Torch Dynamo (torch.compile) support in Neuron will be available in a future release.

---

## How can I install PyTorch NeuronX 2.7?

To install PyTorch NeuronX 2.7 please follow the [PyTorch Neuron \(torch-neuronx\) Setup](#) guides for Amazon Linux 2023 and Ubuntu 22 AMI. Please also refer to the Neuron multi-framework DLAMI setup guide for Ubuntu 22 with a pre-installed virtual environment for PyTorch NeuronX 2.7 that you can use to get started. PyTorch NeuronX 2.7 can be installed using the following:

```
python -m pip install --upgrade neuronx-cc==2.* torch-neuronx==2.7.* torchvision
```

---

**Note:** PyTorch NeuronX 2.7 is currently available for Python 3.9, 3.10, 3.11.

---

## Migrate your application to PyTorch 2.7

Please make sure you have first installed the PyTorch NeuronX 2.7 as described above in [installation guide](#)

### Migrating training scripts

To migrate the training scripts from PyTorch NeuronX 2.5/2.6 to PyTorch NeuronX 2.7, implement the following changes:

---

**Note:** `xm` below refers to `torch_xla.core.xla_model`, `xr` refers to `torch_xla.runtime`, and `xmp` refers to `torch_xla.distributed.xla_multiprocessing`

---

- The environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used) and will be removed in an upcoming release. Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to convert model to BF16 format. (see [migration\\_from\\_xla\\_downcast\\_bf16](#))
- The functions `xm.xrt_world_size()`, `xm.xla_model.get_ordinal()`, and `xm.xla_model.get_local_ordinal()` are deprecated and removed so there's error when used. Please switch to `xr.world_size()`, `xr.global_ordinal()`, and `xr.local_ordinal()` respectively as replacements.
- The default behavior of `torch.load` parameter `weights_only` is changed from `False` to `True`. Leaving `weights_only` as `True` can cause issues with pickling.
- If using `xmp.spawn`, the `nprocs` argument limited to 1 or `None` since v2.1. Previously, passing a value `> 1` would result in a warning. In torch-xla 2.6+, passing a value `> 1` would result in an error with an actionable message to use `NEURON_NUM_DEVICES` to set the number of NeuronCores to use.

See [v2.6 migration guide](#) for additional changes needed if you are migrating from PyTorch NeuronX 2.5. See [v2.5 migration guide](#) for additional changes needed if you are migrating from PyTorch NeuronX 2.1.

## Migrating inference scripts

There are no code changes required in the inference scripts.

## Troubleshooting and Known Issues

### **TypeError: AdamW.\_\_init\_\_() got an unexpected keyword argument 'decoupled\_weight\_decay'**

AdamW now has an additional argument “`decoupled_weight_decay`” which is default to `False`. If you get “`TypeError: AdamW.__init__() got an unexpected keyword argument 'decoupled_weight_decay'`” with NeuronX Distributed, please update to the latest version.

### **Tensor split on second dimension of 2D array not working**

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`. Another work-around is to use `torch.tensor_split`.

### **Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`**

Currently, BERT pretraining performance is ~11% lower when switching to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16` due to <https://github.com/pytorch/xla/issues/8545>. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 and 2.7 although there will be deprecation warnings (as noted below).

### Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.6 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see migration\_from\_xla\_downcast\_bf16)

### AttributeError: <module ‘torch\_xla.core.xla\_model’ ... does not have the attribute ‘xrt\_world\_size’

This is an error that `torch_xla.core.xla_model.xrt_world_size()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.world_size()` instead. If using Hugging Face transformers/accelerate libraries, please use `transformers==4.53.*` and `accelerate==1.7.*`.

### AttributeError: <module ‘torch\_xla.core.xla\_model’ ... does not have the attribute ‘get\_ordinal’

This is an error that `torch_xla.core.xla_model.xla_model.get_ordinal()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.global_ordinal()` instead. If using Hugging Face transformers/accelerate libraries, please use `transformers==4.53.*` and `accelerate==1.7.*`.

### AttributeError: <module ‘torch\_xla.core.xla\_model’ ... does not have the attribute ‘get\_local\_ordinal’

This is an error that `torch_xla.core.xla_model.xla_model.get_local_ordinal()` is removed in torch-xla version 2.7. Please switch to using `torch_xla.runtime.local_ordinal()` instead. If using Hugging Face transformers/accelerate libraries, please use `transformers==4.53.*` and `accelerate==1.7.*`.

### Socket Error: Socket failed to bind

In PyTorch 2.7, there needs to be a socket available for both `torchrun` and the `init_process_group` to bind. Both of these, by default, will be set to unused sockets. If you plan to use a `MASTER_PORT` environment variable then this error may occur, if the port you set it to is already in use.

```
[W socket.cpp:426] [c10d] The server socket has failed to bind to [::]:2.700 (errno: 98 -
↪ Address already in use).
[W socket.cpp:426] [c10d] The server socket has failed to bind to ?UNKNOWN? (errno: 98 -
↪ Address already in use).
[E socket.cpp:462] [c10d] The server socket has failed to listen on any local network
↪ address.
RuntimeError: The server socket has failed to listen on any local network address.
The server socket has failed to bind to ?UNKNOWN? (errno: 98 - Address already in use).
```

To resolve the issue, please ensure if you are setting `MASTER_PORT` that the port you’re setting it to is not used anywhere else in your scripts. Otherwise, you can leave `MASTER_PORT` unset, and `torchrun` will set the default port for you.

**AttributeError: module 'torch' has no attribute 'xla' Failure**

In PyTorch 2.7, training scripts might fail during activation checkpointing with the error shown below.

```
AttributeError: module 'torch' has no attribute 'xla'
```

The solution is to use `torch_xla.utils.checkpoint.checkpoint` instead of `torch.utils.checkpoint.checkpoint` as the checkpoint function while wrapping pytorch modules for activation checkpointing. Refer to the [pytorch/xla discussion](#) regarding this [issue](#). Also set `use_reentrant=True` while calling the `torch_xla` checkpoint function. Failure to do so will lead to `XLA currently does not support use_reentrant==False` error. For more details on checkpointing, refer the [documentation](#).

**Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API**

While using HuggingFace Transformers Trainer API to train (i.e. [HuggingFace Trainer API fine-tuning tutorial](#)), you may see the error “Attempted to access the data pointer on an invalid python storage”. This is a known [issue](#) and has been fixed in the version 4.37.3 of HuggingFace Transformers.

**ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023**

torch-xla version 2.5+ now requires `libcrypt.so.1` shared library. Currently, Amazon Linux 2023 includes `libcrypt.so.2` shared library by default so you may see `ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory` when using torch-neuronx 2.1+ on Amazon Linux 2023. To install `libcrypt.so.1` on Amazon Linux 2023, please run the following installation command (see also <https://github.com/amazonlinux/amazon-linux-2023/issues/182> for more context):

```
sudo yum install libxcrypt-compat
```

**FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path' Failure**

In PyTorch 2.7, users might face the error shown below due to incompatible `libneuronxla` and `torch-neuronx` versions being installed.

```
FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path'
```

Check that the version of `libneuronxla` that support PyTorch NeuronX 2.7 is 2.2.\*. If not, then uninstall `libneuronxla` using `pip uninstall libneuronxla` and then reinstall the packages following the [installation guide](#)

**Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile**

When running Neuron Parallel Compile with HF Trainer API, you may see the error `Status: INVALID_ARGUMENT: Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range` in Accelerator's `pad_across_processes` function. This is due to data-dependent operation in evaluation metrics computation. Data-dependent operations would result in undefined behavior with Neuron Parallel Compile trial execution (execute empty graphs with zero outputs). To work-around this error, please disable `compute_metrics` when `NEURON_EXTRACT_GRAPHS_ONLY` is set to 1:

```
compute_metrics=None if os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY") else compute_metrics
```

### Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.7 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.23, please disable gradient accumulation in torch-neuronx 2.7.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.9/concurrent/futures/
↳process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]
```

### Frequently Asked Questions (FAQ)

#### Do I need to recompile my models with PyTorch 2.7?

Yes.

#### Do I need to update my scripts for PyTorch 2.7?

Please see the [migration guide](#)

#### What environment variables will be changed with PyTorch NeuronX 2.7 ?

The environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

### What features will be missing with PyTorch NeuronX 2.7?

PyTorch NeuronX 2.7 has all of the supported features in PyTorch NeuronX 2.6, with known issues listed above, and unsupported features as listed in *PyTorch Neuron (torch-neuronx) release notes*.

### Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.7?

Yes, NeuronX Distributed and Transformers NeuronX are supported by PyTorch NeuronX 2.7. AWS Neuron Reference for NeMo Megatron has reached end-of-support in release 2.23.

### Can I still use PyTorch 2.6 version?

PyTorch 2.6 is supported since release 2.23.

### Can I still use PyTorch 2.5 version?

PyTorch 2.5 is supported for releases 2.21 to 2.24 and will reach end-of-life in a future release. Additionally, the CVE [CVE-2025-32434](#) affects PyTorch version 2.5. We recommend upgrading to the new version of Torch-NeuronX by following *PyTorch Neuron (torch-neuronx) Setup*.

### Can I still use PyTorch 2.1 version?

PyTorch 2.1 is supported for release 2.21 and has reached end-of-life in release 2.22. Additionally, the CVEs [CVE-2024-31583](#) and [CVE-2024-31580](#) affect PyTorch versions 2.1 and earlier. We recommend upgrading to the new version of Torch-NeuronX by following *PyTorch Neuron (torch-neuronx) Setup*.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 8.3.9 Introducing PyTorch 2.6 Support

### Table of contents

- *What are we introducing?*
- *How is PyTorch NeuronX 2.6 different compared to PyTorch NeuronX 2.5?*
- *How can I install PyTorch NeuronX 2.6?*
- *Migrate your application to PyTorch 2.6*
  - *Migrating training scripts*
  - *Migrating inference scripts*
- *Troubleshooting and Known Issues*
  - *Tensor split on second dimension of 2D array not working*
  - *Lower BERT pretraining performance with torch-neuronx 2.6 compared to torch-neuronx 2.5*

- Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`
- Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.6 release, please downcast your model directly”
- WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.
- WARNING:torch\_xla.core.xla\_model.xla\_model.get\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.global\_ordinal instead.
- WARNING:torch\_xla.core.xla\_model.xla\_model.get\_local\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.local\_ordinal instead.
- Socket Error: Socket failed to bind
- AttributeError: module 'torch' has no attribute 'xla' Failure
- Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API
- ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023
- FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path' Failure
- Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile
- Compiler assertion error when running Stable Diffusion training
- Frequently Asked Questions (FAQ)
  - Do I need to recompile my models with PyTorch 2.6?
  - Do I need to update my scripts for PyTorch 2.6?
  - What environment variables will be changed with PyTorch NeuronX 2.6 ?
  - What features will be missing with PyTorch NeuronX 2.6?
  - Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.6?
  - Can I still use PyTorch 2.5 version?
  - Can I still use PyTorch 2.1 version?

## What are we introducing?

Starting with the Neuron 2.23 release, customers will be able to upgrade to PyTorch NeuronX(`torch-neuronx`) supporting PyTorch 2.6.

*PyTorch Neuron (torch-neuronx) Setup* is updated to include installation instructions for PyTorch NeuronX 2.6 for Amazon Linux 2023 and Ubuntu 22. Note that PyTorch NeuronX 2.6 is supported on Python 3.9, 3.10, and 3.11.

Please review *migration guide* for possible changes to training scripts. No code changes are required for inference scripts.

## How is PyTorch NeuronX 2.6 different compared to PyTorch NeuronX 2.5?

PyTorch NeuronX 2.6 uses Torch-XLA 2.6 which has improved support for Automatic Mixed Precision and buffer aliasing. Additionally:

- Reintroduced `XLA_USE_32BIT_LONG` to give customers the flexibility to use INT32 for their workloads. This flag was removed in v2.5.
- Added `xm.xla_device_kind()` to return XLA device kind string ('NC\_v2' for Trainium1, 'NC\_v3' and 'NC\_v3d' for Trainium2). See [Logical NeuronCore configuration](#) for more info.

See [Torch-XLA 2.6 release](#) for a full list.

See [Migrate your application to PyTorch 2.6](#) for changes needed to use PyTorch NeuronX 2.6.

---

**Note:** GSPMD and Torch Dynamo (`torch.compile`) support in Neuron will be available in a future release.

---

## How can I install PyTorch NeuronX 2.6?

To install PyTorch NeuronX 2.6 please follow the [PyTorch Neuron \(torch-neuronx\) Setup](#) guides for Amazon Linux 2023 and Ubuntu 22 AMI. Please also refer to the Neuron multi-framework DLAMI setup guide for Ubuntu 22 with a pre-installed virtual environment for PyTorch NeuronX 2.6 that you can use to get started. PyTorch NeuronX 2.6 can be installed using the following:

```
python -m pip install --upgrade neuronx-cc==2.* torch-neuronx==2.6.* torchvision
```

---

**Note:** PyTorch NeuronX 2.6 is currently available for Python 3.9, 3.10, 3.11.

---

## Migrate your application to PyTorch 2.6

Please make sure you have first installed the PyTorch NeuronX 2.6 as described above in [installation guide](#)

### Migrating training scripts

To migrate the training scripts from PyTorch NeuronX 2.5 to PyTorch NeuronX 2.6, implement the following changes:

---

**Note:** `xm` below refers to `torch_xla.core.xla_model`, `xr` refers to `torch_xla.runtime`, and `xmp` refers to `torch_xla.distributed.xla_multiprocessing`

---

- The environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used) and will be removed in an upcoming release. Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to convert model to BF16 format. (see `migration_from_xla_downcast_bf16`)
- The functions `xm.xrt_world_size()`, `xm.xla_model.get_ordinal()`, and `xm.xla_model.get_local_ordinal()` are deprecated (warning when used). Please switch to `xr.world_size()`, `xr.global_ordinal()`, and `xr.local_ordinal()` respectively as replacements.
- The default behavior of `torch.load` parameter `weights_only` is changed from `False` to `True`. Leaving `weights_only` as `True` can cause issues with pickling.



- If using `xmp.spawn`, the `nprocs` argument limited to 1 or None since v2.1. Previously, passing a value > 1 would result in a warning. In torch-xla 2.6, passing a value > 1 would result in an error with an actionable message to use `NEURON_NUM_DEVICES` to set the number of NeuronCores to use.

See v2.5 migration guide for additional changes needed if you are migrating from PyTorch NeuronX 2.1.

## Migrating inference scripts

There are no code changes required in the inference scripts.

## Troubleshooting and Known Issues

### Tensor split on second dimension of 2D array not working

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`. Another work-around is to use `torch.tensor_split`.

### Lower BERT pretraining performance with torch-neuronx 2.6 compared to torch-neuronx 2.5

Currently, BERT pretraining performance is ~10% lower with torch-neuronx 2.6 compared to torch-neuronx 2.5. This is due to a known regression in torch-xla <https://github.com/pytorch/xla/issues/9037> and can affect other models with high graph tracing overhead. To work-around this issue, please build the `r2.6_aws_neuron` branch of torch-xla as follows (see `pytorch-neuronx-install-cxx11` for C++11 ABI version):

```
# Setup build env (make sure you are in a python virtual env). Replace "apt" with "yum"
↳ on AL2023.
sudo apt install cmake
pip install yapf==0.30.0
wget https://github.com/bazelbuild/bazelisk/releases/download/v1.20.0/bazelisk-linux-
↳ amd64
sudo cp bazelisk-linux-amd64 /usr/local/bin/bazel
# Clone repos
git clone --recursive https://github.com/pytorch/pytorch --branch v2.6.0
cd pytorch/
git clone --recursive https://github.com/pytorch/xla.git --branch r2.6_aws_neuron
_GLIBCXX_USE_CXX11_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist
cd xla/
CXX_ABI=0 python setup.py bdist_wheel
# pip wheel will be present in ./dist and can be installed instead of the torch-xla
↳ released in pypi.org
```

### Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`

Currently, BERT pretraining performance is ~11% lower when switching to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16` due to <https://github.com/pytorch/xla/issues/8545>. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 and 2.6 although there will be deprecation warnings (as noted below).

### Warning “`XLA_DOWNCAST_BF16` will be deprecated after the 2.6 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see migration\_from\_xla\_downcast\_bf16)

### **WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.**

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size` instead.

### **WARNING:torch\_xla.core.xla\_model.xla\_model.get\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.global\_ordinal instead.**

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal` instead.

### **WARNING:torch\_xla.core.xla\_model.xla\_model.get\_local\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.local\_ordinal instead.**

This is a warning that `torch_xla.core.xla_model.xla_model.get_local_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.local_ordinal` instead.

### Socket Error: Socket failed to bind

In PyTorch 2.6, there needs to be a socket available for both `torchrun` and the `init_process_group` to bind. Both of these, by default, will be set to unused sockets. If you plan to use a `MASTER_PORT` environment variable then this error may occur, if the port you set it to is already in use.

