# Rapid Innovation on AWS Trainium with Triton-Inspired Programming

Jonathan M. Henson    Principal Software Engineer    Amazon Web Services
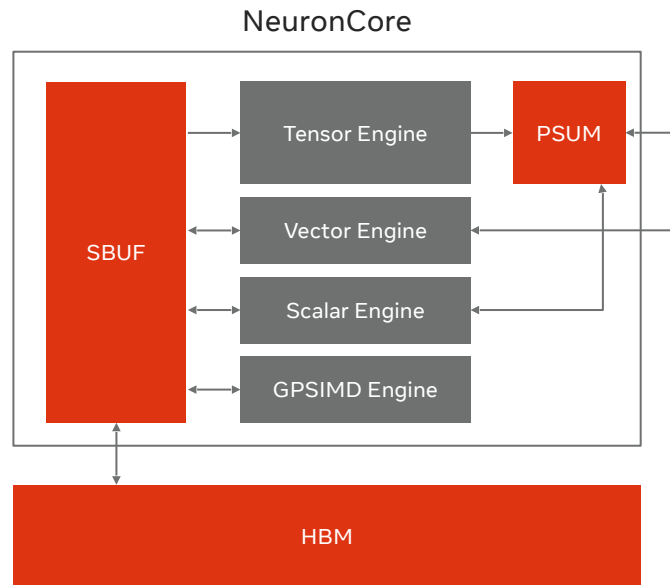
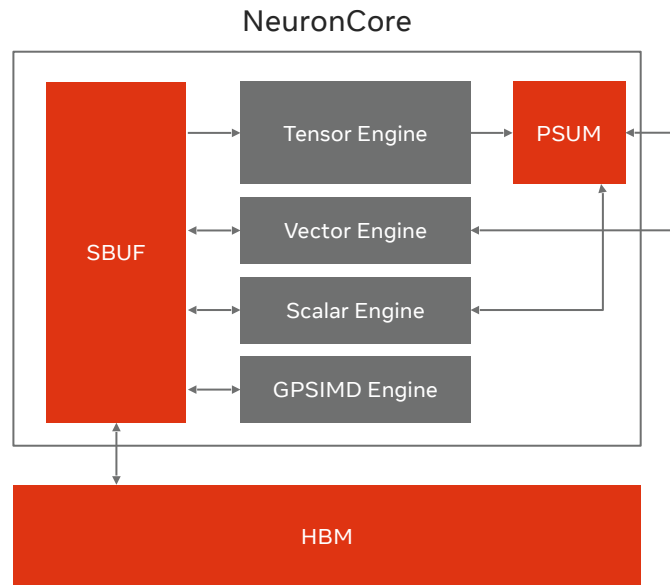# Overview of the Neuron Device Architecture



- Powered by NeuronCores on the Amazon EC2 Trainium and Inferentia instance-type families.
- 2 NeuronCores (v2).
- 2 HBM stacks with a total 32GiB at 820 GB/s.
- 32 DMA engines to move data within and across devices.
- 6 CC-Cores
- 2 (Inferentia2) or 4 (Trainium) NeuronLink-v2 for device-to-device collective communication.

# Compute Engines



- **Memory**: two software-managed on-chip SRAMs
- **Tensor**: matrix-multiplication, tensor-tensor ops, transpose
- **Vector**: accelerates vector operations
- **Scalar**: includes hardware acceleration for non-linear activation functions
- **GPSIMD**: general purpose compute with custom-ops

# Compute Engines



- **Memory**: two software-managed on-chip SRAMs
- **Tensor**: matrix-multiplication, tensor-tensor ops, transpose
- **Vector**: accelerates vector operations
- **Scalar**: includes hardware acceleration for non-linear activation functions
- **GPSIMD**: general purpose compute with custom-ops

