intel

# TECH tour.TW

# Lunar Lake AI Hardware Accelerators

**TAP**
Intel Fellow

# Unmatched AI Compute

## With our Multi-Engine approach

Up to

# 120

platform TOPS

**GPU**
Creator & gamer AI

**NPU**
AI assistants & gen AI

**CPU**
Light "embedded" AI

intel. TECH. tour.TW

Lunar Lake

# CPU
## AI Engine

| **P-core & E-core** | **VNNI & AVX** | **5** |
|---|---|---|
| CPU architecture | AI instructions | peak TOPs |

All tops are Int8 on high end SKU, will vary based on SKU

intel

TECH
tour.TW

# Lunar Lake
# NPU Deepdive

**Darren Crews**
Sr Principal Engineer, NPU Lead Architect

Lunar Lake

# NPU
## AI Engine

| | | |
|---|---|---|
| **NPU 4** | **2x** | **48** |
| Architecture | Power efficiency | Peak TOPs |

All tops are Int8 on high end SKU, will vary based on SKU

# Continuous NPU Improvements

Across 4 generations of IP

**0.5 pTOPs**

**NPU 1**
2018

**7 pTOPs**

**NPU 2**
2021

**11.5 pTOPs**

**NPU 3**
2023

**48 pTOPs**

**NPU 4**
2024

All tops are Int8 on high end SKU, will vary based on SKU

intel. TECH tour.TW

**NPU 4**
2024

**Proven foundations**
based on three prior generations

**Higher compute capacity**
to support growing number of use cases

**Increased efficiency**
to support longer battery life

intel. TECH tour.TW

# What is a TOP?

| **Trillions** | of | **Operations** | per | **Second** |

**Operations**
|
1**M**ultiply &
1**A**ccumulate
(MAC)

**Second**
|
Clock
Frequency
(Hz)

intel. TECH tour.TW

# How Many AI TOPS?

**Operations**   per   **Second**   ÷   **Trillions**

1 Multiply & 1 Accumulate (MAC)

Clock frequency (Hz)

$10^{12}$

intel. TECH. tour.TW

# How Many AI TOPS?

**Operations**  ✕  **Frequency**  ÷  **Trillions**

1 Multiply &
1 Accumulate
(MAC)

$10^{12}$

intel. TECH. tour.TW

# Operation Types
## Overview

| | Scalar | Vector | Matrix |
|---|---|---|---|
| Complexity | 1 | N | $N^2$ |
| Example functions | Conditional, Looping | SoftMax, Activation functions | Convolution, Matrix multiplication |
| Occurrence in AI | Low | Very high | Very high |

# Operation Types
## Overview

| | Scalar | Vector | Matrix |
|---|---|---|---|
| **Complexity** | 1 | N | $N^2$ |
| **Example functions** | Conditional / Looping | SoftMax / Activation functions | Convolution / Matrix multiplication |
| **Occurrence in AI** | Low | Very high | Very high |

**TOPs**

# Scaling the NPU

Increase number of engines

Increase frequency

Improve architecture

NPU 3

NPU 4

NPU 3 · 4K MACs · 2 NCEs → NPU 4 · 12K MACs · 6 NCEs

**NPU 3**

Global Control

| MMU | DMA |

Scratchpad RAM · L2$

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
Mac Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
Mac Array
Activation Function · Load/Store
Data Conversion

**NPU 4**

Global Control

| MMU | DMA |

Scratchpad RAM · L2$

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

Neural Compute Engine
SHAVE DSP · SHAVE DSP
Inference Pipeline
MAC Array
Activation Function · Load/Store
Data Conversion

intel. TECH tour.TW

# Scaling the NPU

Increase number of engines

Increase frequency

Improve architecture

NPU 3

NPU 4

# Increased Efficiency & Increased Performance

- Increased clock
- New node
- Architecture improvements



2x perf at ISO power[1]