```
[W socket.cpp:426] [c10d] The server socket has failed to bind to [::]:2.600 (errno: 98 -  
→ Address already in use).  
[W socket.cpp:426] [c10d] The server socket has failed to bind to ?UNKNOWN? (errno: 98 -  
→ Address already in use).  
[E socket.cpp:462] [c10d] The server socket has failed to listen on any local network  
→ address.  
RuntimeError: The server socket has failed to listen on any local network address.  
The server socket has failed to bind to ?UNKNOWN? (errno: 98 - Address already in use).
```

To resolve the issue, please ensure if you are setting `MASTER_PORT` that the port you’re setting it to is not used anywhere else in your scripts. Otherwise, you can leave `MASTER_PORT` unset, and `torchrun` will set the default port for you.

**AttributeError: module 'torch' has no attribute 'xla' Failure**

In PyTorch 2.6, training scripts might fail during activation checkpointing with the error shown below.

```
AttributeError: module 'torch' has no attribute 'xla'
```

The solution is to use `torch_xla.utils.checkpoint.checkpoint` instead of `torch.utils.checkpoint.checkpoint` as the checkpoint function while wrapping pytorch modules for activation checkpointing. Refer to the [pytorch/xla discussion](#) regarding this [issue](#). Also set `use_reentrant=True` while calling the `torch_xla` checkpoint function. Failure to do so will lead to `XLA currently does not support use_reentrant==False` error. For more details on checkpointing, refer the [documentation](#).

**Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API**

While using HuggingFace Transformers Trainer API to train (i.e. [HuggingFace Trainer API fine-tuning tutorial](#)), you may see the error “Attempted to access the data pointer on an invalid python storage”. This is a known [issue](#) and has been fixed in the version 4.37.3 of HuggingFace Transformers.

**ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023**

torch-xla version 2.6+ now requires `libcrypt.so.1` shared library. Currently, Amazon Linux 2023 includes `libcrypt.so.2` shared library by default so you may see *ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory* when using torch-neuronx 2.1+ on Amazon Linux 2023. To install `libcrypt.so.1` on Amazon Linux 2023, please run the following installation command (see also <https://github.com/amazonlinux/amazon-linux-2023/issues/182> for more context):

```
sudo yum install libxcrypt-compat
```

**FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path' Failure**

In PyTorch 2.6, users might face the error shown below due to incompatible `libneuronxla` and `torch-neuronx` versions being installed.

```
FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path'
```

Check that the version of `libneuronxla` that support PyTorch NeuronX 2.6 is 2.2.\*. If not, then uninstall `libneuronxla` using `pip uninstall libneuronxla` and then reinstall the packages following the [installation guide](#)

**Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile**

When running Neuron Parallel Compile with HF Trainer API, you may see the error `Status: INVALID_ARGUMENT: Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range` in Accelerator's `pad_across_processes` function. This is due to data-dependent operation in evaluation metrics computation. Data-dependent operations would result in undefined behavior with Neuron Parallel Compile trial execution (execute empty graphs with zero outputs). To work-around this error, please disable `compute_metrics` when `NEURON_EXTRACT_GRAPHS_ONLY` is set to 1:

```
compute_metrics=None if os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY") else compute_metrics
```

### Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.6 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.23, please disable gradient accumulation in torch-neuronx 2.6.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.9/concurrent/futures/
↳process.py at line 239 with exception:
too many partition dims! {{0,+,960}[10],+,10560}[10]
```

### Frequently Asked Questions (FAQ)

#### Do I need to recompile my models with PyTorch 2.6?

Yes.

#### Do I need to update my scripts for PyTorch 2.6?

Please see the [migration guide](#)

#### What environment variables will be changed with PyTorch NeuronX 2.6 ?

The environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

### What features will be missing with PyTorch NeuronX 2.6?

PyTorch NeuronX 2.6 has all of the supported features in PyTorch NeuronX 2.5, with known issues listed above, and unsupported features as listed in *PyTorch Neuron (torch-neuronx) release notes*.

### Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.6?

Yes, NeuronX Distributed, and Transformers NeuronX, and AWS Neuron Reference for NeMo Megatron libraries will work with PyTorch NeuronX 2.6.

### Can I still use PyTorch 2.5 version?

PyTorch 2.5 is supported for releases 2.21/2.22/2.23 and will reach end-of-life in a future release. Additionally, the CVE [CVE-2025-32434](#) affects PyTorch version 2.5. We recommend upgrading to the new version of Torch-NeuronX by following *PyTorch Neuron (torch-neuronx) Setup*.

### Can I still use PyTorch 2.1 version?

PyTorch 2.1 is supported for release 2.21 and has reached end-of-life in release 2.22. Additionally, the CVEs [CVE-2024-31583](#) and [CVE-2024-31580](#) affect PyTorch versions 2.1 and earlier. We recommend upgrading to the new version of Torch-NeuronX by following *PyTorch Neuron (torch-neuronx) Setup*.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 8.3.10 Introducing PyTorch 2.5 Support

### Table of contents

- *What are we introducing?*
- *How is PyTorch NeuronX 2.5 different compared to PyTorch NeuronX 2.1?*
- *How can I install PyTorch NeuronX 2.5?*
- *Migrate your application to PyTorch 2.5*
  - *Migrating training scripts*
  - *Migrating inference scripts*
- *Troubleshooting and Known Issues*
  - *Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with OSError when the Neuron Cache is placed on FSx mount*
  - *Running in-place update operations (e.g. all\_reduce) on 0-dimensional tensors result in buffer aliasing errors in torch 2.5 and earlier*
  - *Tensor split on second dimension of 2D array not working*
  - *Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)*

- *Certain sequence of operations with `xm.save()` could corrupt tensors*
- *Lower BERT pretraining performance when switch to using `model.to(torch.bfloat16)`*
- *Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.5 release, please downcast your model directly”*
- *WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.*
- *WARNING:torch\_xla.core.xla\_model.xla\_model.get\_ordinal() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.global\_ordinal instead.*
- *AttributeError: module ‘torch\_xla.runtime’ has no attribute ‘using\_pjrt’*
- *Socket Error: Socket failed to bind*
- *AttributeError: module ‘torch’ has no attribute ‘xla’ Failure*
- *Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API*
- *ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023*
- *FileNotFoundError: [Errno 2] No such file or directory: ‘libneuronpjrt-path’ Failure*
- *Glibc error on Amazon Linux 2*
- *Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile*
- *Compiler assertion error when running Stable Diffusion training*
- *Frequently Asked Questions (FAQ)*
  - *Do I need to recompile my models with PyTorch 2.5?*
  - *Do I need to update my scripts for PyTorch 2.5?*
  - *What environment variables will be changed with PyTorch NeuronX 2.5 ?*
  - *What features will be missing with PyTorch NeuronX 2.5?*
  - *Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.5?*
  - *Can I still use PyTorch 2.1 version?*

## What are we introducing?

Starting with the Neuron 2.21 release, customers will be able to upgrade to PyTorch NeuronX(`torch-neuronx`) supporting PyTorch 2.5.

*PyTorch Neuron (torch-neuronx) Setup* is updated to include installation instructions for PyTorch NeuronX 2.5 for Amazon Linux 2023 and Ubuntu 22. Note that PyTorch NeuronX 2.5 does not support Python 3.8 which is default in Ubuntu 20. To use Ubuntu 20, customers will need to install Python 3.9+.

Please review *migration guide* for possible changes to training scripts. No code changes are required for inference scripts.

## How is PyTorch NeuronX 2.5 different compared to PyTorch NeuronX 2.1?

PyTorch NeuronX 2.5 uses Torch-XLA 2.5 which has improved support for eager debug mode, Automatic Mixed Precision, PJRT device auto-detection, FP8, and others. See [Torch-XLA 2.5 release](#) for a full list.

See [Migrate your application to PyTorch 2.5](#) for changes needed to use PyTorch NeuronX 2.5.

---

**Note:** GSPMD and Torch Dynamo (torch.compile) support in Neuron will be available in a future release.

---

## How can I install PyTorch NeuronX 2.5?

To install PyTorch NeuronX 2.5 please follow the [PyTorch Neuron \(torch-neuronx\) Setup](#) guides for Amazon Linux 2023 and Ubuntu 22 AML. Please also refer to the Neuron multi-framework DLAMI setup guide for Ubuntu 22 with a pre-installed virtual environment for PyTorch NeuronX 2.5 that you can use to get started. PyTorch NeuronX 2.5 can be installed using the following:

```
python -m pip install --upgrade neuronx-cc==2.* torch-neuronx==2.5.* torchvision
```

---

**Note:** PyTorch NeuronX 2.5 is currently available for Python 3.9, 3.10, 3.11.

---

## Migrate your application to PyTorch 2.5

Please make sure you have first installed the PyTorch NeuronX 2.5 as described above in [installation guide](#)

### Migrating training scripts

To migrate the training scripts from PyTorch NeuronX 2.1 to PyTorch NeuronX 2.5, implement the following changes:

---

**Note:** `xm` below refers to `torch_xla.core.xla_model` and `xr` refers to `torch_xla.runtime`

---

- The environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to convert model to BF16 format. (see [migration\\_from\\_xla\\_downcast\\_bf16](#))
- The `torch_xla.experimental.pjrt` module which was replaced by `torch_xla.runtime` in Torch-XLA 2.1, has been removed in Torch-XLA 2.5. Users should now utilize the `torch_xla.runtime` module as a replacement.
- `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime.
- `xm.all_reduce` no longer operates in-place for single tensors. To fix this, please convert the single tensor to an array (e.g., `[single_tensor]`) or assign the output of `xm.all_reduce` to a variable.
- The functions `xm.xrt_world_size()`, `xm.xla_model.get_ordinal()`, and `xm.xla_model.get_local_ordinal()` are deprecated (warning when used). Please switch to `xr.world_size`, `xr.global_ordinal`, and `xr.local_ordinal` respectively as replacements.
- `torch_xla.experimental.xla_sharding` is now replaced by `torch_xla.distributed.spmd.xla_sharding`.

- Class `ZeroRedundancyOptimizer` now has two new arguments that replaces the optional boolean argument `coalesce_cc`:
  - `bucket_cap_mb_all_gather` (int, Optional): Number of MegaBytes of the tensor bucket to fill before doing all-gather. Default: 0 (disable all gather coalescing).
  - `bucket_cap_mb_reduce_scatter` (int, Optional): Number of MegaBytes of the tensor bucket to fill before doing reduce-scatter. Default: 0 (disable reduce scatter coalescing).

## Migrating inference scripts

There are no code changes required in the inference scripts.

## Troubleshooting and Known Issues

### Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with `OSError` when the Neuron Cache is placed on FSx mount

Currently, the Neuronx-Distributed Training Llama 3.1 70B 8-node tutorial failed with `OSError` (Errno 61) when the Neuron Cache is placed on FSx mount:

```
[rank197]: RuntimeError: Bad StatusOr access: INVALID_ARGUMENT: RunNeuronCCImpl: error_
↳ condition !(error != 400): <class 'OSError'>: [Errno 61] No data available: '/fsx1/
↳ neuron_cache/neuronxcc-2.16.372.0+4a9b2326/MODULE_3540044791706521849+4eb52b03/model.
↳ neff' -> '/tmp/tmpx7bvfpmm/model.neff'
```

We found that the error is due to FSx failing during file copy when there are multiple readers (13 workers fail to copy out of 256). This issue doesn't affect simpler models like BERT.

To work-around the issue, please use the shared NFS mount (/home directory on a Parallel Cluster) instead of FSx to store Neuron Cache. This will be fixed in an upcoming release.

### Running in-place update operations (e.g. `all_reduce`) on 0-dimensional tensors result in buffer aliasing errors in torch 2.5 and earlier

Torch's lazy tensor core has a feature where 0-dimensional tensors are stored in a device cache, so scalar constant values can be transferred once and then reused. The values in the device cache are supposed to be marked read-only and never participate in parameter aliasing. However, due to a bug in torch-xla 2.5 (#8499), sometimes the read-only flag can be dropped, allowing these tensors to be donated, resulting in aliasing errors later when the cached value is used again.

A work-around is to avoid using 0-dimensional tensors by changing them to be 1d tensor of length 1 (example). If modifying library code is not possible, disable XLA parameter aliasing by setting environment variable `XLA_ENABLE_PARAM_ALIASING=0`



### Tensor split on second dimension of 2D array not working

Currently, when using tensor split operation on a 2D array in the second dimension, the resulting tensors don't have the expected data (<https://github.com/pytorch/xla/issues/8640>). The work-around is to set `XLA_DISABLE_FUNCTIONALIZATION=0`. Another work-around is to use `torch.tensor_split`.

### Import torch\_xla crashed with TypeError: must be called with a dataclass type or instance with torch-xla 2.5 and torch 2.5.1+cpu (CPU flavor)

When using torch 2.5.1+cpu (CPU flavor) on python 3.10, importing torch\_xla crashed with `TypeError: must be called with a dataclass type or instance` due to installed triton version 3.2.0 (<https://github.com/pytorch/xla/issues/8560>). To work-around, please remove the installed triton package or downgrade to triton==3.1.0 or use the regular torch 2.5.1 (GPU flavor).

### Certain sequence of operations with xm.save() could corrupt tensors

When using the `xm.save` function to save tensors, please use `xm.mark_step()` before `xm.save` to avoid the error described in <https://github.com/pytorch/xla/issues/8422> where parameter aliasing could corrupt other tensor values. This issue will be fixed in a future release.

(Here `xm` is `torch_xla.core.xla_model` following PyTorch/XLA convention)

### Lower BERT pretraining performance when switch to using model.to(torch.bfloat16)

Currently, BERT pretraining performance is ~11% lower when switching to using `model.to(torch.bfloat16)` as part of migration away from the deprecated environment variable `XLA_DOWNCAST_BF16` due to <https://github.com/pytorch/xla/issues/8545>. As a work-around to recover the performance, you can set `XLA_DOWNCAST_BF16=1` which would still work in torch-neuronx 2.5 and 2.6 although there will be deprecation warnings (as noted below).

### Warning “XLA\_DOWNCAST\_BF16 will be deprecated after the 2.5 release, please downcast your model directly”

Environment variables `XLA_DOWNCAST_BF16` and `XLA_USE_BF16` are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

### WARNING:root:torch\_xla.core.xla\_model.xrt\_world\_size() will be removed in release 2.7. is deprecated. Use torch\_xla.runtime.world\_size instead.

This is a warning that `torch_xla.core.xla_model.xrt_world_size()` will be removed in a future release. Please switch to using `torch_xla.runtime.world_size` instead.

**WARNING:** `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in release 2.7. is deprecated. Use `torch_xla.runtime.global_ordinal` instead.

This is a warning that `torch_xla.core.xla_model.xla_model.get_ordinal()` will be removed in a future release. Please switch to using `torch_xla.runtime.global_ordinal` instead.

#### **AttributeError: module 'torch\_xla.runtime' has no attribute 'using\_pjrt'**

In Torch-XLA 2.5, `torch_xla.runtime.using_pjrt` is removed because PJRT is the sole Torch-XLA runtime. See [commit PR](#).

#### **Socket Error: Socket failed to bind**

In PyTorch 2.5, there needs to be a socket available for both `torchrun` and the `init_process_group` to bind. Both of these, by default, will be set to unused sockets. If you plan to use a `MASTER_PORT` environment variable then this error may occur, if the port you set it to is already in use.

```
[W socket.cpp:426] [c10d] The server socket has failed to bind to [::]:29500 (errno: 98 -  
↪ Address already in use).  
[W socket.cpp:426] [c10d] The server socket has failed to bind to ?UNKNOWN? (errno: 98 -  
↪ Address already in use).  
[E socket.cpp:462] [c10d] The server socket has failed to listen on any local network  
↪ address.  
RuntimeError: The server socket has failed to listen on any local network address.  
The server socket has failed to bind to ?UNKNOWN? (errno: 98 - Address already in use).
```

To resolve the issue, please ensure if you are setting `MASTER_PORT` that the port you're setting it to is not used anywhere else in your scripts. Otherwise, you can leave `MASTER_PORT` unset, and `torchrun` will set the default port for you.

#### **AttributeError: module 'torch' has no attribute 'xla' Failure**

In PyTorch 2.5, training scripts might fail during activation checkpointing with the error shown below.

```
AttributeError: module 'torch' has no attribute 'xla'
```

The solution is to use `torch_xla.utils.checkpoint.checkpoint` instead of `torch.utils.checkpoint.checkpoint` as the checkpoint function while wrapping pytorch modules for activation checkpointing. Refer to the [pytorch/xla](#) discussion regarding this [issue](#). Also set `use_reentrant=True` while calling the `torch_xla` checkpoint function. Failure to do so will lead to XLA currently does not support `use_reentrant==False` error. For more details on checkpointing, refer the [documentation](#).