# Triton Leads to the Neuron Kernel Interface (NKI)

```python
import triton
import triton.language as tl


@triton.jit
def add_kernel(x_ptr,  # *Pointer* to first input vector.
               y_ptr,  # *Pointer* to second input vector.
               output_ptr,  # *Pointer* to output vector.
               n_elements,  # Size of the vector.
               BLOCK_SIZE: tl.constexpr,  # Number of elements each program should process.
               # NOTE: `constexpr` so it can be used as a shape value.
               ):
    # There are multiple 'programs' processing different data. We identify which program
    # we are here:
    pid = tl.program_id(axis=0)  # We use a 1D launch grid so axis is 0.
    # This program will process inputs that are offset from the initial data.
    # For instance, if you had a vector of length 256 and block_size of 64, the programs
    # would each access the elements [0:64, 64:128, 128:192, 192:256].
    # Note that offsets is a list of pointers:
    block_start = pid * BLOCK_SIZE
    offsets = block_start + tl.arange(0, BLOCK_SIZE)
    # Create a mask to guard memory operations against out-of-bounds accesses.
    mask = offsets < n_elements
    # Load x and y from DRAM, masking out any extra elements in case the input is not a
    # multiple of the block size.
    x = tl.load(x_ptr + offsets, mask=mask)
    y = tl.load(y_ptr + offsets, mask=mask)
    output = x + y
    # Write x + y back to DRAM.
    tl.store(output_ptr + offsets, output, mask=mask)
```

```python
import neuronxcc.nki as nki
import neuronxcc.nki.language as nl

def nki_tensor_add_kernel_(a_input, b_input, c_output):
    # Calculate tile offsets based on current 'program'
    offset_i_x = nl.program_id(0) * 128
    offset_i_y = nl.program_id(1) * 512

    # Generate tensor indices to index tensors a and b
    ix = offset_i_x + nl.arange(128)[:, None]
    iy = offset_i_y + nl.arange(512)[None, :]

    # Load input data from device memory (HBM) to on-chip memory (SBUF)
    # We refer to an indexed portion of a tensor as an intermediate tensor
    a_tile = nl.load(a_input[ix, iy])
    b_tile = nl.load(b_input[ix, iy])

    # compute a + b
    c_tile = a_tile + b_tile

    # store the addition results back to device memory (c_output)
    nl.store(c_output[ix, iy], value=c_tile)
```

- Tile semantics natural fit for NeuronCore memory model
- Adapted to target baremetal Neuron ISA access.