4x peak performance

NPU 3

NPU 4

Performance

Power

intel. TECH tour.TW

# Scaling the NPU

Increase number of engines

Increase frequency

Improve architecture

NPU 3

NPU 4

intel® TECH tour.TW

# NPU 4
## Architecture overview

# NPU 4
## Neural compute engine

**Specialized engines**
Matrix + Vector

**Inference pipeline**
MAC arrays + fixed function

**Programmable DSPs**

Global Control

| MMU | DMA |

Scratchpad RAM

L2$

| Neural Compute Engine | Neural Compute Engine | Neural Compute Engine | Neural Compute Engine | Neural Compute Engine | Neural Compute Engine |
|---|---|---|---|---|---|
| SHAVE DSP / SHAVE DSP | SHAVE DSP / SHAVE DSP | SHAVE DSP / SHAVE DSP | SHAVE DSP / SHAVE DSP | SHAVE DSP / SHAVE DSP | SHAVE DSP / SHAVE DSP |
| Inference Pipeline | Inference Pipeline | Inference Pipeline | Inference Pipeline | Inference Pipeline | Inference Pipeline |
| MAC Array | MAC Array | MAC Array | MAC Array | MAC Array | MAC Array |
| Activation Function / Load/Store / Data Conversion | Activation Function / Load/Store / Data Conversion | Activation Function / Load/Store / Data Conversion | Activation Function / Load/Store / Data Conversion | Activation Function / Load/Store / Data Conversion | Activation Function / Load/Store / Data Conversion |

intel. TECH tour.TW

# NPU 4
## MAC array

**Matrix multiplication & convolution**

**2048** MAC/cycle int8
**1024** MAC/cycle FP16

**Up to 2x[1] efficiency**
driving better perf/watt

**INT8**
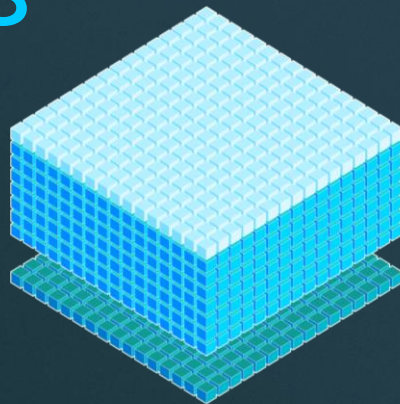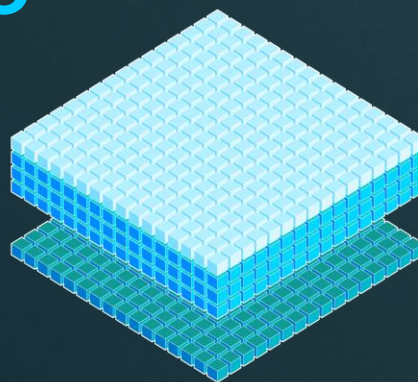16x16x8

**FP16**
16x16x4

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DS

Inference Pipeline

MAC Array

Activation Function

Load Stor

Data Conversion

[1]See details in backup

# NPU 4
## Data conversion

**Datatype conversion**

**Fused operations**

**Output data re-layout**

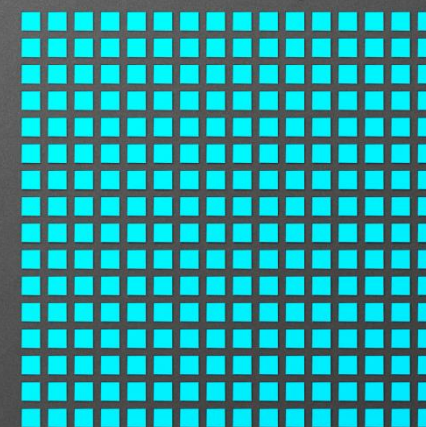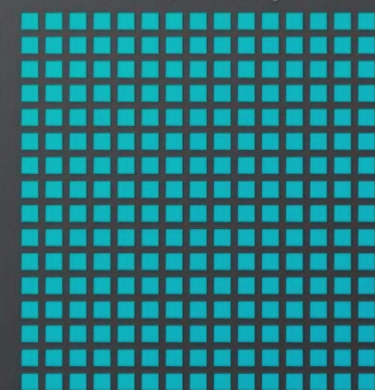min $(X_f)$      0      max $(X_f)$