## Error Attempted to access the data pointer on an invalid python storage when using HF Trainer API

While using HuggingFace Transformers Trainer API to train (i.e. *HuggingFace Trainer API fine-tuning tutorial*), you may see the error “Attempted to access the data pointer on an invalid python storage”. This is a known [issue](#) and has been fixed in the version 4.37.3 of HuggingFace Transformers.

## ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory on Amazon Linux 2023

torch-xla version 2.5+ now requires libcrypt.so.1 shared library. Currently, Amazon Linux 2023 includes libcrypt.so.2 shared library by default so you may see *ImportError: libcrypt.so.1: cannot open shared object file: No such file or directory* when using torch-neuronx 2.1+ on Amazon Linux 2023. To install libcrypt.so.1 on Amazon Linux 2023, please run the following installation command (see also <https://github.com/amazonlinux/amazon-linux-2023/issues/182> for more context):

```
sudo yum install libxcrypt-compat
```

## FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path' Failure

In PyTorch 2.5, users might face the error shown below due to incompatible libneuronxla and torch-neuronx versions being installed.

```
FileNotFoundError: [Errno 2] No such file or directory: 'libneuronpjrt-path'
```

Check that the version of libneuronxla that support PyTorch NeuronX 2.5 is 2.1.\*. If not, then uninstall libneuronxla using `pip uninstall libneuronxla` and then reinstall the packages following the [installation guide](#)

## Glibc error on Amazon Linux 2

If using Torch-NeuronX 2.5 on Amazon Linux 2, you will see a Glibc error below. Please switch to a newer supported OS such as Ubuntu 22 or Amazon Linux 2023.

```
ImportError: /lib64/libc.so.6: version `GLIBC_2.27' not found (required by /tmp/debug/_
↳ XLAC.cpython-38-x86_64-linux-gnu.so)
```

## Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range error during Neuron Parallel Compile

When running Neuron Parallel Compile with HF Trainer API, you may see the error `Status: INVALID_ARGUMENT: Input dimension should be either 1 or equal to the output dimension it is broadcasting into or IndexError: index out of range` in Accelerator’s `pad_across_processes` function. This is due to data-dependent operation in evaluation metrics computation. Data-dependent operations would result in undefined behavior with Neuron Parallel Compile trial execution (execute empty graphs with zero outputs). To work-around this error, please disable `compute_metrics` when `NEURON_EXTRACT_GRAPHS_ONLY` is set to 1:

```
compute_metrics=None if os.environ.get("NEURON_EXTRACT_GRAPHS_ONLY") else compute_metrics
```

## Compiler assertion error when running Stable Diffusion training

Currently, with PyTorch 2.5 (torch-neuronx), we are seeing the following compiler assertion error with Stable Diffusion training when gradient accumulation is enabled. This will be fixed in an upcoming release. For now, if you would like to run Stable Diffusion training with Neuron SDK release 2.21/2.22, please disable gradient accumulation in torch-neuronx 2.5.

```
ERROR 222163 [NeuronAssert]: Assertion failure in usr/lib/python3.9/concurrent/futures/
↳ process.py at line 239 with exception:
too many partition dims! {{0,+,960}}[10],+,10560}}[10]
```

## Frequently Asked Questions (FAQ)

### Do I need to recompile my models with PyTorch 2.5?

Yes.

### Do I need to update my scripts for PyTorch 2.5?

Please see the *migration guide*

### What environment variables will be changed with PyTorch NeuronX 2.5 ?

The environment variables XLA\_DOWNCAST\_BF16 and XLA\_USE\_BF16 are deprecated (warning when used). Please switch to automatic mixed-precision or use `model.to(torch.bfloat16)` command to cast model to BF16. (see `migration_from_xla_downcast_bf16`)

### What features will be missing with PyTorch NeuronX 2.5?

PyTorch NeuronX 2.5 now has most of the supported features in PyTorch NeuronX 2.1, with known issues listed above, and unsupported features as listed in *PyTorch Neuron (torch-neuronx) release notes*.

### Can I use Neuron Distributed and Transformers Neuron libraries with PyTorch NeuronX 2.5?

Yes, NeuronX Distributed, and Transformers NeuronX, and AWS Neuron Reference for NeMo Megatron libraries will work with PyTorch NeuronX 2.5.

### Can I still use PyTorch 2.1 version?

PyTorch 2.1 is supported for release 2.21 and will reach end-of-life in a future release. Additionally, the CVEs [CVE-2024-31583](#) and [CVE-2024-31580](#) affect PyTorch versions 2.1 and earlier. We recommend upgrading to the new version of Torch-NeuronX by following *PyTorch Neuron (torch-neuronx) Setup*.

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 8.4 Neuron FAQ

### Table of contents

- *Neuron 2.x FAQ*
- *Training Only FAQ*
- *Inference Only FAQ*
- *Runtime FAQ*
- *Compiler FAQ*
- *Neuron Containers*
- *ONNX FAQ*
- *Support*

### 8.4.1 Neuron 2.x FAQ

- [neuron2-intro-faq](#)

### 8.4.2 Training Only FAQ

- [neuron-training-faq](#)

### 8.4.3 Inference Only FAQ

- [neuron-f1-faq](#)
- [trouble-shooting-inf1-faq](#)
- [tf1\\_faq](#)
- [tf2\\_faq](#)
- [NeuronPerf](#)

### 8.4.4 Runtime FAQ

- *Neuron Runtime FAQ*

### 8.4.5 Compiler FAQ

- *Neuron Compiler FAQ (neuronx-cc)*
- *Neuron Compiler FAQ (neuron-cc)*

### 8.4.6 Neuron Containers

- *Neuron Containers FAQ*

### 8.4.7 ONNX FAQ

- onnx-faq

### 8.4.8 Support

- neuron\_roadmap\_faq
- contribute-faq

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## 8.5 Troubleshooting Guide

### Table of contents

- *Training Only Troubleshooting*
- *Inference Only Troubleshooting*
- *Runtime Troubleshooting*
- *Containers Troubleshooting*
- *Setup Troubleshooting*

### 8.5.1 Training Only Troubleshooting

- *PyTorch Neuron for Training*

## 8.5.2 Inference Only Troubleshooting

- *PyTorch Neuron for Inference*
- *NeuronPerf*
- *MXNet Neuron*

## 8.5.3 Runtime Troubleshooting

- *Neuron Runtime Troubleshooting on Inf1 and Trn1*

## 8.5.4 Containers Troubleshooting

- *Containers*

## 8.5.5 Setup Troubleshooting

- *neuron-setup-troubleshooting*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 8.6 Neuron Glossary

### Table of contents

- *Terms*
  - *Neuron Devices (Accelerated Machine Learning chips)*
  - *Neuron powered Instances*
  - *NeuronCore terms*
  - *Neuron SDK terms*
- *Abbreviations*

### 8.6.1 Terms

**Neuron Devices (Accelerated Machine Learning chips)**

| Term                        | Description   |
|-----------------------------|---|
| <b>Inferentia</b>           | AWS first generation accelerated machine learning chip supporting inference only          |
| <b>Trainium/Inferentia2</b> | AWS second generation accelerated machine learning chip supporting training and inference |
| <b>Trainium2</b>            | AWS second generation accelerated machine learning chip supporting training and inference |
| <b>Neuron Device</b>        | Accelerated machine learning chip (e.g. Inferentia or Trainium)                           |

**Neuron powered Instances**

| Term        | Description  |
|-------------|--|
| <b>Inf1</b> | Inferentia powered accelerated compute EC2 instance  |
| <b>Trn1</b> | Trainium powered accelerated compute EC2 instance    |
| <b>Inf2</b> | Inferentia2 powered accelerated compute EC2 instance |
| <b>Trn2</b> | Trainium2 powered accelerated compute EC2 instance   |



## NeuronCore terms

| Term                                   | Description  |
|--|--|
| <b>NeuronCore</b>                      | The machine learning compute cores within Inferentia/Trainium  |
| <b>NeuronCore-v1</b>                   | Neuron Core within Inferentia  |
| <b>NeuronCore-v2</b>                   | Neuron Core within Trainium1/Inferentia2   |
| <b>NeuronCore-v3</b>                   | Neuron Core within Trainium2   |
| <b>Tensor Engine</b>                   | 2D systolic array (within the NeuronCore), used for matrix computations  |
| <b>Scalar Engine</b>                   | A scalar-engine within each NeuronCore, which can accelerate element-wise operations (e.g. GELU, ReLU, reciprocal, etc)              |
| <b>Vector Engine</b>                   | A vector-engine with each NeuronCore, which can accelerate spatial operations (e.g. layerNorm, TopK, pooling, etc)                   |
| <b>GPSIMD Engine</b>                   | Embedded General Purpose SIMD cores, within each NeuronCore, to accelerate custom-operators  |
| <b>Sync Engine</b>                     | The SP engine, which is integrated inside NeuronCore. Used for synchronization and DMA triggering.                                   |
| <b>Collective Communication Engine</b> | Dedicated engine for collective communication, allows for overlapping computation and communication                                  |
| <b>High Bandwidth Memory</b>           | <a href="#">High Bandwidth Memory</a> , used as device memory for NeuronCore-v2 and beyond.  |
| <b>State Buffer</b>                    | The main software-managed on-chip memory in NeuronCore-v1 and beyond.  |
| <b>Partial Sum Buffer</b>              | A second software-managed on-chip memory in NeuronCore-v1 and beyond, with near-memory accumulation support for TensorE output data. |
| <b>NeuronLink</b>                      | Interconnect between NeuronCores   |
| <b>NeuronLink-v1</b>                   | Interconnect between NeuronCores in Inferentia device  |
| <b>NeuronLink-v2</b>                   | Interconnect between NeuronCores in Trainium1/Inferentia2 device   |
| <b>NeuronLink-v3</b>                   | Interconnect between NeuronCores in Trainium2 device   |

## Neuron SDK terms

| Term                           | Description  |
|--------------------------------|--|
| <b>Neuron Kernel Interface</b> | A bare-metal language and compiler for directly programming Neuron devices available on AWS Trainium/Inferentia2 and beyond devices. |

## 8.6.2 Abbreviations

| Abbreviation         | Description                           |
|----------------------|---------------------------------------|
| <b>NxD Core</b>      | NeuronX Distributed Core Library      |
| <b>NxD Training</b>  | NeuronX Distributed Training Library  |
| <b>NxD Inference</b> | NeuronX Distributed Inference Library |
| <b>NC</b>            | Neuron Core                           |
| <b>NeuronCore</b>    | Neuron Core                           |
| <b>ND</b>            | Neuron Device                         |
| <b>NeuronDevice</b>  | Neuron Device                         |
| <b>TensorE</b>       | Tensor Engine                         |
| <b>ScalarE</b>       | Scalar Engine                         |
| <b>VectorE</b>       | Vector Engine                         |
| <b>GpSimdE</b>       | GpSimd Engine                         |
| <b>CCE</b>           | Collective Communication Engine       |
| <b>HBM</b>           | High Bandwidth Memory                 |

continues on next page

Table 8.1 – continued from previous page

| Abbreviation     | Description             |
|------------------|-------------------------|
| <b>SBUF</b>      | State Buffer            |
| <b>PSUM</b>      | Partial Sum Buffer      |
| <b>FP32</b>      | Float32                 |
| <b>TF32</b>      | TensorFloat32           |
| <b>FP16</b>      | Float16                 |
| <b>BF16</b>      | Bfloat16                |
| <b>cFP8</b>      | Configurable Float8     |
| <b>RNE</b>       | Round Nearest Even      |
| <b>SR</b>        | Stochastic Rounding     |
| <b>NKI</b>       | Neuron Kernel Interface |
| <b>CustomOps</b> | Custom Operators        |
| <b>RT</b>        | Neuron Runtime          |
| <b>DP</b>        | Data Parallel           |
| <b>DPr</b>       | Data Parallel degree    |
| <b>TP</b>        | Tensor Parallel         |
| <b>TPr</b>       | Tensor Parallel degree  |

continues on next page

Table 8.1 – continued from previous page

| Abbreviation | Description              |
|--------------|--------------------------|
| <b>PP</b>    | Pipeline Parallel        |
| <b>PPr</b>   | Pipeline Parallel degree |

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

---

## LEGACY SOFTWARE

*This document is relevant for: Inf1*

### 9.1 MXNet Neuron (maintenance)

MXNet Neuron unlocks high-performance and cost-effective deep learning acceleration on AWS Trainium-based and Inferentia-based Amazon EC2 instances.

MXNet Neuron enables native MXNet models to be accelerated on Neuron devices, so you can use your existing framework application and get started easily with minimal code changes.

*This document is relevant for: Inf1*

#### 9.1.1 MXNet Neuron Setup

MxNet Neuron (`mxnet-neuron`) Setup for Inf1 Instances *This document is relevant for: Inf1*

*This document is relevant for: Inf1*

#### 9.1.2 Inference (`mxnet-neuron`) (maintenance)

*This document is relevant for: Inf1*

#### Tutorials (`mxnet-neuron`)

*This document is relevant for: Inf1*

#### Computer Vision Tutorials (`mxnet-neuron`)

- ResNet-50 tutorial [\[html\]](#) [\[notebook\]](#)
- Model Serving tutorial [\[html\]](#)
- Getting started with Gluon tutorial [\[html\]](#) [\[notebook\]](#)

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

### Natural Language Processing (NLP) Tutorials (mxnet-neuron)

- MXNet 1.8: Using data parallel mode tutorial [\[html\]](#) [\[notebook\]](#)

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

### Utilizing Neuron Capabilities Tutorials (mxnet-neuron)

- NeuronCore Groups tutorial [\[html\]](#) [\[notebook\]](#)

*This document is relevant for:* Inf1

### Computer Vision Tutorials

- ResNet-50 tutorial [\[html\]](#) [\[notebook\]](#)
- Model Serving tutorial [\[html\]](#)
- Getting started with Gluon tutorial [\[html\]](#) [\[notebook\]](#)

### Natural Language Processing (NLP) Tutorials

- MXNet 1.8: Using data parallel mode tutorial [\[html\]](#) [\[notebook\]](#)

### Utilizing Neuron Capabilities Tutorials

- NeuronCore Groups tutorial [\[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

### API Reference Guide (mxnet-neuron)

*This document is relevant for:* Inf1

## Neuron Apache MXNet Compilation Python API

The MXNet-Neuron compilation Python API provides a method to compile model graph for execution on Inferentia.

### Description

Within the graph or subgraph, the compile method selects and sends Neuron-supported operations to Neuron-Compiler for compilation and saves the compiled artifacts in the graph. Uncompilable operations are kept as original operations for framework execution.

The compiled graph can be saved using the MXNet `save_checkpoint` and served using MXNet Model Serving. Please see `mxnet-neuron-model-serving` for more information about exporting to saved model and serving using MXNet Model Serving.

Options can be passed to Neuron compiler via the compile function. For example, the “`--neuroncore-pipeline-cores`” option directs Neuron compiler to compile each subgraph to fit in the specified number of NeuronCores. This number can be less than the total available NeuronCores on an Inf1 instance. See *Neuron compiler CLI Reference Guide (neuron-cc)* for more information about compiler options.

For debugging compilation, use `SUBGRAPH_INFO=1` environment setting before calling the compilation script. The extract subgraphs are preserved as hidden files in the run directory. For more information, see *Using Neuron GatherInfo Tool to collect debug and support information*

### MXNet 1.5

#### Method

```
from mxnet.contrib import neuron
neuron.compile(sym, args, aux, inputs, **compile_args)
```

#### Arguments

- **sym** - Symbol object loaded from symbol.json file
- **args** - args/params dictionary loaded from params file
- **aux** - aux/params dictionary loaded from params file
- **inputs** - a dictionary with key/value mappings for input name to input numpy arrays
- **kwargs** (optional) - a dictionary with key/value mappings for MXNet-Neuron compilation and Neuron Compiler options.
  - For example, to limit the number of NeuronCores per subgraph, use `compile_args={'--neuroncore-pipeline-cores' : N}` where N is an integer representing the maximum number of NeuronCores per subgraph.
  - Additional compiler flags can be passed using `'flags' : [<flags>]` where is a comma separated list of strings. See *Using Neuron GatherInfo Tool to collect debug and support information* for example of passing debug flags to compiler.
  - Advanced option to exclude node names: `compile_args={'excl_node_names' : [<node names>]}` where is a comma separated list of node name strings.