# Triton Leads to the Neuron Kernel Interface (NKI)

```python
import triton
import triton.language as tl


@triton.jit
def add_kernel(x_ptr,  # *Pointer* to first input vector.
               y_ptr,  # *Pointer* to second input vector.
               output_ptr,  # *Pointer* to output vector.
               n_elements,  # Size of the vector.
               BLOCK_SIZE: tl.constexpr,  # Number of elements each program should process.
               # NOTE: `constexpr` so it can be used as a shape value.
               ):
    # There are multiple 'programs' processing different data. We identify which program
    # we are here:
    pid = tl.program_id(axis=0)  # We use a 1D launch grid so axis is 0.
    # This program will process inputs that are offset from the initial data.
    # For instance, if you had a vector of length 256 and block_size of 64, the programs
    # would each access the elements [0:64, 64:128, 128:192, 192:256].
    # Note that offsets is a list of pointers:
    block_start = pid * BLOCK_SIZE
    offsets = block_start + tl.arange(0, BLOCK_SIZE)
    # Create a mask to guard memory operations against out-of-bounds accesses.
    mask = offsets < n_elements
    # Load x and y from DRAM, masking out any extra elements in case the input is not a
    # multiple of the block size.
    x = tl.load(x_ptr + offsets, mask=mask)
    y = tl.load(y_ptr + offsets, mask=mask)
    output = x + y
    # Write x + y back to DRAM.
    tl.store(output_ptr + offsets, output, mask=mask)
```

```python
import neuronxcc.nki as nki
import neuronxcc.nki.language as nl

def nki_tensor_add_kernel_(a_input, b_input, c_output):
    # Calculate tile offsets based on current 'program'
    offset_i_x = nl.program_id(0) * 128
    offset_i_y = nl.program_id(1) * 512

    # Generate tensor indices to index tensors a and b
    ix = offset_i_x + nl.arange(128)[:, None]
    iy = offset_i_y + nl.arange(512)[None, :]

    # Load input data from device memory (HBM) to on-chip memory (SBUF)
    # We refer to an indexed portion of a tensor as an intermediate tensor
    a_tile = nl.load(a_input[ix, iy])
    b_tile = nl.load(b_input[ix, iy])

    # compute a + b
    c_tile = a_tile + b_tile

    # store the addition results back to device memory (c_output)
    nl.store(c_output[ix, iy], value=c_tile)
```
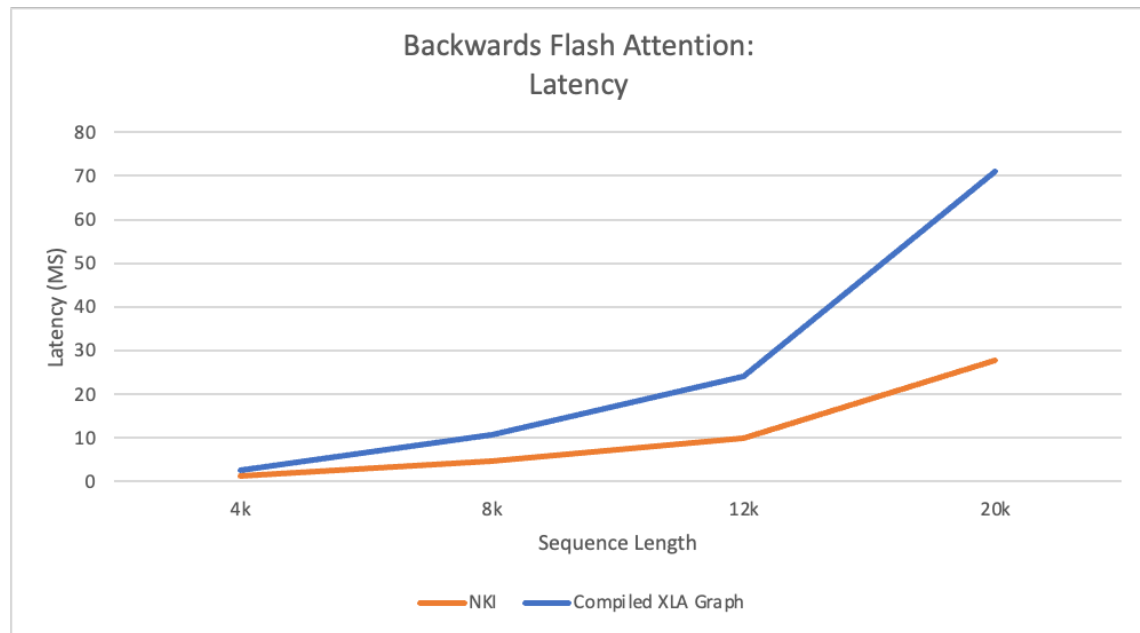
- Tile semantics natural fit for NeuronCore memory model
- Adapted to target baremetal Neuron ISA access.

# NKI Already Enabling Community Delivery of Novel Model Architectures



Backwards Flash Attention: Latency

- Flash Attention, 2.5x+ improvement!
- Mamba

# Available in Latest Neuron SDK Release

- Integrated with Jax and PyTorch
- Launch and Run on Amazon EC2 in seconds.
- Open Sourced Kernels on GitHub
- Documentation on GitHub and AWS Docs

github.com/aws-neuron/nki-samples.git

https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/nki

Thank you!