0      255

Neural Compute Engine

SHAVE DSP    SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/ Store

Data Conversion

Neural Compute Engine

SHAVE DSP    SHAVE DS

Inference Pipeline

MAC Array

Activation Function

Loa Sto

Data Conversion

# NPU 4
## SHAVE DSP

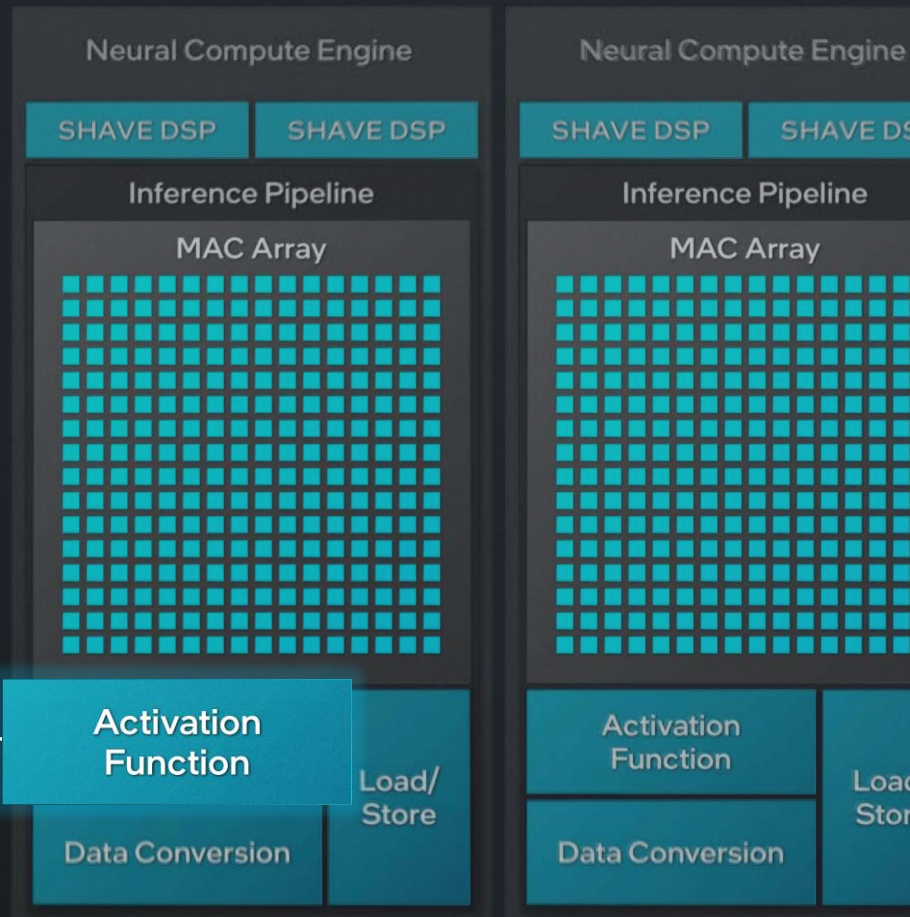