## Returns

- **sym** - new partitioned symbol
- **args** - modified args/params
- **auxs** - modified aux/params

## Example Usage: Compilation

The following is an example usage of the compilation, with default compilation arguments:

```
from mxnet.contrib import neuron
...
neuron.compile(sym, args, aux, inputs={'data' : img})
```

## MXNet 1.8

### Method

```
import mx_neuron as neuron
neuron.compile(obj, args=None, aux=None, inputs=None, **compile_args)
```

## Arguments

- **obj** - Symbol object loaded from symbol.json file or gluon.HybridBlock object
- **args** (optional) - args/params dictionary loaded from params file. Only needed in case of Symbol object
- **aux** (optional) - aux/params dictionary loaded from params file. Only needed in case of Symbol object
- **inputs** - a dictionary with key/value mappings for input name to input numpy arrays.
- **kwargs** (optional) - a dictionary with key/value mappings for MXNet-Neuron compilation and Neuron Compiler options.
  - For example, to limit the number of NeuronCores per subgraph, use `compile_args={'--neuroncore-pipeline-cores' : N}` where N is an integer representing the maximum number of NeuronCores per subgraph.
  - Additional compiler flags can be passed using `'flags' : [<flags>]` where is a comma separated list of strings. See *Using Neuron GatherInfo Tool to collect debug and support information* for example of passing debug flags to compiler.
  - Advanced option to exclude node names: `compile_args={'excl_node_names' : [<node names>]}` where is a comma separated list of node name strings.
  - `work_dir`: relative or absolute path for storing compiler artifacts (including params and jsons) generated during compilation when `SUBGRAPH_INFO=1`.



## Returns

- **(sym, args, auxs)** - for symbol object as input. sym, args and auxs are new partitioned symbol, modified args/params and modified aux/params respectively.
- **(obj)** - for gluon.HybridBlock object as input. obj is the partitioned and optimized gluon.Hybrid block object for Neuron backend.

## Example Usage: Compilation

The following is an example usage of the compilation, with default compilation arguments for symbol object:

```
import mx_neuron as neuron
...
neuron.compile(sym, args, aux, inputs={'data' : img})
```

The following is an example usage of the compilation, with default compilation arguments for gluon.HybridBlock object (only supported in MXNet-Neuron 1.8):

```
import mx_neuron as neuron
...
neuron.compile(obj, inputs={'data' : img})
```

## Example Usage: Extract Compilation Statistics

To extract operation counts, insert the following code after compile step (assume csym is the compiled MXNet symbol):

```
import json

# Return list of nodes from MXNet symbol
def sym_nodes(sym):
    return json.loads(sym.tojson())['nodes']

# Return number of operations in node list
def count_ops(graph_nodes):
    return len([x['op'] for x in graph_nodes if x['op'] != 'null'])

# Return triplet of compile statistics
# - count of operations in symbol database
# - number of Neuron subgraphs
# - number of operations compiled to Neuron runtime
def get_compile_stats(sym):
    cnt = count_ops(sym_nodes(sym))
    neuron_subgraph_cnt = 0
    neuron_compiled_cnt = 0
    for g in sym_nodes(sym):
        if g['op'] == '_neuron_subgraph_op':
            neuron_subgraph_cnt += 1
            for sg in g['subgraphs']:
                neuron_compiled_cnt += count_ops(sg['nodes'])
    return (cnt, neuron_subgraph_cnt, neuron_compiled_cnt)
```

(continues on next page)

(continued from previous page)

```
original_cnt = count_ops(sym_nodes(sym))
post_compile_cnt, neuron_subgraph_cnt, neuron_compiled_cnt = get_compile_stats(csym)
print("INFO:mxnet: Number of operations in original model: ", original_cnt)
print("INFO:mxnet: Number of operations in compiled model: ", post_compile_cnt)
print("INFO:mxnet: Number of Neuron subgraphs in compiled model: ", neuron_subgraph_cnt)
print("INFO:mxnet: Number of operations placed on Neuron runtime: ", neuron_compiled_cnt)
```

```
INFO:mxnet: Number of operations in original model: 67
INFO:mxnet: Number of operations in compiled model: 4
INFO:mxnet: Number of Neuron subgraphs in compiled model: 2
INFO:mxnet: Number of operations placed on Neuron runtime: 65
```

*This document is relevant for:* Inf1

- *Neuron Apache MXNet Compilation Python API*

*This document is relevant for:* Inf1

*This document is relevant for:* Inf1

## Developer Guide

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

## Flexible Execution Group (FlexEG) in Neuron-MXNet

### Introduction

Inf1 instances are available with a different number of Inferentia chips, each Inferentia chip is combined of 4 Neuron-Cores and an Inf1 instance includes 4 to 64 NeuronCores depending on the instance size. With Neuron Runtime 1.x (neuron-rtd server), NeuronCores could be combined into NeuronCore Groups (NCG), which were basic scheduling units of compiled neural network in Neuron. Creation of desired sized NCGs was done at the start of the application and could not be modified afterwards.

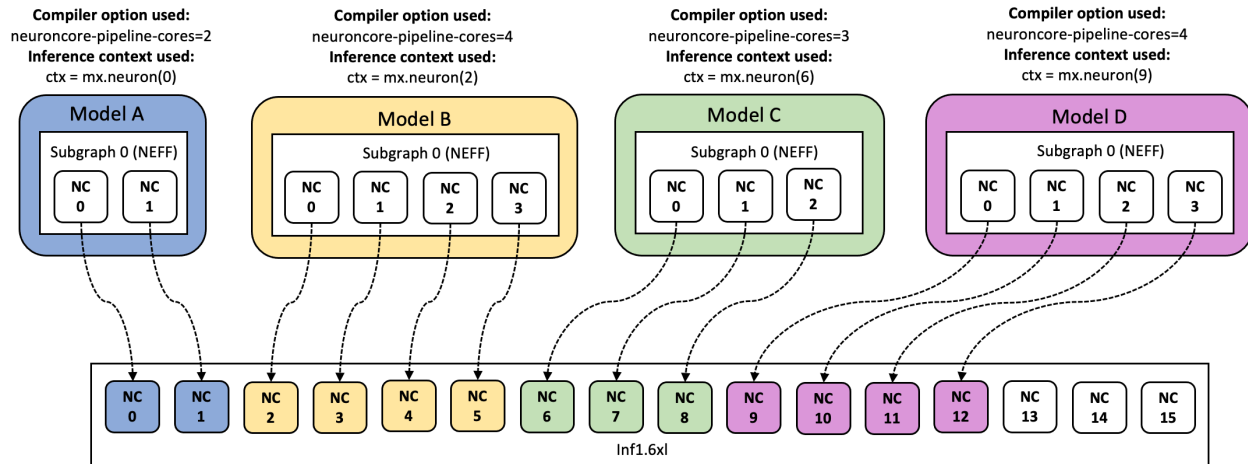
Starting with Neuron SDK 1.16.0, and with the introduction of Neuron Runtime 2.x, MXNet Neuron 1.8 introduces Flexible Execution Groups (FlexEG) feature. With FlexEG, you do not have to create NCGs at the start of the process, instead you will set the index of the first NeuronCore you want to load models onto, and FlexEG feature will enable the flexibility of loading models onto any available NeuronCore on the inf1 instance starting from the first NeuronCore you set. This guide will show you how to efficiently utilize NeuronCores using FlexEG feature in NeuronMXNet.

### FlexEG

With the introduction of FlexEG, you don't need to create NCGs and can load models onto a group of consecutive NeuronCores by providing the index of the first NeuronCore in the group. Neuron runtime takes care of figuring out the number of NeuronCores required for the given compiled model and loads the model using the required number of cores (sequentially starting with the NeuronCore index provided by the user).

For example, assuming that you have an Inf1.6x1 machine and there are 4 models A, B, C, D compiled to 2, 4, 3, and 4 NeuronCores respectively, you can map any model to any core by context `mx.neuron(neuron_core_index)` where `neuron_core_index` is the NeuronCore index (0,1,2,3,4 ... ).

In the example below, you map model A to `mx.neuron(0)` context, model B to `mx.neuron(2)` context, model C to `mx.neuron(6)` context and model D to `mx.neuron(9)` context.

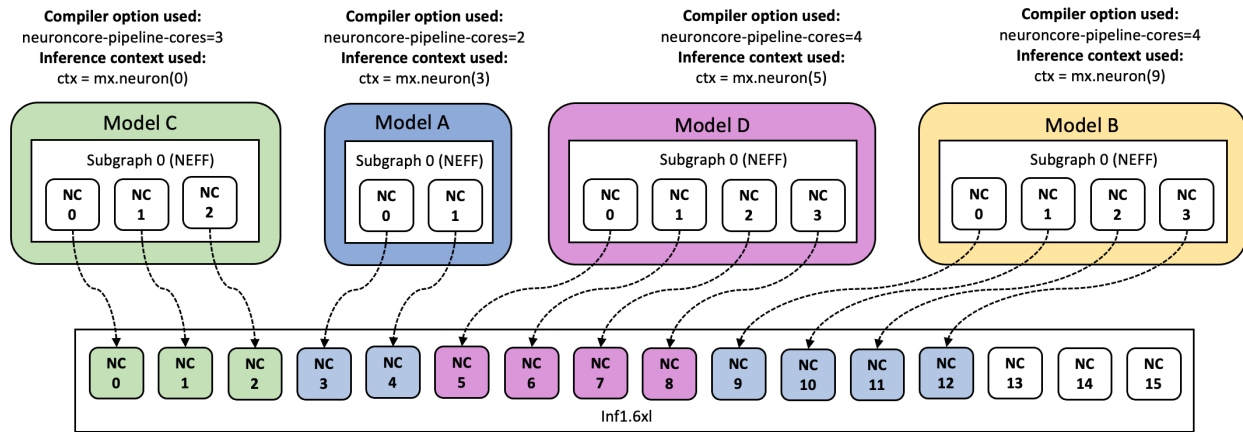


The above configuration is achieved by using application code similar to below:

```
# Load models (MXNet)
# loaded into the 2 cores starting with core 0
sym, args, aux = mx.model.load_checkpoint(mx_model0_file, 0)
model0 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')
# loaded into the 4 cores starting with core 2
sym, args, aux = mx.model.load_checkpoint(mx_model1_file, 0)
model1 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')
# loaded into the 3 cores starting with core 6
sym, args, aux = mx.model.load_checkpoint(mx_model2_file, 0)
model2 = sym.bind(ctx=mx.neuron(6), args=args, aux_states=aux, grad_req='null')
# loaded into the 4 cores starting with core 9
sym, args, aux = mx.model.load_checkpoint(mx_model3_file, 0)
model3 = sym.bind(ctx=mx.neuron(9), args=args, aux_states=aux, grad_req='null')

# run inference by simply calling the loaded model
results0 = model0.forward(data=inputs0)
results1 = model1.forward(data=inputs1)
results2 = model2.forward(data=inputs2)
results3 = model3.forward(data=inputs3)
```

Since there is no NCG creation at the start of the process, you can load the same four models but in a different configuration by changing the context being used for inference. For example, you could map model C to `mx.neuron(0)` context, model A to `mx.neuron(3)` context, model D to `mx.neuron(5)` context and model B to `mx.neuron(9)` context.



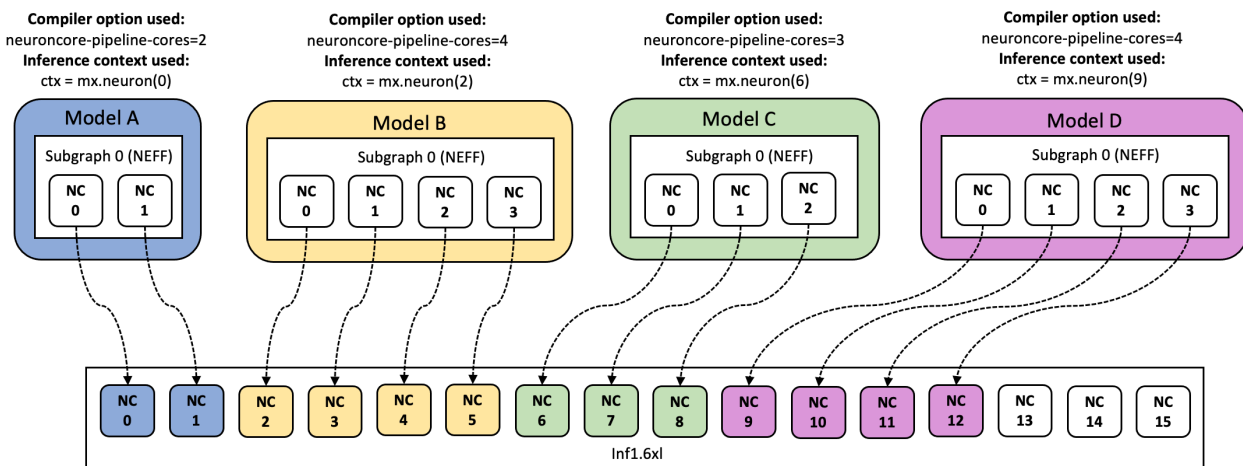
## Migration from NeuronCore Groups to FlexEG

NeuronCore Groups are defined by setting the environment variable `NEURONCORE_GROUP_SIZES` with a comma separated list of number of cores in each group. In this mode of operation, number of devices (defined in `NEURONCORE_GROUP_SIZES`) are grouped together to create a single entity.

`NEURONCORE_GROUP_SIZES` environment variable is set at runtime:

```
#!/bin/bash
export NEURONCORE_GROUP_SIZES=2,4,3,4
python your_neuron_application.py
```

NeuronCore groups are created once at the start of the application and cannot be modified / re-created till the application process runs. The above flow creates 4 neuron devices with 2,4,3 and 4 devices each. In order to get the same configuration as the example from before, you map model A to `mx.neuron(0)` context, model B to `mx.neuron(1)` context, model C to `mx.neuron(2)` context and model D to `mx.neuron(3)` context.



This can be achieved programmatically as shown below:

```
# Set Environment
os.environ['NEURONCORE_GROUP_SIZES']='2,4,3,4'

# Load models (MXNet)
```

(continues on next page)

(continued from previous page)

```

# loaded into the first group of NC0-NC1
sym, args, aux = mx.model.load_checkpoint(mx_model0_file, 0)
model0 = sym.bind(ctx=mx.neuron(0), args=args, aux_states=aux, grad_req='null')
# loaded into the second group of NC2-NC5
sym, args, aux = mx.model.load_checkpoint(mx_model1_file, 0)
model1 = sym.bind(ctx=mx.neuron(1), args=args, aux_states=aux, grad_req='null')
# loaded into the third group of NC6-NC8
sym, args, aux = mx.model.load_checkpoint(mx_model2_file, 0)
model2 = sym.bind(ctx=mx.neuron(2), args=args, aux_states=aux, grad_req='null')
# loaded into the fourth group of NC9-NC12
sym, args, aux = mx.model.load_checkpoint(mx_model3_file, 0)
model3 = sym.bind(ctx=mx.neuron(3), args=args, aux_states=aux, grad_req='null')

# run inference by simply calling the loaded model
results0 = model0.forward(data=inputs0)
results1 = model1.forward(data=inputs1)
results2 = model2.forward(data=inputs2)
results3 = model3.forward(data=inputs3)

```

So comparing to FlexEG, we see that in case of NCGs neuron context requires the index of the execution group, while in FlexEG neuron context requires the NeuronCore index of the first NeuronCore on which the model is supposed to be loaded and executed. For example, with `NEURONCORE_GROUP_SIZES='2,4,3,4'`, `ctx=mx.neuron(1)` loads the model on execution group 1 which effectively loads the model on the 2nd NCG group which has 4 NeuronCores.

## Best practices when using FlexEG

FlexEG gives the user most flexibility in terms of accessing cores and loading models on specific cores. With this the users can effortlessly load and execute new models on NeuronCores without closing the application. Here we shall outline some of the best practices that should be kept in mind while using FlexEG.

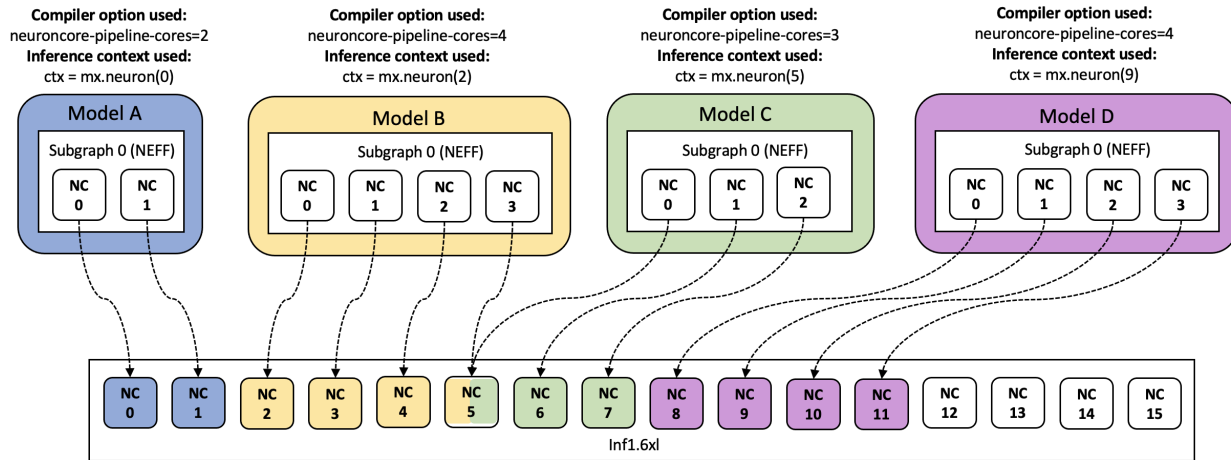
### Choosing starting core

FlexEG tries to use the required number of cores (based on the input model) starting with the core index provided by the user. In case the system, doesn't have the required number of cores after the starting core index, model load will fail. For example: We have a model X which needs 2 cores and an `inf1.xl` machine with 4 NeuronCores (NeuronCore indexes are: 0, 1, 2 and 3). As the model needs at least 2 cores, valid start indexes for this model are: 0, 1, 2. However if the user gives 3 as the neuron context, then there are no 2 cores available starting from core 3. So it will fail.