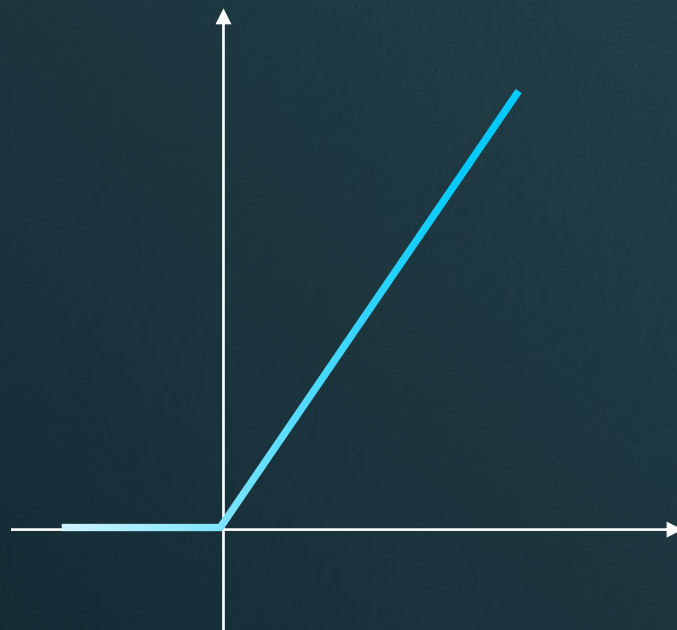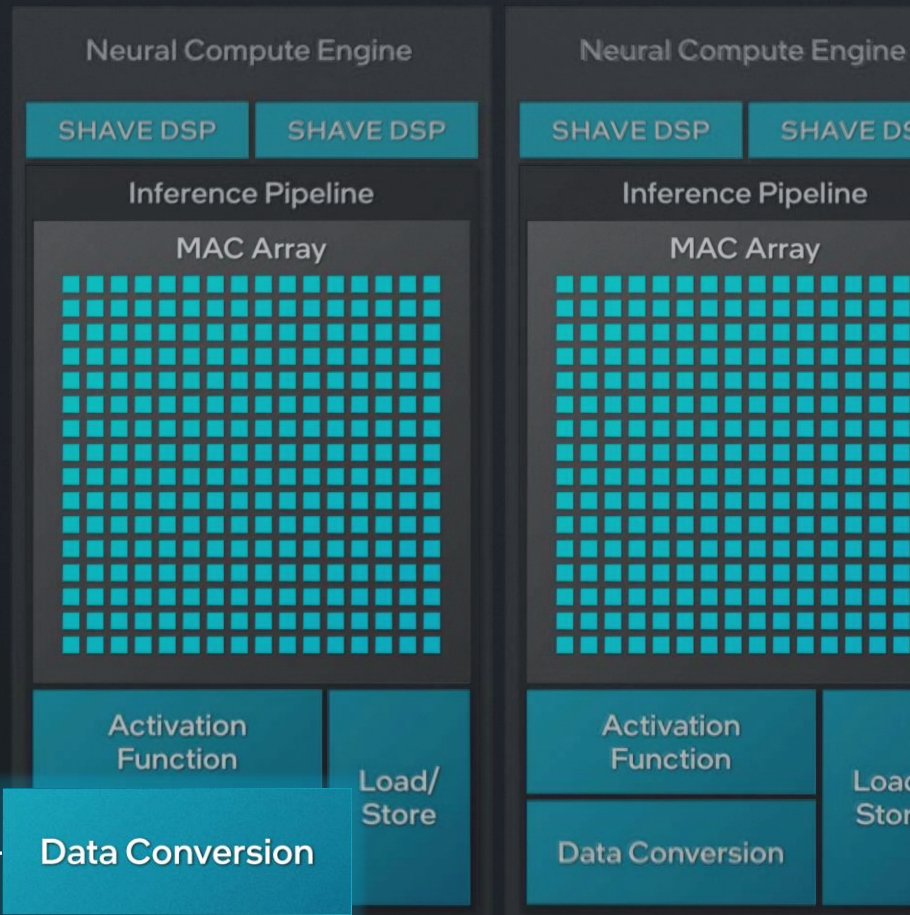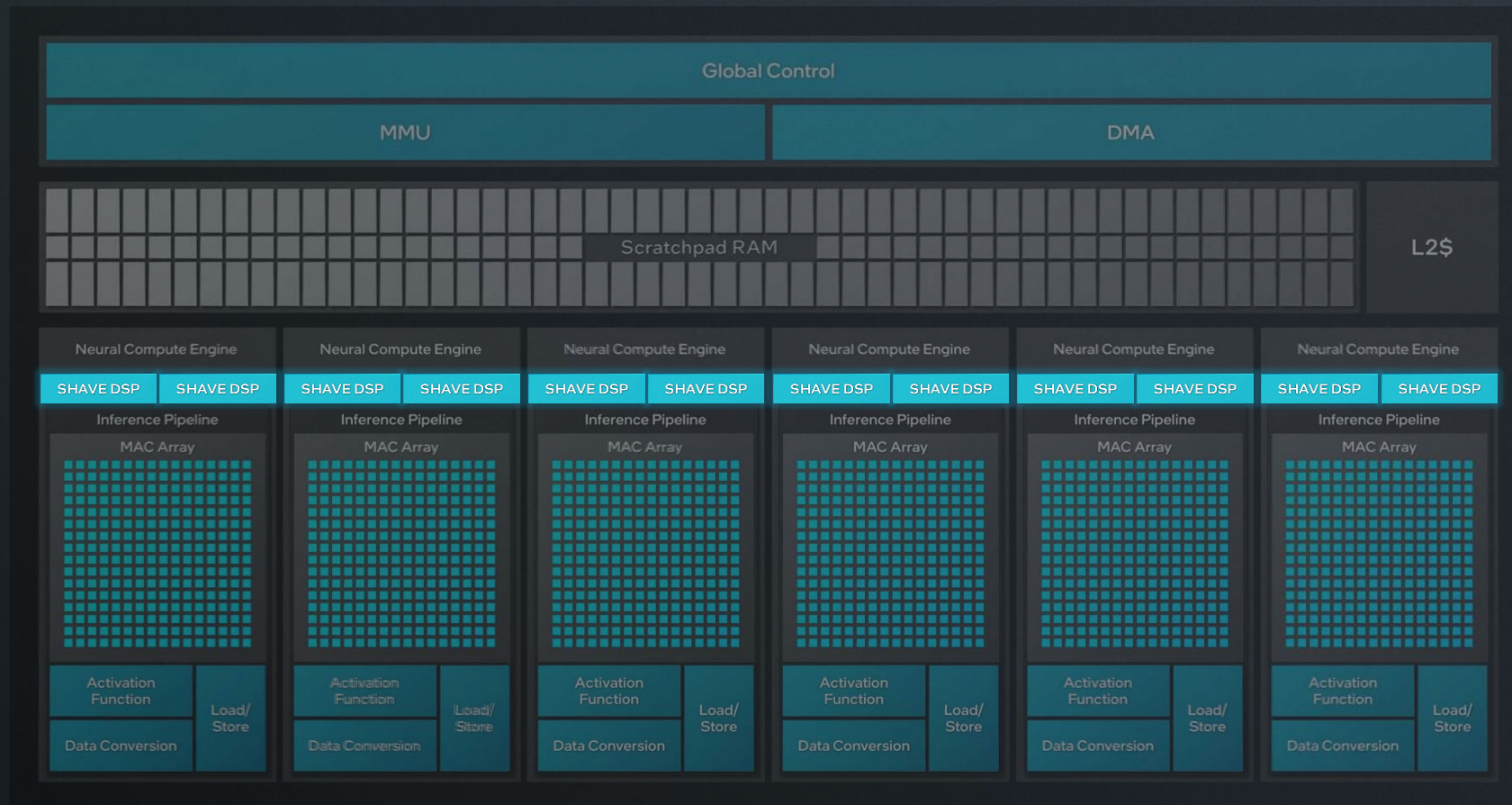**Upgraded SHAVE DSP**
4x vector compute