### Performance vs. Flexibility tradeoff

While using data parallel model of operation (where models are executed in parallel), for optimal performance the user should make sure that the models are not sharing any cores. That is because NeuronCores can execute one model at a time, when two or more models are executed on the same core (assuming that they are already loaded), it executes the first model, stops it, starts the second model and then executes it. This is called model switching and involves additional overhead and prevents execution on model in parallel. For example: assuming that you have an `Inf1.6x1` machine and there are 4 models A, B, C, D compiled to 2, 4, 3, and 4 NeuronCores respectively. Loading model A to `mx.neuron(0)` context, model B to `mx.neuron(2)` context, model C to `mx.neuron(6)` context and model D to `mx.neuron(9)` context is a good configuration because no two models are sharing NeuronCores and thus can be executed in parallel. However, Loading model A to `mx.neuron(0)` context, model B to `mx.neuron(2)` context, model C to

`mx.neuron(5)` context and model D to `mx.neuron(9)` context is not a good configuration as models B and C share NeuronCore 5 and thus cannot be executed in parallel.



This document is relevant for: Inf1, Inf2, Trn1, Trn2

- *Flexible Execution Group (FlexEG) in Neuron-MXNet*

This document is relevant for: Inf1

This document is relevant for: Inf1

## Misc (mxnet-neuron)

This document is relevant for: Inf1

## Troubleshooting Guide for Neuron Apache MXNet

### Table of Contents

- *Inference Runtime Error*
  - *Out-of-memory error when calling Symbol API bind() too many times*
  - *Inference crashed with MXNetError: InferShapeKeyword argument name xyz not found*
  - *Inference crashed at mx.nd.waitall() with MXNetError: Check failed: bin.dtype() == mshadow::kUInt8*
  - *Inference crashed with NRTD error 1002*
- *Multi-Model Server*
  - *Failed to create NEURONCORE Group with GRPC Error. Status Error: 14, Error message: "Connect Failed"*
  - *Multiple MMS workers die with "Backend worker process die." message*
  - *MMS throws a "mxnet.base.MXNetError: array::at" error*
  - *MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded*

- Pipeline mode is not able to execute inferences requests in parallel
- Features only in MXNet-Neuron 1.5
- Features only in MXNet-Neuron 1.8

## Inference Runtime Error

### Out-of-memory error when calling Symbol API bind() too many times

**Important:** NEURONCORE\_GROUP\_SIZES will no longer be supported starting Neuron 1.19.0 release if your application is using NEURONCORE\_GROUP\_SIZES please see [Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#) and `eol-ncgs-env_2` for more details.

If you see out-of-memory error when using Symbol API's `bind()` function, please ensure that the `bind()` function is called once for each desired model instance. For example, on `inf1.xlarge`, use Symbol API to create 4 parallel instances of a model that was compiled to 1 NeuronCore (`-neuroncore-pipeline-cores=1`), each is bound to an different `mx.neuron(i)` context where `i` is the NeuronCore Group index ranging from 0 to 3. Then use 4 threads to feed the 4 instances in parallel. For example:

```
NUM_PARALLEL = 4
os.environ['NEURONCORE_GROUP_SIZES'] = ','.join('1' for _ in range(NUM_PARALLEL))

data_iter = []
for i in range(NUM_PARALLEL):
    data_iter.append(mx.io.ImageRecordIter(
        path_imgrec=recfile_base, data_shape=(3, 224, 224), batch_size=1,
        prefetch_buffer=1,
        num_parts=NUM_PARALLEL, part_index=i))

sym, args, auxs = mx.model.load_checkpoint('resnet-50_compiled', 0)

exec_list = []
for i in range(NUM_PARALLEL):
    exec = sym.bind(ctx=mx.neuron(i), args=args, aux_states=auxs, grad_req='null')
    exec_list.append(exec)

def single_thread_infer(i):
    for batch in data_iter[i]:
        img = batch.data[0]
        label = batch.label
        feed_dict = {'data': img}
        exe = exec_list[i]
        exe.copy_params_from(feed_dict)
        exe.forward()
        out = exe.outputs[0]

future_list = []
with futures.ThreadPoolExecutor(max_workers=NUM_PARALLEL) as executor:
    for i in range(NUM_PARALLEL):
        future_list.append(executor.submit(single_thread_infer, i))
```

### Inference crashed with MXNetError: InferShapeKeyword argument name xyz not found

If you see MXNetError:

```
mxnet.base.MXNetError: [11:55:39] src/c_api/c_api_symbolic.cc:508: InferShapeKeyword_
↳ argument name xyz not found."
```

This is followed by a list of “Candidate arguments”. This list shows all the input argument names that the model knows about, and ‘xyz’ is not in the list. To fix this, remove entry xyz from the feed dictionary.

### Inference crashed at mx.nd.waitall() with MXNetError: Check failed: bin.dtype() == mshadow::kUInt8

When executing Symbol API’s forward function followed by mx.nd.waitall(), where MXNetError exception occurs with ‘Check failed: bin.dtype() == mshadow::kUInt8’.

### Inference crashed with NRTD error 1002

During inference, the user may encounter an error with details “[NRTD:infer\_wait] error: 1002”:

```
mxnet.base.MXNetError: [11:26:56] src/operator/subgraph/neuron/./neuron_util.h:1175:
↳ Check failed: rsp_wait.status().code() == 0 || rsp_wait.status().code() == 1003: Failed
Infer Wait with Neuron-RTD Error. Neuron-RTD Status Code: 1002, details: "[NRTD:infer_
↳ wait] error: 1002
"
```

Runtime errors are listed in rtd-return-codes. In particular, 1002 means that some invalid input has been submitted to infer, e.g. missing some of the input tensors, incorrect input tensor sizes. Please examine /var/log/syslog to see more details on the error. For example, you may see:

```
Oct 30 19:13:39 ip-172-31-93-131 nrted[1125]: [TDRV:io_queue_prepare_input_nonhugetlb]
↳ Unexpected input size, for data00, expected: 2097152, received: 33554432
```

This means that the input tensor size is larger than what the model was compiled for (i.e. the example input tensor shapes passed during compilation).

### Multi-Model Server

#### Failed to create NEURONCORE Group with GRPC Error. Status Error: 14, Error message: “Connect Failed”

NOTE: This error only applies to MXNet 1.5.

If the client is unable to start workers and you get a message that MMS is unable to create NeuronCore Group, please check that Neuron RTD is running (neuron-rtd process).

```
{
"code": 500,
"type": "InternalServerError",
"message": "Failed to start workers"
}
```



```
2019-10-23 19:56:23,187 [INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.
↳ mms.wlm.WorkerLifeCycle - [19:56:23] src/operator/subgraph/inferentia/./inferentia_
↳ util.h:218: Check failed: status.ok() Failed to create NeuronCore Group with GRPC_
↳ Error. Status Error: 14, Error message: "Connect Failed"
```

## Multiple MMS workers die with “Backend worker process die.” message

**Important:** NEURONCORE\_GROUP\_SIZES will no longer be supported starting Neuron 1.19.0 release if your application is using NEURONCORE\_GROUP\_SIZES please see [Migrate your application to Neuron Runtime 2.x \(libnrt.so\)](#) and [eol-ncgs-env\\_2](#) for more details.

If you run inference with MMS and get multiple messages “Backend worker process die”, please ensure that the number of workers (“initial\_workers”) passed during load model is less than or equal to number of NeuronCores available divided by number of NeuronCores required by model.

```
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - Backend worker process die.
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - Traceback (most recent call last):
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - File "/usr/local/lib/python3.6/site-packages/
↳ mxnet/symbol/symbol.py", line 1524, in simple_bind
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - ctypes.byref(exe_handle)))
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - File "/usr/local/lib/python3.6/site-packages/
↳ mxnet/base.py", line 252, in check_call
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - raise MXNetError(py_str(_LIB.
↳ MXGetLastError()))
com.amazonaws.ml.mms.wlm.WorkerLifeCycle - mxnet.base.MXNetError: [00:26:32] src/
↳ operator/subgraph/neuron/./neuron_util.h:221: Check failed: 0 == create_eg_rsp.
↳ status().code() Failed to create NeuronCore Group with KRTD Error. KRTD Status Code: 4,
↳ details: ""
```

As indicated in [Performance Tuning](#), for greater flexibility user can use NEURONCORE\_GROUP\_SIZES to specify the groupings of NeuronCores into Neuron devices, each device consisting of one or more NeuronCores. Each worker would take a device. The total number of NeuronCores taken by all the workers should be less than or equal the total number of NeuronCores visible to neuron-rt. This situation should be considered at full load (MMS scales up to max\_workers). Additionally, to properly assign model to Neuron device, the environment NEURONCORE\_GROUP\_SIZES must be specified within the model server class (ie. mxnet\_model\_service.py in the example above). For example, add the following line within mxnet\_model\_service.py for model compiled to 1 NeuronCore:

```
os.environ['NEURONCORE_GROUP_SIZES'] = '1'
```

More information about max\_worker limit setting can be found at [MMS Management API Documentation](#). For example, to run up to 4 workers in inf1.xlarge where 4 NeuronCores are available by default to Neuron-RTD, set max\_workers to 4:

```
curl -v -X PUT "http://localhost:8081/models/squeezenet_v1.1_compiled?min_worker=1?max_
↳ worker=4"
```

## MMS throws a “mxnet.base.MXNetError: array::at” error

If you see “mxnet.base.MXNetError: array::at” when running MMS please check that NDAarray/Gluon API is not used as they are not supported in MXNet-Neuron. If you would like to use NDAarray or Gluon API, please upgrade to MXNet 1.8.

```
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ array::at
[INFO ] W-9000-squeezenet_v1.1_compiled com.amazonaws.ml.mms.wlm.WorkerThread - Backend
↳ response time: 30
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ Traceback (most recent call last):
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ File "/tmp/models/6606fa046f68a34df87f15362a7a2d9a49749878/model_handler.py", line
↳ 82, in handle
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ data = self.inference(data)
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ File "/tmp/models/6606fa046f68a34df87f15362a7a2d9a49749878/mxnet_model_service.py",
↳ line 153, in inference
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ d.wait_to_read()
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ File "/home/user/regression_venv_p3.6/lib/python3.6/site-packages/mxnet/ndarray/
↳ ndarray.py", line 1819, in wait_to_read
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ check_call(_LIB.MXNDArrayWaitToRead(self.handle))
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ File "/home/user/regression_venv_p3.6/lib/python3.6/site-packages/mxnet/base.py",
↳ line 253, in check_call
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ raise MXNetError(py_str(_LIB.MXGetLastError()))
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ mxnet.base.MXNetError: array::at
[INFO ] W-9000-squeezenet_v1.1_compiled-stdout com.amazonaws.ml.mms.wlm.WorkerLifeCycle -
↳ Invoking custom service failed.
```

## MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded

NOTE: This issue is resolved in version 1.5.1.1.1.88.0 released 11/17/2020 and only applies for MXNet 1.5.

MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server. Restarting the model server may fail with “Failed to create NEURONCORE\_GROUP” error:

```
mxnet.base.MXNetError: [00:26:59] src/operator/subgraph/neuron/./neuron_util.h:348:
↳ Check failed: 0 == create_eg_rsp.status().code(): Failed to create NEURONCORE_GROUP
↳ with Neuron-RTD Error. Neuron-RTD Status Code: 9, details: ""
```

The workaround is to run “/opt/aws/neuron/bin/neuron-cli reset” to clear Neuron RTD states after all models are unloaded and server is shut down before restarting the model server.

## Pipeline mode is not able to execute inferences requests in parallel

If you see that multiple executors in a neuron pipeline setup (one model compiled for more than one neuron-cores using `-neuroncore-pipeline-cores` option during compilation) are not running in parallel, please set the following MXNet's environment variables before inference to allow mxnet to execute the CPU ops in parallel. Otherwise it will be sequential and stall the executors.

`MXNET_CPU_WORKER_NTHREADS` is used to do that. Setting its value to `__subgraph_opt_neuroncore__` in the compiled model json will ensure that all the executors (threads) can be run in parallel.

## Features only in MXNet-Neuron 1.5

- Shared memory for IFMaps transfer to neuron runtime (has higher performance compared to GRPC mode)
- Neuron profiling using MXNet

## Features only in MXNet-Neuron 1.8

- Gluon API support
- Library mode neuron runtime

*This document is relevant for: Inf1*

*This document is relevant for: Inf1*

## Apache MXNet Neuron Release Notes

### Table of contents

- [Apache MXNet Neuron release \[1.8.0.2.4.40.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.4.25.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.4.10.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.4.9.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.4.1.0\]](#)
- [\[1.5.1.1.10.39.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.2.127.0\]](#)
- [\[1.5.1.1.10.37.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.2.43.0\]](#)
- [\[1.5.1.1.10.11.0\]](#)
- [\[1.5.1.1.10.0.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.2.2.0\]](#)
- [\[1.5.1.1.9.0.0\]](#)
- [Apache MXNet Neuron release \[1.8.0.2.1.5.0\]](#)

- *Apache MXNet Neuron release [1.8.0.2.0.276.0]*
- *Apache MXNet Neuron release [1.8.0.2.0.271.0]*
- *[1.5.1.1.7.0.0]*
- *[1.5.1.1.6.5.0]*
- *[1.8.0.1.3.4.0]*
- *[1.5.1.1.6.1.0]*
- *[1.8.0.1.3.0.0]*
- *[1.8.0.1.2.1.0]*
- *[1.8.0.1.1.2.0]*
- *[1.5.1.1.4.x.x]*
- *[1.5.1.1.4.4.0]*
- *[1.5.1.1.3.8.0]*
- *[1.5.1.1.3.7.0]*
- *[1.5.1.1.3.2.0]*
- *[1.5.1.1.2.1.0]*
- *[1.5.1.1.1.88.0]*
- *[1.5.1.1.1.52.0]*
- *[1.5.1.1.1.1.0]*
- *[1.5.1.1.0.2101.0]*
- *[1.5.1.1.0.2093.0]*
- *[1.5.1.1.0.2033.0]*
- *[1.5.1.1.0.1900.0]*
- *[1.5.1.1.0.1596.0]*
- *[1.5.1.1.0.1498.0]*
- *[1.5.1.1.0.1401.0]*
- *[1.5.1.1.0.1325.0]*
- *[1.5.1.1.0.1349.0]*
- *[1.5.1.1.0.1325.0]*
- *[1.5.1.1.0.1260.0]*

This document lists the release notes for MXNet-Neuron framework.

**Apache MXNet Neuron release [1.8.0.2.4.40.0]**

Date: 12/21/2023

**Summary**

Minor updates.

**Apache MXNet Neuron release [1.8.0.2.4.25.0]**

Date: 10/15/2023

**Summary**

Minor updates.

**Apache MXNet Neuron release [1.8.0.2.4.10.0]**

Date: 7/19/2023

**Summary**

Minor bug fixes and enhancements for MXNet 1.8 Neuron.

**Apache MXNet Neuron release [1.8.0.2.4.9.0]**

Date: 6/14/2023

**Summary**

Minor bug fixes and enhancements for MXNet 1.8 Neuron.

**Apache MXNet Neuron release [1.8.0.2.4.1.0]**

Date: 5/1/2023

### New in this release

- Updated Neuron Runtime library to version 2.12
- Added missing LICENSE.txt

### Known Issues and Limitations

- Bert-base in 16 NeuronCores pipeline mode has 50% lower performance when running 16 inferences in parallel with Runtime version 2.12.

### [1.5.1.1.10.39.0]

Date: 5/1/2023

### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

This is the last released version. Please use neuron-cc version 1.15.0 only for this mxnet-neuron version. Also, this version is limited to python 3.9 or below only.

```
python -m pip install mxnet_neuron==1.5.1.* neuron-cc==1.15.0
```

### Apache MXNet Neuron release [1.8.0.2.2.127.0]

Date: 3/28/2023

### Summary

Minor bug fixes and enhancements for MXNet 1.8 Neuron.

### [1.5.1.1.10.37.0]

Date: 3/28/2023

### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

## Apache MXNet Neuron release [1.8.0.2.2.43.0]

Date: 11/23/2022

### Summary

Minor bug fixes and enhancements for MXNet 1.8 Neuron.

## [1.5.1.1.10.11.0]

Date: 11/23/2022

### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

## [1.5.1.1.10.0.0]

Date: 04/28/2022

### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

## Apache MXNet Neuron release [1.8.0.2.2.2.0]

Date: 03/25/2022

### New in this release

- Added support for unloading models from a NeuronDevice by deleting the model instance in user application. Users can now call `del` in Python on an executor and to unload the model from a NeuronDevice (provided the deleted executor is the last executor pointing to the given model). This requires the latest `aws-mx-1.8` package from [https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws\\_mx-1.8.0.2-py2.py3-none-manylinux2014\\_x86\\_64.whl](https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws_mx-1.8.0.2-py2.py3-none-manylinux2014_x86_64.whl).

### Bug fixes

- Fixed a memory leak caused by stale unloaded models in NeuronDevice memory. For this fix to take effect please install `aws-mx` package from [https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws\\_mx-1.8.0.2-py2.py3-none-manylinux2014\\_x86\\_64.whl](https://aws-mx-pypi.s3.us-west-2.amazonaws.com/1.8.0/aws_mx-1.8.0.2-py2.py3-none-manylinux2014_x86_64.whl) along with the latest `mx-neuron` package.

### [1.5.1.1.9.0.0]

Date: 03/25/2022

#### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

### Apache MXNet Neuron release [1.8.0.2.1.5.0]

Date: 01/20/2022

#### New in this release

- Added support of `mx_neuron.__version__` to get the build version of MXNet Neuron plugin

#### Bug fixes

- Fixed assertion errors when inference was completed with NaNs. The expected behavior is to complete inference successfully and warn the user that ``NaN``s were seen during the current inference.
- Fixed compile issue when individual output nodes have multiple output nodes. Because the output index was being dropped, fewer number of output feature maps were being considered and that caused failures during inference.

### Apache MXNet Neuron release [1.8.0.2.0.276.0]

Date: 11/05/2021

- Updated Neuron Runtime (which is integrated within this package) to `libnrt 2.2.18.0` to fix a container issue that was preventing the use of containers when `/dev/neuron0` was not present. See details here [neuron-runtime-release-notes](#).

### Apache MXNet Neuron release [1.8.0.2.0.271.0]

Date 10/27/2021

#### New in this release

- MXNet Neuron 1.8 now support Neuron Runtime 2.x (`libnrt.so` shared library) only.

---

#### Important:

- You must update to the latest Neuron Driver (`aws-neuron-dkms` version 2.1 or newer) for proper functionality of the new runtime library.
- Read *Introducing Neuron Runtime 2.x (libnrt.so)* application note that describes *why are we making this change* and how *this change will affect the Neuron SDK* in detail.



- Read *Migrate your application to Neuron Runtime 2.x (libnrt.so)* for detailed information of how to migrate your application.
- 

- Introducing Flexible Execution Groups (FlexEG) feature. See *Flexible Execution Group (FlexEG) in Neuron-MXNet* application note.

## Resolved Issues

- Fixed a bug that prevented compilation of gluon models with multiple cpu and neuron nodes.
- Added more debug logic to help with profiling of model load timing.

### [1.5.1.1.7.0.0]

Date 10/27/2021

## New in this release

- MXNet 1.5 enters maintenance mode. Please visit `maintenance_mxnet_1_5` for more information.

## Resolved Issues

- Minor bug fixes.

### [1.5.1.1.6.5.0]

Date 08/12/2021

## Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

### [1.8.0.1.3.4.0]

Date 08/12/2021

## Summary

Minor bug fixes and enhancements for MXNet 1.8 Neuron.

### [1.5.1.1.6.1.0]

Date 07/02/2021

#### Summary

Minor bug fixes and enhancements for MXNet 1.5 Neuron.

### [1.8.0.1.3.0.0]

Date 07/02/2021

#### Summary

Support for Autoloop, Cpredict API and minor bug fixes and enhancements for MXNet 1.8 Neuron.

#### Major New Features

- Added support for Autoloop feature for MXNet 1.8 Neuron.

#### Resolved Issues

- Added support for CPredict API.

### [1.8.0.1.2.1.0]

Date 5/28/2021

#### Summary

Minor bug fixes and enhancements for MXNet 1.8 Neuron

#### Resolved Issues

- Added support for Neuron profiler

**[1.8.0.1.1.2.0]**

Date 4/30/2021

**Summary**

Initial release of Apache MXNet 1.8 for Neuron

**Major New Features**

- Gluon API and Neuron support for NLP BERT models
- Neuron is now a plugin
- Please note new API changes to support plugin mode: *Neuron Apache MXNet Compilation Python API*

**[1.5.1.1.4.x.x]**

Date 5/28/2021

**Summary**

- Minor enhancements.

**[1.5.1.1.4.4.0]**

Date 4/30/2021

**Summary**

- Resolve an issue with Neuron profiling.

**Resolved Issues**

- Issue: when Neuron profiling is enabled in MXNet-Neuron 1.5.1 (using NEURON\_PROFILE=<dir>), and TensorBoard is used to read in the profiled data, user would see an error message “panic: runtime error: index out of range”. This issue is resolved in this release.

**[1.5.1.1.3.8.0]**

Date 3/4/2021

**Summary**

Minor enhancements.

**[1.5.1.1.3.7.0]**

Date 2/24/2021

**Summary**

Fix for CVE-2021-3177.

**[1.5.1.1.3.2.0]**

Date 1/30/2021

**Summary**

Various minor improvements

**[1.5.1.1.2.1.0]**

Date 12/23/2020

**Summary**

Various minor improvements

**[1.5.1.1.1.88.0]**

Date 11/17/2020

## Summary

This release includes the bug fix for MXNet Model Server not being able to clean up Neuron RTD states after model is unloaded (deleted) from model server.

## Resolved Issues

- Issue: MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server.
  - Workaround for earlier versions: run “/opt/aws/neuron/bin/neuron-cli reset” to clear Neuron RTD states after all models are unloaded and server is shut down.

### [1.5.1.1.1.52.0]

Date 09/22/2020

## Summary

Various minor improvements.

## Major New Features

## Resolved Issues

- Issue: When first importing MXNet into python process and subprocess call is invoked, user may get an OS-Error exception “OSError: [Errno 14] Bad address” during subprocess call (see <https://github.com/apache/incubator-mxnet/issues/13875> for more details). This issue is fixed with a mitigation patch from MXNet for Open-MP fork race conditions.
  - Workaround for earlier versions: Export KMP\_INIT\_AT\_FORK=false before running python process.

### [1.5.1.1.1.1.0]

Date 08/08/2020

## Summary

Various minor improvements.

### Major New Features

### Resolved Issues

[1.5.1.1.0.2101.0]

Date 08/05/2020

### Summary

Various minor improvements.

### Major New Features

### Resolved Issues

[1.5.1.1.0.2093.0]

Date 07/16/2020

### Summary

This release contains a few bug fixes and user experience improvements.

### Major New Features

### Resolved Issues

- User can specify NEURONCORE\_GROUP\_SIZES without brackets (for example, “1,1,1,1”), as can be done in TensorFlow-Neuron and PyTorch-Neuron.
- Fixed a memory leak when inferring neuron subgraph properties
- Fixed a bug dealing with multi-input subgraphs

[1.5.1.1.0.2033.0]

Date 6/11/2020

## Summary

- Added support for profiling during inference

## Major New Features

- Profiling can now be enabled by specifying the profiling work directory using `NEURON_PROFILE` environment variable during inference. For an example of using profiling, see `tensorboard-neuron`. (Note that graph view of MXNet graph is not available via TensorBoard).

## Resolved Issues

## Known Issues and Limitations

## Other Notes

[1.5.1.1.0.1900.0]

Date 5/11/2020

## Summary

Improved support for shared-memory communication with Neuron-Runtime.

## Major New Features

- Added support for the BERT-Base model (base: L-12 H-768 A-12), max sequence length 64 and batch size of 8.
- Improved security for usage of shared-memory for data transfer between framework and Neuron-Runtime
- Improved allocation and cleanup of shared-memory resource
- Improved container support by automatic falling back to GRPC data transfer if shared-memory cannot be allocated by Neuron-Runtime

## Resolved Issues

- User is unable to allocate Neuron-Runtime shared-memory resource when using MXNet-Neuron in a container to communicate with Neuron-Runtime in another container. This is resolved by automatic falling back to GRPC data transfer if shared-memory cannot be allocated by Neuron-Runtime.
- Fixed issue where some large models could not be loaded on inferentia.

### Known Issues and Limitations

### Other Notes

[1.5.1.1.0.1596.0]

Date 3/26/2020

### Summary

No major changes or fixes

### Major New Features

### Resolved Issues

### Known Issues and Limitations

### Other Notes

[1.5.1.1.0.1498.0]

Date 2/27/2020

### Summary

No major changes or fixes.

### Major New Features

### Resolved Issues

The issue(s) below are resolved:

- Latest pip version 20.0.1 breaks installation of MXNet-Neuron pip wheel which has py2.py3 in the wheel name.

### Known Issues and Limitations

- User is unable to allocate Neuron-Runtime shared-memory resource when using MXNet-Neuron in a container to communicate with Neuron-Runtime in another container. To work-around, please set environment variable NEURON\_RTD\_USE\_SHM to 0.



## Other Notes

[1.5.1.1.0.1401.0]

Date 1/27/2020

## Summary

No major changes or fixes.

## Major New Features

## Resolved Issues

- The following issue is resolved when the latest multi-model-server with version  $\geq 1.1.0$  is used with MXNet-Neuron. You would still need to use “`/opt/aws/neuron/bin/neuron-cli reset`” to clear all Neuron RTD states after multi-model-server is exited:
  - Issue: MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server and previous workaround “`/opt/aws/neuron/bin/neuron-cli reset`” is unable to clear all Neuron RTD states.

## Known Issues and Limitations

- Latest pip version 20.0.1 breaks installation of MXNet-Neuron pip wheel which has py2.py3 in the wheel name. This breaks all existing released versions. The error looks like:

```
Looking in indexes: https://pypi.org/simple, https://pip.repos.neuron.amazonaws.com
ERROR: Could not find a version that satisfies the requirement mxnet-neuron (from
  versions: none)
ERROR: No matching distribution found for mxnet-neuron
```

- Work around: install the older version of pip using “`pip install pip==19.3.1`”.

## Other Notes

[1.5.1.1.0.1325.0]

Date 12/1/2019

### Summary

### Major New Features

### Resolved Issues

- Issue: Compiler flags cannot be passed to compiler during compile call. The fix: compiler flags can be passed to compiler during compile call using “flags” option followed by a list of flags.
- Issue: Advanced CPU fallback option is a way to attempt to improve the number of operators on Inferentia. The default is currently set to on, which may cause failures. The fix: This option is now off by default.

### Known Issues and Limitations

- Issue: MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server and previous workaround “/opt/aws/neuron/bin/neuron-cli reset” is unable to clear all Neuron RTD states.
  - Workaround: run “sudo systemctl restart neuron-rtd” to clear Neuron RTD states after all models are unloaded and server is shut down.

### Other Notes

[1.5.1.1.0.1349.0]

Date 12/20/2019

### Summary

No major changes or fixes. Released with other Neuron packages.

[1.5.1.1.0.1325.0]

Date 12/1/2019

### Summary

### Major New Features

### Resolved Issues

- Issue: Compiler flags cannot be passed to compiler during compile call. The fix: compiler flags can be passed to compiler during compile call using “flags” option followed by a list of flags.
- Issue: Advanced CPU fallback option is a way to attempt to improve the number of operators on Inferentia. The default is currently set to on, which may cause failures. The fix: This option is now off by default.

## Known Issues and Limitations

- Issue: MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server and previous workaround “`/opt/aws/neuron/bin/neuron-cli reset`” is unable to clear all Neuron RTD states.
  - Workaround: run “`sudo systemctl restart neuron-rtd`” to clear Neuron RTD states after all models are unloaded and server is shut down.

## Other Notes

[1.5.1.1.0.1260.0]

Date: 11/25/2019

## Summary

This version is available only in released DLAMI v26.0 and is based on MXNet version 1.5.1. Please dlami-rn-known-issues to latest version.

## Major new features

## Resolved issues

## Known issues and limitations

- Issue: Compiler flags cannot be passed to compiler during compile call.
- Issue: Advanced CPU fallback option is a way to attempt to improve the number of operators on Inferentia. The default is currently set to on, which may cause failures.
  - Workaround: explicitly turn it off by setting compile option `op_by_op_compiler_retry` to 0.
- Issue: Temporary files are put in current directory when debug is enabled.
  - Workaround: create a separate work directory and run the process from within the work directory
- Issue: MXNet Model Server is not able to clean up Neuron RTD states after model is unloaded (deleted) from model server.
  - Workaround: run “`/opt/aws/neuron/bin/neuron-cli reset`” to clear Neuron RTD states after all models are unloaded and server is shut down.
- Issue: MXNet 1.5.1 may return inconsistent node names for some operators when they are the primary outputs of a Neuron subgraph. This causes failures during inference.
  - Workaround : Use the `excl_node_names` compilation option to change the partitioning of the graph during compile so that these nodes are not the primary output of a neuron subgraph. See [Neuron Apache MXNet Compilation Python API](#)

```
compile_args = { 'excl_node_names': ["node_name_to_exclude"] }
```

### Models Supported

The following models have successfully run on neuron-inferentia systems

1. Resnet50 V1/V2
2. Inception-V2/V3/V4
3. Parallel-WaveNet
4. Tacotron 2
5. WaveRNN

### Other Notes

- Python versions supported:
  - 3.5, 3.6, 3.7
- Linux distribution supported:
  - Ubuntu 18, Amazon Linux 2

*This document is relevant for: **Inf1***

- *[Troubleshooting Guide for Neuron Apache MXNet](#)*
- *[What's New](#)*
- *[Neuron Apache MXNet Supported operators](#)*

*This document is relevant for: **Inf1***

Setup (`mxnet-neuron`)

### Tutorials

#### Computer Vision Tutorials

- ResNet-50 tutorial [\[html\]](#) [\[notebook\]](#)
- Model Serving tutorial [\[html\]](#)
- Getting started with Gluon tutorial [\[html\]](#) [\[notebook\]](#)

#### Natural Language Processing (NLP) Tutorials

- MXNet 1.8: Using data parallel mode tutorial [\[html\]](#) [\[notebook\]](#)

## Utilizing Neuron Capabilities Tutorials

- [NeuronCore Groups tutorial \[html\]](#) [\[notebook\]](#)

---

**Note:** To use Jupyter Notebook see:

- [setup-jupyter-notebook-steps-troubleshooting](#)
  - [running-jupyter-notebook-as-script](#)
- 

## API Reference Guide

- *Neuron Apache MXNet Compilation Python API*

## Developer Guide

- *Flexible Execution Group (FlexEG) in Neuron-MXNet*

## Misc

- *Troubleshooting Guide for Neuron Apache MXNet*
- *What's New*
- *Neuron Apache MXNet Supported operators*

*This document is relevant for:* **Inf1**

MxNet Neuron(`mxnet-neuron`) for Inference on **Inf1**

*This document is relevant for:* **Inf1**



## ABOUT NEURON

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 10.1 Release Details

#### 10.1.1 Latest Release

- *What's New*
- Release Artifacts

#### 10.1.2 Previous Releases

- prev-rn
- pre-release-content
- prev-n1-rn

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 10.2 Roadmap

The AWS Neuron feature roadmap provides visibility onto what we are working on in terms of functional and performance in the near future. We hope this will help you better plan how to use Neuron with your products. We'd love to get our customers feedback as well, to help us ensure we are working on the most important requests.

Neuron Roadmap items are classified into one of the following status:

| Roadmap Status        | Description   |
|-----------------------|---|
| No Status             | Items that are in the roadmap backlog               |
| Todo                  | Items that are planned for future releases          |
| In Progress           | Items that are planned for the next release         |
| Done                  | Items that have been completed in previous releases |
| Not to be Implemented | Items that are not planned for future releases      |

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

## 10.3 Support

*This document is relevant for: Inf1, Inf2, Trn1, Trn2*

### 10.3.1 Neuron Software Maintenance policy

#### Table of Contents

- *Overview*
- *Neuron Software Definitions*
  - *Neuron SDK*
  - *Neuron components*
    - \* *Neuron extension components*
    - \* *Neuron standalone components*
  - *Neuron Model Classes*
  - *Neuron features*
  - *Neuron APIs*
  - *Dependency software components*
  - *Neuron Deep Learning AMIs and Deep Learning Containers*
- *Neuron Software Life-cycle*
  - *Neuron Software Regular Updates*
  - *Neuron Software Maintenance Updates*
  - *Neuron DLAMIs and DLCs Updates*
  - *Neuron Extension Components Updates*
  - *Communication methods*
- *Neuron Software Versioning*
  - *Neuron SDK Documentation Versioning*
  - *Neuron SDK Versioning*
  - *Neuron extension components Versioning*
  - *Neuron Standalone Component Versioning*
- *Neuron Software Release Types*
  - *Major release*
  - *Minor release*
  - *Patch release*
  - *Pre-releases*
- *Neuron Software Classification*
  - *Neuron SDK and Neuron components*



- *Neuron APIs*
- *Neuron Features*
- *Neuron Supported Model Classes*

## Overview

This document outlines software maintenance policy for AWS Neuron Software Development Kit (SDK), Neuron Components, both extension and standalone components, supported model classes, features, APIs, DLAMIs and DLCs, and dependency software. AWS Neuron is the SDK for Amazon EC2 [Inferentia](#) and Amazon EC2 [Trainium](#) based instances purpose-built for deep learning. Neuron integrates with popular Machine Learning (ML) frameworks like PyTorch, JAX, and TensorFlow and includes a compiler, runtime, driver, profiling tools, and libraries to support high performance training of generative AI models on Trainium and Inferentia powered instances.