**12x overall vector perf**
improves transformer /LLM performance

Global Control

MMU

DMA

Scratchpad RAM

L2$

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Data Conversion

Load/ Store

See details in backup

# SHAVE DSP
## Vector increase

**NPU3 SHAVE DSP** · FP16

128bit

| FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 |

×

| FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 |

**8** FP16 vector ops/clock per **DSP**

VRF → VAU → VRF

---

**NPU4 SHAVE DSP** · FP16

512bit

| FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 |

×

| FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 | FP16 |

**32** FP16 vector ops/clock per **DSP**

VRF → VAU → VRF
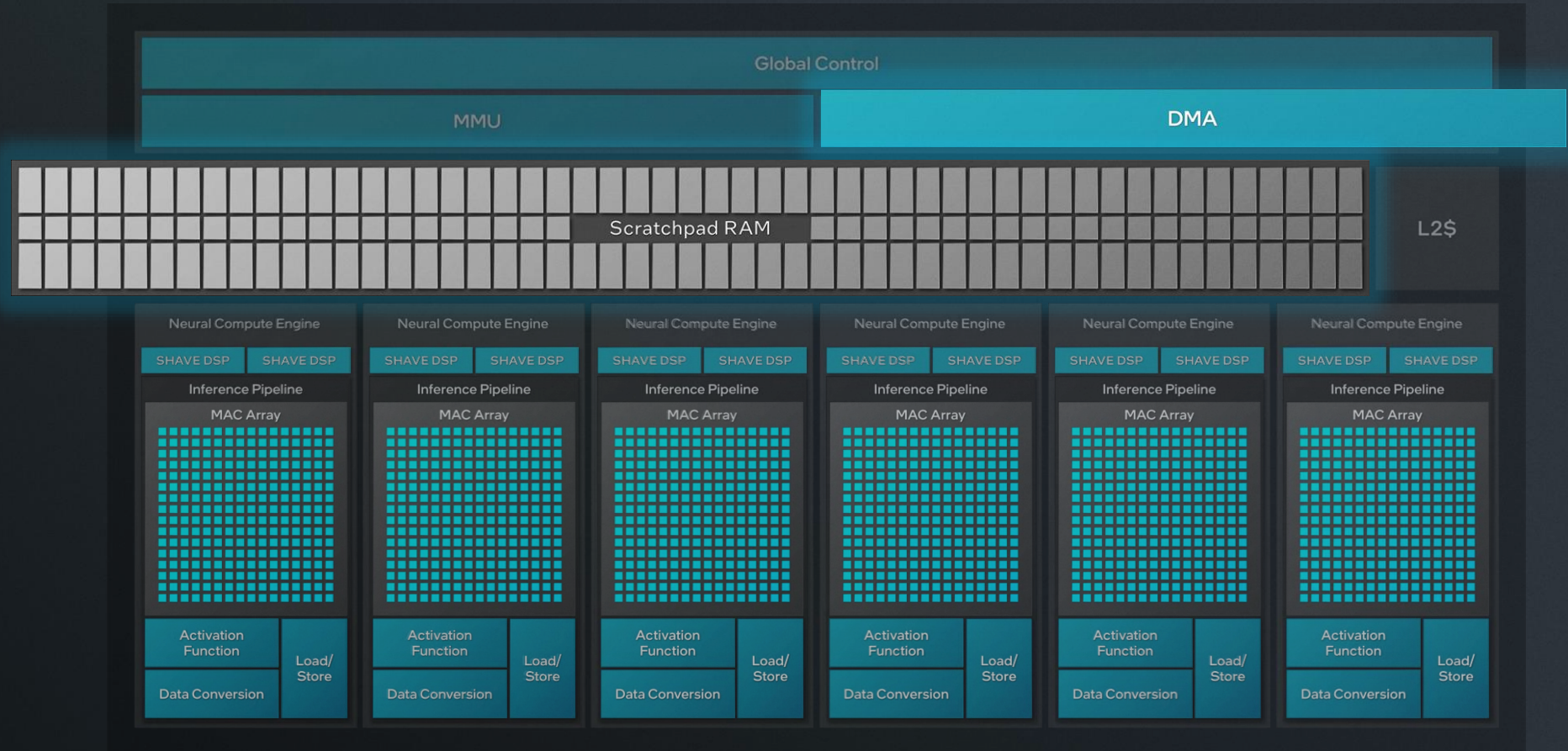
# NPU 4
## DMA engine

**2x DMA bandwidth**
improves network performance especially LLMs