This document addresses Neuron Software life-cycle and the Neuron SDK release versioning.

## Neuron Software Definitions

Neuron Software refers to the complete set of software elements provided by AWS Neuron, including:

### Neuron SDK

The core software development kit that enables users to build, train, and deploy machine learning models on Inferentia and Trainium based instances. The Neuron SDK encompasses the entire set of components, features, APIs, and other elements that are bundled together and made available in a particular version of the Neuron SDK release.

### Neuron components

Neuron components refer to any packages or libraries within the Neuron SDK that offer specific functionality. These components are typically accessible through PIP, RPM, or Debian packages for easy installation and usage. There are two main categories of Neuron components: Neuron extension components and Neuron standalone components.

### Neuron extension components

Neuron extension components are components that integrate Neuron support into open source machine learning frameworks, libraries or tools enhancing their functionality and extending their capabilities as necessary. When referring to Neuron extension components, we are also referring to the parts of the open source machine learning framework or library that are supported by Neuron. The software life-cycle of the open source machine learning frameworks, libraries or tools that are extended by Neuron is managed and maintained by their respective communities or the vendors responsible for those specific components. Examples for Neuron extension components are:

- **Third party ML Library:** Examples include Neuron Nemo Megatron.
- **Third party ML Framework:** Examples include PyTorch NeuronX and TensorFlow Neuron.

### Neuron standalone components

Neuron standalone components are self-contained components within the Neuron SDK. Examples of such components are Neuron Compiler, Neuron Tools and Neuron Runtime.

### Neuron Model Classes

A Neuron supported model class is tightly coupled with a specific Neuron extension component (e.g. PyTorch NeuronX) or Neuron library (e.g. NeuronX Distributed) and the workload type (e.g. Training or Inference). For example a model can be supported at Beta level in PyTorch NeuronX for training and Stable level in PyTorch NeuronX for inference.

### Neuron features

A Neuron feature refers to any functionality or attribute that is part of the Neuron SDK, whether it belongs to the entire Neuron SDK or to one of its specific components.

For example, a Neuron feature is [Neuron Persistent Cache](#) in the *Transformers Neuronx library*

### Neuron APIs

A Neuron API refers to any API, CLI, environment variables, or flag that belong to the entire Neuron SDK or to one of the Neuron components. A Neuron API allows developers to interact with and leverage the capabilities of the Neuron SDK and its components.

Examples include *Neuron Trace API* and *Neuron Compiler flags*

### Dependency software components

External software components or frameworks that the Neuron SDK and its components rely on for proper functioning and compatibility, such as language runtimes or operating systems.

The software life-cycle of the dependency software components, is managed and maintained by their respective communities or the vendors responsible for those specific dependency software components. The following terms are examples of underlying dependency software components:

- **Operating System (OS):** Examples include Ubuntu 22 and Amazon Linux 2023
- **Language Runtime:** Examples include Python 3.10

### Neuron Deep Learning AMIs and Deep Learning Containers

*Neuron Deep Learning AMIs (DLAMIs)* and *Neuron Deep Learning Containers (DLCs)* are pre-configured Amazon Machine Images and Docker container that come with the Neuron SDK and necessary dependencies pre-installed, providing a ready-to-use environment for machine learning development.

## Neuron Software Life-cycle

The typical life-cycle for Neuron software consists of several phases, though not all phases are applicable to every type of Neuron software. The phases are as follows:

- **Developer Preview or Beta** (these terms are used interchangeably in Neuron collaterals)
- **Release Candidate (RC)**
- **General Availability (GA) or Stable** (these terms are used interchangeably in Neuron collaterals)
- **Maintenance**
- **End-of-Support (EOS)**

The following table outlines the details for each phase for Neuron software:

|                               | Description  | Comments  |
|-------------------------------|--|---|
| Developer Preview (Beta)      | In this phase, Neuron Software is not supported, should not be used in production environments, and is meant for early access and feedback purposes only. It is possible for future releases to introduce breaking changes. See <i>Neuron Software Classification</i> for more information   |   |
| Release Candidate (RC)        | Once AWS identifies a release to be a stable product, it may be marked as a Release Candidate (RC). This phase is usually short and during it AWS will provide for Neuron Software on an as-needed basis.  | This phase applies only to Neuron SDK and Neuron components   |
| General Availability (Stable) | During this phase, AWS releases regular for Neuron Software on an as-needed basis. See <i>Neuron Software Classification</i> for more information  |   |
| Maintenance                   | During the maintenance phase, AWS will provide <i>maintenance updates</i> for Neuron Software on an as-needed basis. Any new PIP, RPM, and Debian packages for the Neuron Software, as well as updated versions of the Neuron DLAMIs and Neuron DLCs, will be released only when deemed necessary by the AWS Neuron team. Users can expect updates to be less frequent compared to <i>regular</i> as the focus will be on addressing critical issues and ensuring the stability of the software.<br>Maintenance Announcement: AWS will make a public <i>announcement</i> at least one month before the Neuron Software enters Maintenance phase.   | This phase does not apply to Dependency Software Components, Neuron DLCs, Neuron DLAMIs, Neuron Features and APIs |
| End of Support (EOS)          | When Neuron Software reaches the end of its support lifecycle, it will no longer receive <i>regular</i> updates and <i>maintenance updates</i> (including security updates). While AWS will continue to provide access to all previously released PIP, RPM, and Debian packages for the Neuron Software, as well as earlier versions of the Neuron DLAMIs and Neuron DLCs, it's important to note that these older versions will not receive any updates or support. Customers can still use these resources at their own discretion, but it is highly recommended to upgrade to the latest available versions<br>End of Support Announcement: AWS will make a public <i>announcement</i> at least one month before a Neuron Software enters End of Support. |   |

## Neuron Software Regular Updates

Regular updates for Neuron Software address the following areas: new features, feature improvements, performance enhancements, bug resolution, security vulnerability fixes, upgrades to Neuron dependency software components and upgrades to Neuron extension components. To handle these regular updates, AWS will release a new version of the Neuron SDK, incrementing the minor version (the second digit in the version number) for a minor release or incrementing the major version (the first digit in the version number) for a major release when significant changes that break compatibility are introduced. It's important to note that any bug-fixes or security issues in regular updates are not applied retroactively to previous versions of the Neuron SDK. To benefit from these updates, users must adopt the latest release.

For more information see:

- [Neuron DLAMIs and DLCs Updates](#)
- [Neuron Extension Components Updates](#)
- [Neuron Software Versioning](#)

**Neuron SDK Installation and Update instructions** To install and update to the latest Neuron packages, customers need to pin the major version of the Neuron package. For example, to install latest Neuron tools package, call `sudo apt-get install aws-neuronx-tools=2.*` and to install latest PyTorch Neuron package for Trn1, call `pip install torch-neuronx==2.1.0.1.*`. This is done to future-proof instructions for new, backwards-incompatible major version releases.

## Neuron Software Maintenance Updates

Maintenance updates for Neuron Software address three key areas: resolving bugs, fixing security vulnerabilities, and upgrading dependency software components. At AWS discretion, additional critical features or performance enhancement may also be included. To handle these maintenance updates, AWS will release a new version of the Neuron SDK, incrementing the patch number (the last digit in the version number) to indicate a patch release. Major or minor releases may also contain maintenance updates. It's important to note that these maintenance updates are not applied retroactively to previous versions of the Neuron SDK. To take advantage of these updates, users must adopt the latest patch release.

For more information see:

- [Neuron DLAMIs and DLCs Updates](#)
- [Neuron Extension Components Updates](#)
- [Neuron Software Versioning](#)

## Neuron DLAMIs and DLCs Updates

AWS will address *regular* updates, life-cycle changes, maintenance updates, and security issues related to any third-party software included in the Neuron DLAMI or DLCs by releasing new versions of the Neuron DLAMI or DLCs. However, updates won't be applied retroactively to older versions of the Neuron DLAMI or DLCs. Instead, users will need to use the new versions to get the latest updates. Generally, Neuron DLAMIs and Deep Learning Containers (DLCs) will support one latest LTS Linux Distribution version (Ubuntu, Amazon Linux, and Rocky9), with exceptions. Neuron Base DLAMIs (which come pre-installed with Neuron driver, EFA, and Neuron tools) will support the two latest versions of LTS Linux Distributions.

For more information see:

- [Neuron Extension Components Updates](#)
- [Neuron Software Versioning](#)

## Neuron Extension Components Updates

When a new version of an open source ML framework (e.g. PyTorch) is supported by a Neuron extension component (e.g., PyTorch NeuronX), the Neuron extension component for the latest supported ML framework version will become the default for installation. If users wish to use a Neuron extension component for an earlier supported ML framework version, they will need to explicitly specify the desired version during installation. After upgrading a Neuron extension component to support a newer version of an ML framework, AWS will continue to provide *regular updates* for the Neuron extension component that supports the earlier ML framework version for a minimum of 6 months. After the 6 months period, the Neuron extension component for the earlier supported ML framework version may transition into a maintenance mode. In the maintenance mode, updates for the older Neuron extension component versions will be provided on an as-needed basis, focusing on critical bug fixes and security patches. For more information see: [Neuron extension component versioning](#)

## Communication methods

Neuron software classification and lifecycle announcements are communicated as follows:

- Neuron SDK documentation under [Announcements](#)

To see the list of available Neuron SDK versions and supported dependency software components versions:

- Neuron SDK documentation under [Release Content](#)
- Neuron SDK documentation under [What's New](#)

## Neuron Software Versioning

### Neuron SDK Documentation Versioning

Neuron SDK documentation is versioned and maps to the corresponding Neuron SDK version. Users can switch to earlier versions of the Neuron SDK documentation by selecting the version from the dropdown in bottom left portion of the side bar.

### Neuron SDK Versioning

The AWS SDK release versions are in the form of [A.B.C] where (A) represents the major version, (B) represents the minor version, and (C) represents the patch version.

### Neuron extension components Versioning

Neuron extension components versioning (like PyTorch NeuronX) is in the form [X.Y.Z] . [A.B.C], where [X.Y.Z] represents the third party component's major (X), minor (Y), and patch (Z) versions and [A.B.C] represents the Neuron extension components (A), minor (B), and patch (C) versions.

### Neuron Standalone Component Versioning

Neuron Component versioning (except of Neuron extension components like PyTorch NeuronX) is in the form [A.B.C.D], where A represents the major version, B represents the minor version, and C.D represents the patch version.

### Neuron Software Release Types

#### Major release

Increasing the major version indicates that the Neuron software underwent significant and substantial changes in an incompatible manner. Applications need to be updated in order for them to work with the newest SDK version. It is important to update major versions carefully and in accordance with the upgrade guidelines provided by AWS. After increasing the major version, the Neuron software may not maintain compatibility with previous supported versions of *Neuron Runtime*, *Neuron Compiler*, and *NEFF*.

#### Minor release

Increasing the minor version indicates that the Neuron software added functionality in a backwards compatible manner.

#### Patch release

Increasing the patch version indicates that the Neuron software added backward compatible bug or security fixes. A bug fix is defined as an internal change that fixes incorrect behavior.

#### Pre-releases

- **Developer Preview (Beta):** During this phase, the Neuron software is not supported, should not be used in production environments, and is meant for early access and feedback purposes only. It is possible for future releases to introduce breaking changes. In the case of a Developer Preview (Beta) release, the minor version will include a lower case b along with a (Beta) tag.
- **Release Candidate (RC):** Once Neuron identifies a release to be a stable product, it may mark it as a Release Candidate. Release Candidates are ready for GA release unless significant bugs emerge, and will receive full AWS Neuron support. In the case of a RC release, the minor version will include a lower case rc along with a (RC) tag.

### Neuron Software Classification

This section explains the Neuron software classification for APIs, libraries, packages, features, and Neuron supported model classes mentioned in the Neuron documentation.

## Neuron SDK and Neuron components

|                          | Testing                  | Features   | Performance |
|--------------------------|--------------------------|--|-------------|
| Developer Preview (Beta) | Basic                    | Minimal Viable Product (MVP) *                   |             |
| Release Candidate (RC)   | Basic                    | Minimal Viable Product (MVP)*                    | Tested      |
| GA (Stable)              | Standard Product Testing | Incremental additions or changes in new releases | Tested      |

\* A minimum viable product (MVP) for a Neuron Component contains just enough features to be usable by early customers who can then provide feedback for future development. MVP can be different per use case and depends on the specific package/library of interest. Please note that in many cases, an MVP can also represent an advanced level of features.

## Neuron APIs

|                          | API Contract  | API Backward Compatibility |
|--------------------------|---|----------------------------|
| Alpha                    | Unstable and undocumented   | No                         |
| Developer Preview (Beta) | Major changes may happen  | No                         |
| GA (Stable)              | Incremental changes in new releases (without breaking the API contract) | Yes *                      |

\* In certain cases, when necessary, AWS may introduce API changes that may break compatibility, with notice provided ahead of time.

## Neuron Features

|                          | Testing                  | Functionality  | Performance             |
|--------------------------|--------------------------|--|-------------------------|
| Alpha                    | No formal testing done   | Partial functionality with limited set of core capabilities, far from Minimum Viable Product (MVP) * | Not tested or evaluated |
| Developer Preview (Beta) | Basic                    | Minimum Viable Product (MVP) *   |                         |
| GA (Stable)              | Standard Product Testing | Incremental additions or changes in new releases   | Tested                  |

\* A minimum viable product (MVP) for a Neuron Feature contains just enough functionality to be usable by early customers who can then provide feedback for future development. MVP can be different per use case and depends on the specific feature of interest. Please note that in many cases, an MVP can also represent an advanced level of functionality.

## Neuron Supported Model Classes

|                          | Accuracy / Convergence | Throughput / Latency |
|--------------------------|------------------------|----------------------|
| Developer Preview (Beta) | Validated              | Tested               |
| GA (Stable)              | Validated              | Tested               |

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### 10.3.2 Security Disclosures

If you think you've found a potential security issue, please do not post it in the Issues. Instead, please follow the instructions here (<https://aws.amazon.com/security/vulnerability-reporting/>) or email AWS security directly (<mailto:aws-security@amazon.com>).

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

### 10.3.3 Contact Us

For support please checkout the [Github issues](#) or [Neuron AWS forums](#) for an answer, if none of those resources have an answer to your question please open a ticket.

If you have an urgent need for a feature you can also contact us directly at [aws-neuron-support@amazon.com](mailto:aws-neuron-support@amazon.com).

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2

*This document is relevant for:* Inf1, Inf2, Trn1, Trn2



## PYTHON MODULE INDEX

p

placement, [221](#)



## Symbols

`__init__()` (*nki.isa.nc\_version method*), 1293

## A

`abs` (C++ function), 1514  
`abs()` (*in module nki.language*), 1205  
`abs_out` (C++ function), 1514  
`accessor` (C++ function), 1517  
`activation()` (*in module nki.isa*), 1248  
`activation_reduce()` (*in module nki.isa*), 1250  
`add` (C++ function), 1516  
`add()` (*in module nki.language*), 1197  
`add_out` (C++ function), 1515  
`affine_range()` (*in module nki.language*), 1232  
`affine_select()` (*in module nki.isa*), 1269  
`all()` (*in module nki.language*), 1204  
`all_reduce()` (*in module nki.language*), 1231  
`alloc()` (*in module nki.compiler.psum*), 1297  
`alloc()` (*in module nki.compiler.sbuf*), 1294  
`arange()` (*in module nki.language*), 1228  
`arctan()` (*in module nki.language*), 1211  
`assert_shape()` (*nki.tensor method*), 1184  
`astype()` (*nki.tensor method*), 1184  
`atomic_rmw()` (*in module nki.language*), 1190  
`auto_alloc()` (*in module nki.compiler.psum*), 1300  
`auto_alloc()` (*in module nki.compiler.sbuf*), 1297

## B

`baremetal()` (*in module nki*), 1181  
`benchmark()` (*in module nki*), 1178  
**BF16**, 1635  
`bfloat16` (*in module nki.language*), 1241  
`bitwise_and` (C++ function), 1516, 1517  
`bitwise_and()` (*in module nki.language*), 1219  
`bitwise_and_out` (C++ function), 1516  
`bitwise_not` (C++ function), 1517  
`bitwise_not_out` (C++ function), 1517  
`bitwise_or` (C++ function), 1517  
`bitwise_or()` (*in module nki.language*), 1220  
`bitwise_or_out` (C++ function), 1517  
`bitwise_xor()` (*in module nki.language*), 1220  
`bn_aggr()` (*in module nki.isa*), 1277

`bn_stats()` (*in module nki.isa*), 1275  
`broadcast_to()` (*in module nki.language*), 1191  
`broadcast_to()` (*nki.tensor method*), 1184  
built-in function  
    `torch.neuron.DataParallel()`, 215  
    `torch.neuron.DataParallel.disable_dynamic_batching()`, 90, 216  
    `torch_neuron.trace()`, 210  
    `torch_neuronx.analyze()`, 87  
    `torch_neuronx.bucket_model_trace()`, 77  
    `torch_neuronx.DataParallel()`, 90  
    `torch_neuronx.dynamic_batch()`, 78  
    `torch_neuronx.experimental.multicore_context()`, 86  
    `torch_neuronx.experimental.neuron_cores_context()`, 85  
    `torch_neuronx.experimental.profiler.profile()`, 347  
    `torch_neuronx.experimental.profiler.profile.start()`, 348  
    `torch_neuronx.experimental.set_multicore()`, 85  
    `torch_neuronx.experimental.set_neuron_cores()`, 84  
    `torch_neuronx.move_trace_to_device()`, 75  
    `torch_neuronx.PartitionerConfig()`, 79  
    `torch_neuronx.replace_weights()`, 81  
    `torch_neuronx.trace()`, 71

## C

**CCE**, 1634  
`ceil` (C++ function), 1514  
`ceil()` (*in module nki.language*), 1207  
`ceil_out` (C++ function), 1514  
**cFP8**, 1635  
`clamp` (C++ function), 1515  
`clamp_out` (C++ function), 1515  
`close` (C++ function), 1521  
**Collective Communication Engine**, 1633  
`copy()` (*in module nki.language*), 1190  
`cos` (C++ function), 1515  
`cos()` (*in module nki.language*), 1209

`cos_out` (C++ function), 1515

`CustomOps`, 1635

## D

`device_print`() (in module *nki.language*), 1239

`dge_mode` (class in *nki.isa*), 1292

`div` (C++ function), 1516

`div_out` (C++ function), 1516

`divide`() (in module *nki.language*), 1199

`dma_copy`() (in module *nki.isa*), 1279

`DP`, 1635

`DPr`, 1635

`dropout`() (in module *nki.isa*), 1268

`dropout`() (in module *nki.language*), 1217

`ds`() (in module *nki.language*), 1227

`dtype` (*nki.tensor* property), 1184

## E

`empty` (C++ function), 1514

`empty_like`() (in module *nki.language*), 1193

`enable_stack_allocator`() (in module *nki.compiler*), 1300

`engine` (class in *nki.isa*), 1291

`equal`() (in module *nki.language*), 1222

`erf`() (in module *nki.language*), 1214

`erf_dx`() (in module *nki.language*), 1215

`exp` (C++ function), 1515

`exp`() (in module *nki.language*), 1208

`exp_out` (C++ function), 1515

`expand_dims`() (in module *nki.language*), 1229

`expand_dims`() (*nki.tensor* method), 1184

`eye` (C++ function), 1514

## F

`fill_` (C++ function), 1518

`FLOAT32_TO_FLOAT16` (*torch\_neuron.Optimization* attribute), 211

`float8_e4m3` (in module *nki.language*), 1241

`float8_e5m2` (in module *nki.language*), 1241

`floor` (C++ function), 1514

`floor`() (in module *nki.language*), 1206

`floor_out` (C++ function), 1514

`fmod`() (in module *nki.language*), 1208

`force_auto_alloc`() (in module *nki.compiler*), 1301

`FP16`, 1635

`FP32`, 1635

`fp32` (class in *nki.language*), 1242

`full` (C++ function), 1514

`full`() (in module *nki.language*), 1194

## G

`gather_flattened`() (in module *nki.language*), 1230

`gelu`() (in module *nki.language*), 1213

`gelu_apprx_tanh`() (in module *nki.language*), 1213

`gelu_dx`() (in module *nki.language*), 1213

`get_accessor_coherence_policy` (C++ function), 1517

`get_cpu_count` (C++ function), 1525

`get_cpu_id` (C++ function), 1524

`get_dst_tensor` (C++ function), 1524

`get_nc_version`() (in module *nki.isa*), 1293

`GPSIMD Engine`, 1633

`GpSimdE`, 1634

`greater`() (in module *nki.language*), 1223

`greater_equal`() (in module *nki.language*), 1224

## H

`HBM`, 1634

`hbm` (in module *nki.language*), 1235

`High Bandwidth Memory`, 1633

## I

`Inf1`, 1632

`Inf2`, 1632

`Inferentia`, 1632

`invert`() (in module *nki.language*), 1221

`iota`() (in module *nki.isa*), 1266

`itemsize` (*nki.tensor* property), 1185

## J

`jit`() (in module *nki*), 1178

## L

`left_shift`() (in module *nki.language*), 1221

`less`() (in module *nki.language*), 1224

`less_equal`() (in module *nki.language*), 1225

`load`() (in module *nki.language*), 1186

`load_transpose2d`() (in module *nki.language*), 1189

`local_gather`() (in module *nki.isa*), 1277

`log` (C++ function), 1515

`log`() (in module *nki.language*), 1209

`log10` (C++ function), 1515

`log10_out` (C++ function), 1515

`log2` (C++ function), 1515

`log2_out` (C++ function), 1515

`log_out` (C++ function), 1515

`logical_and`() (in module *nki.language*), 1225

`logical_not`() (in module *nki.language*), 1227

`logical_or`() (in module *nki.language*), 1226

`logical_xor`() (in module *nki.language*), 1226

`loop_reduce`() (in module *nki.language*), 1239

## M

`matmul`() (in module *nki.language*), 1218

`max`() (in module *nki.language*), 1201

`max8`() (in module *nki.isa*), 1283

maximum() (in module *nki.language*), 1200  
 mean() (in module *nki.language*), 1202  
 memset() (in module *nki.isa*), 1274  
 mgrid (in module *nki.language*), 1228  
 min() (in module *nki.language*), 1202  
 minimum() (in module *nki.language*), 1201  
 mish() (in module *nki.language*), 1215  
 mod() (in module *nki.language*), 1207  
 mod\_alloc() (in module *nki.compiler.psum*), 1298  
 mod\_alloc() (in module *nki.compiler.sbuf*), 1295  
 module  
     placement, 221  
 mul (C++ function), 1516  
 mul\_out (C++ function), 1516  
 multiply() (in module *nki.language*), 1199

## N

NC, 1634  
 nc (in module *nki.language*), 1238  
 nc\_find\_index8() (in module *nki.isa*), 1284  
 nc\_match\_replace8() (in module *nki.isa*), 1285  
 nc\_matmul() (in module *nki.isa*), 1244  
 nc\_stream\_shuffle() (in module *nki.isa*), 1289  
 nc\_transpose() (in module *nki.isa*), 1247  
 nc\_version (class in *nki.isa*), 1293  
 ND, 1634  
 ndarray() (in module *nki.language*), 1192  
 ndim (*nki.tensor* property), 1185  
 negative() (in module *nki.language*), 1205  
 Neuron Device, 1632  
 Neuron Kernel Interface, 1634  
 neuron-cc  
     neuron-cc command line option, 1140, 1143  
 neuron-cc command line option  
     neuron-cc, 1140, 1143  
 neuron-ls  
     neuron-ls command line option, 1014  
 neuron-ls command line option  
     neuron-ls, 1014  
 neuron-monitor  
     neuron-monitor command line option, 996  
 neuron-monitor command line option  
     neuron-monitor, 996  
 neuron-profile  
     neuron-profile command line option, 1061  
 neuron-profile command line option  
     neuron-profile, 1061  
 NeuronCore, 1633, 1634  
 NeuronCore-v1, 1633  
 NeuronCore-v2, 1633  
 NeuronCore-v3, 1633  
 NeuronDevice, 1634  
 NeuronLink, 1633  
 NeuronLink-v1, 1633

NeuronLink-v2, 1633  
 NeuronLink-v3, 1633  
 neuronx-cc  
     neuronx-cc command line option, 1124, 1126  
 neuronx-cc command line option  
     neuronx-cc, 1124, 1126  
 NKI, 1635  
 not\_equal() (in module *nki.language*), 1223  
 nrt\_add\_tensor\_to\_tensor\_set (C function), 939  
 nrt\_allocate\_tensor\_set (C function), 939  
 nrt\_close (C function), 933  
 nrt\_destroy\_tensor\_set (C function), 939  
 nrt\_execute (C function), 940  
 nrt\_execute\_repeat (C function), 940  
 nrt\_free\_model\_tensor\_info (C function), 936  
 nrt\_get\_model\_instance\_count (C function), 936  
 nrt\_get\_model\_nc\_count (C function), 935  
 nrt\_get\_model\_tensor\_info (C function), 935  
 nrt\_get\_tensor\_from\_tensor\_set (C function), 939  
 nrt\_get\_total\_nc\_count (C function), 941  
 nrt\_get\_version (C function), 941  
 nrt\_get\_visible\_nc\_count (C function), 941  
 nrt\_init (C function), 932  
 nrt\_load (C function), 934  
 nrt\_load\_collectives (C function), 934  
 nrt\_profile\_start (C function), 940  
 nrt\_profile\_stop (C function), 940  
 nrt\_tensor\_allocate (C function), 937  
 nrt\_tensor\_allocate\_empty (C function), 938  
 nrt\_tensor\_allocate\_slice (C function), 938  
 nrt\_tensor\_attach\_buffer (C function), 938  
 nrt\_tensor\_free (C function), 937  
 nrt\_tensor\_get\_size (C function), 938  
 nrt\_tensor\_get\_va (C function), 939  
 nrt\_tensor\_read (C function), 937  
 nrt\_tensor\_write (C function), 938  
 nrt\_unload (C function), 935  
 num\_programs() (in module *nki.language*), 1236  
 NxD Core, 1634  
 NxD Inference, 1634  
 NxD Training, 1634

## O

ones (C++ function), 1514  
 ones() (in module *nki.language*), 1194  
 operator= (C++ function), 1518

## P

par\_dim (in module *nki.language*), 1235  
 Partial Sum Buffer, 1633  
 placement  
     module, 221  
 pow (C++ function), 1515  
 pow\_out (C++ function), 1515

`power()` (in module `nki.language`), 1200  
`PP`, 1636  
`PPr`, 1636  
`private_hbm` (in module `nki.language`), 1235  
`prod()` (in module `nki.language`), 1204  
`profile()` (in module `nki`), 1180  
`program_id()` (in module `nki.language`), 1236  
`program_ndim()` (in module `nki.language`), 1237  
`PSUM`, 1635  
`psum` (in module `nki.language`), 1235

## R

`rand()` (in module `nki.language`), 1194  
`random_seed()` (in module `nki.language`), 1195  
`range_select()` (in module `nki.isa`), 1270  
`read` (C++ function), 1521  
`read_stream_accessor` (C++ function), 1517  
`reciprocal()` (in module `nki.isa`), 1266  
`reciprocal()` (in module `nki.language`), 1219  
`reduce_cmd` (class in `nki.isa`), 1292  
`relu()` (in module `nki.language`), 1212  
`reshape()` (`nki.tensor` method), 1185  
`right_shift()` (in module `nki.language`), 1222  
`rms_norm()` (in module `nki.language`), 1217  
`RNE`, 1635  
`rsqrt()` (in module `nki.language`), 1211  
`RT`, 1635

## S

`SBUF`, 1635  
`sbuf` (in module `nki.language`), 1235  
Scalar Engine, 1633  
`scalar_tensor_tensor()` (in module `nki.isa`), 1257  
`ScalarE`, 1634  
`sequential_range()` (in module `nki.language`), 1233  
`set_accessor_coherence_policy` (C++ function), 1518  
`shape` (`nki.tensor` property), 1185  
`shared_constant()` (in module `nki.language`), 1195  
`shared_hbm` (in module `nki.language`), 1236  
`shared_identity_matrix()` (in module `nki.language`), 1195  
`sigmoid()` (in module `nki.language`), 1212  
`sign()` (in module `nki.language`), 1205  
`silu()` (in module `nki.language`), 1214  
`silu_dx()` (in module `nki.language`), 1214  
`simulate_kernel()` (in module `nki`), 1182  
`sin` (C++ function), 1515  
`sin()` (in module `nki.language`), 1209  
`sin_out` (C++ function), 1515  
`skip_middle_end_transformations()` (in module `nki.compiler`), 1300  
`softmax()` (in module `nki.language`), 1216  
`softplus()` (in module `nki.language`), 1215

`spmd_dim` (in module `nki.language`), 1237  
`sqrt()` (in module `nki.language`), 1211  
`square()` (in module `nki.language`), 1216  
`SR`, 1635  
State Buffer, 1633  
`static_range()` (in module `nki.language`), 1232  
`store()` (in module `nki.language`), 1187  
`sub` (C++ function), 1516  
`sub_out` (C++ function), 1516  
`subtract()` (in module `nki.language`), 1198  
`sum()` (in module `nki.language`), 1203  
Sync Engine, 1633

## T

`tan` (C++ function), 1515  
`tan()` (in module `nki.language`), 1210  
`tan_out` (C++ function), 1515  
`tanh()` (in module `nki.language`), 1210  
`tcu_accessor` (C++ function), 1518  
`tcu_to_tensor` (C++ function), 1523  
`tensor` (class in `nki`), 1183  
Tensor Engine, 1633  
`tensor_copy()` (in module `nki.isa`), 1261  
`tensor_copy_dynamic_dst()` (in module `nki.isa`), 1264  
`tensor_copy_dynamic_src()` (in module `nki.isa`), 1262  
`tensor_copy_predicated()` (in module `nki.isa`), 1264  
`tensor_partition_reduce()` (in module `nki.isa`), 1253  
`tensor_reduce()` (in module `nki.isa`), 1251  
`tensor_scalar()` (in module `nki.isa`), 1258  
`tensor_scalar_reduce()` (in module `nki.isa`), 1260  
`tensor_tensor()` (in module `nki.isa`), 1254  
`tensor_tensor_scan()` (in module `nki.isa`), 1255  
`tensor_to_tcu` (C++ function), 1523  
`TensorE`, 1634  
TF32, 1635  
`tf32` (in module `nki.language`), 1240  
`tile_size` (class in `nki.language`), 1241  
`torch.neuron.DataParallel()`  
    built-in function, 215  
`torch.neuron.DataParallel.disable_dynamic_batching()`  
    built-in function, 90, 216  
`torch::neuron::tcu_free` (C++ function), 1523  
`torch::neuron::tcu_malloc` (C++ function), 1523  
`torch_neuron.experimental.multicore_context()`  
    (in module `placement`), 221  
`torch_neuron.experimental.neuron_cores_context()`  
    (in module `placement`), 221  
`torch_neuron.experimental.set_multicore()` (in module `placement`), 223  
`torch_neuron.experimental.set_neuron_cores()`  
    (in module `placement`), 223

`torch_neuron.Optimization` (*built-in class*), 211

`torch_neuron.trace()`  
built-in function, 210

`torch_neuronx.analyze()`  
built-in function, 87

`torch_neuronx.bucket_model_trace()`  
built-in function, 77

`torch_neuronx.BucketModelConfig` (*built-in class*),  
76

`torch_neuronx.DataParallel()`  
built-in function, 90

`torch_neuronx.dynamic_batch()`  
built-in function, 78

`torch_neuronx.experimental.multicore_context()`  
built-in function, 86

`torch_neuronx.experimental.neuron_cores_context()`  
built-in function, 85

`torch_neuronx.experimental.profiler.profile()`  
built-in function, 347

`torch_neuronx.experimental.profiler.profile.start()`  
built-in function, 348

`torch_neuronx.experimental.set_multicore()`  
built-in function, 85

`torch_neuronx.experimental.set_neuron_cores()`  
built-in function, 84

`torch_neuronx.move_trace_to_device()`  
built-in function, 75

`torch_neuronx.PartitionerConfig()`  
built-in function, 79

`torch_neuronx.replace_weights()`  
built-in function, 81

`torch_neuronx.trace()`  
built-in function, 71

TP, 1635

TPr, 1635

Trainium/Inferentia2, 1632

Trainium2, 1632

`transpose()` (*in module nki.language*), 1218

Trn1, 1632

Trn2, 1632

`trunc()` (*in module nki.language*), 1206

## V

`var()` (*in module nki.language*), 1203

Vector Engine, 1633

VectorE, 1634

`view()` (*nki.tensor method*), 1185

## W

`where()` (*in module nki.language*), 1229

`write` (C++ function), 1521

`write_stream_accessor` (C++ function), 1517

## Z

`zeros` (C++ function), 1514

`zeros()` (*in module nki.language*), 1193

`zeros_like()` (*in module nki.language*), 1193