**New functions**
Embedding tokenization

Global Control

MMU | DMA

Scratchpad RAM

L2$

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

Neural Compute Engine

SHAVE DSP | SHAVE DSP

Inference Pipeline

MAC Array

Activation Function

Load/Store

Data Conversion

intel. TECH tour.TW

See backup for details.
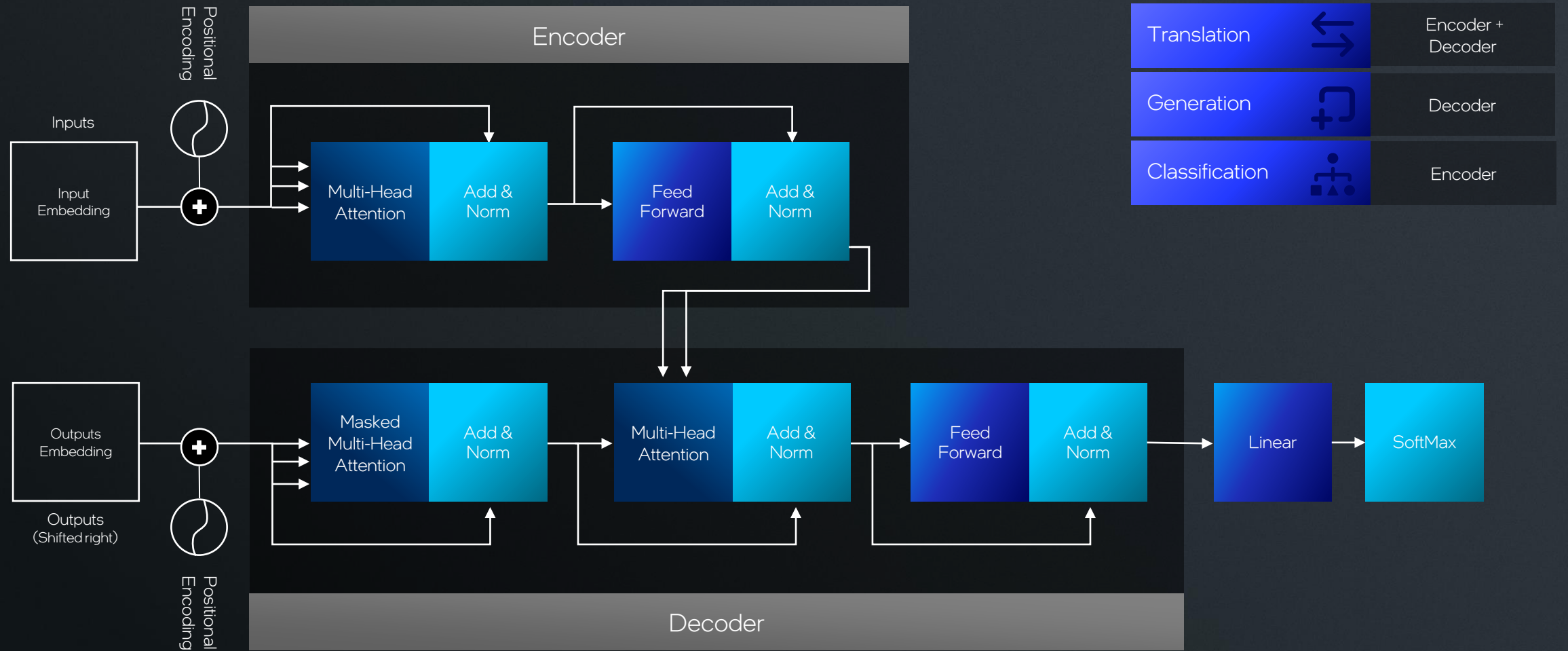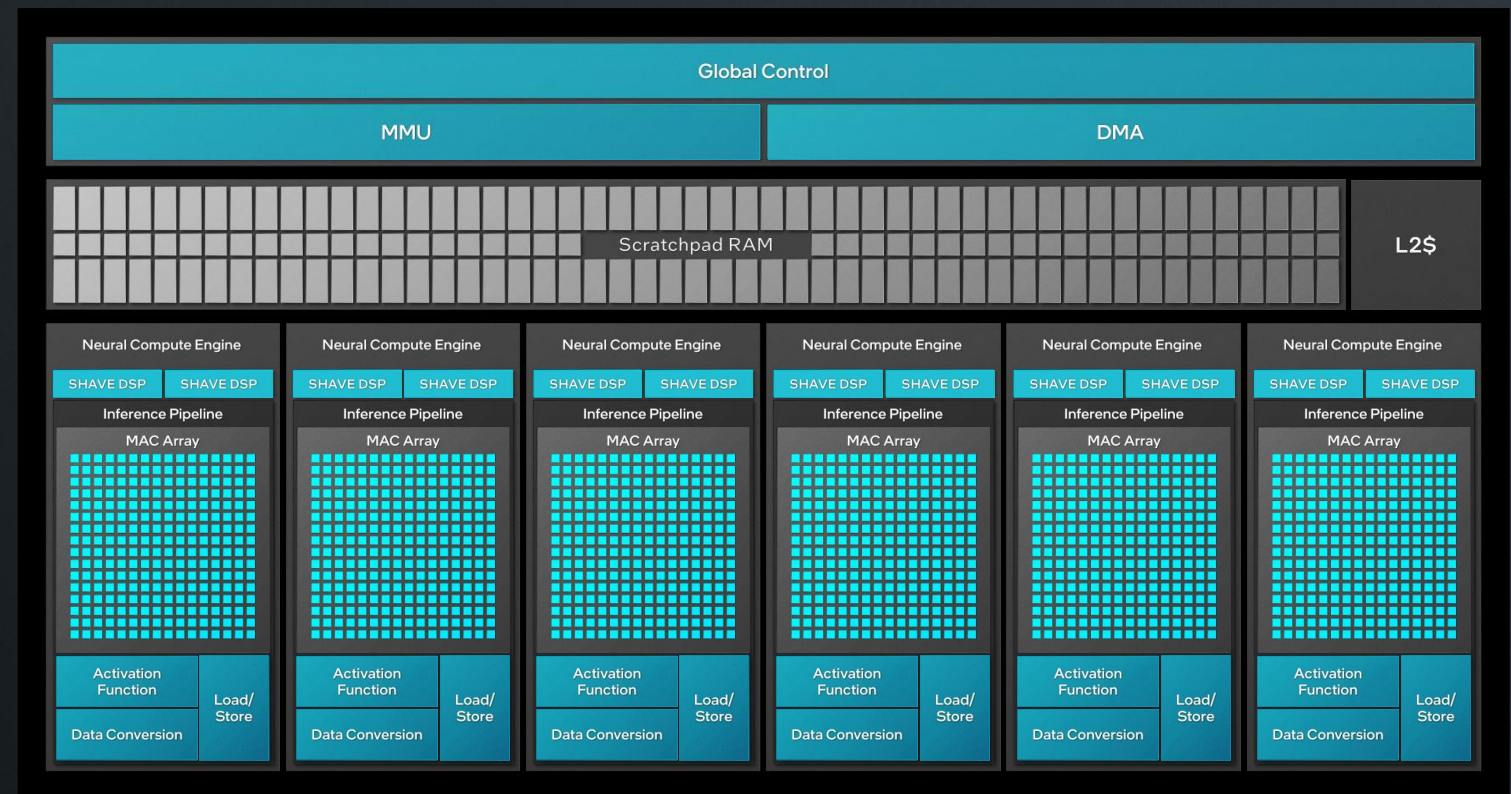
# Transformer Model Architecture
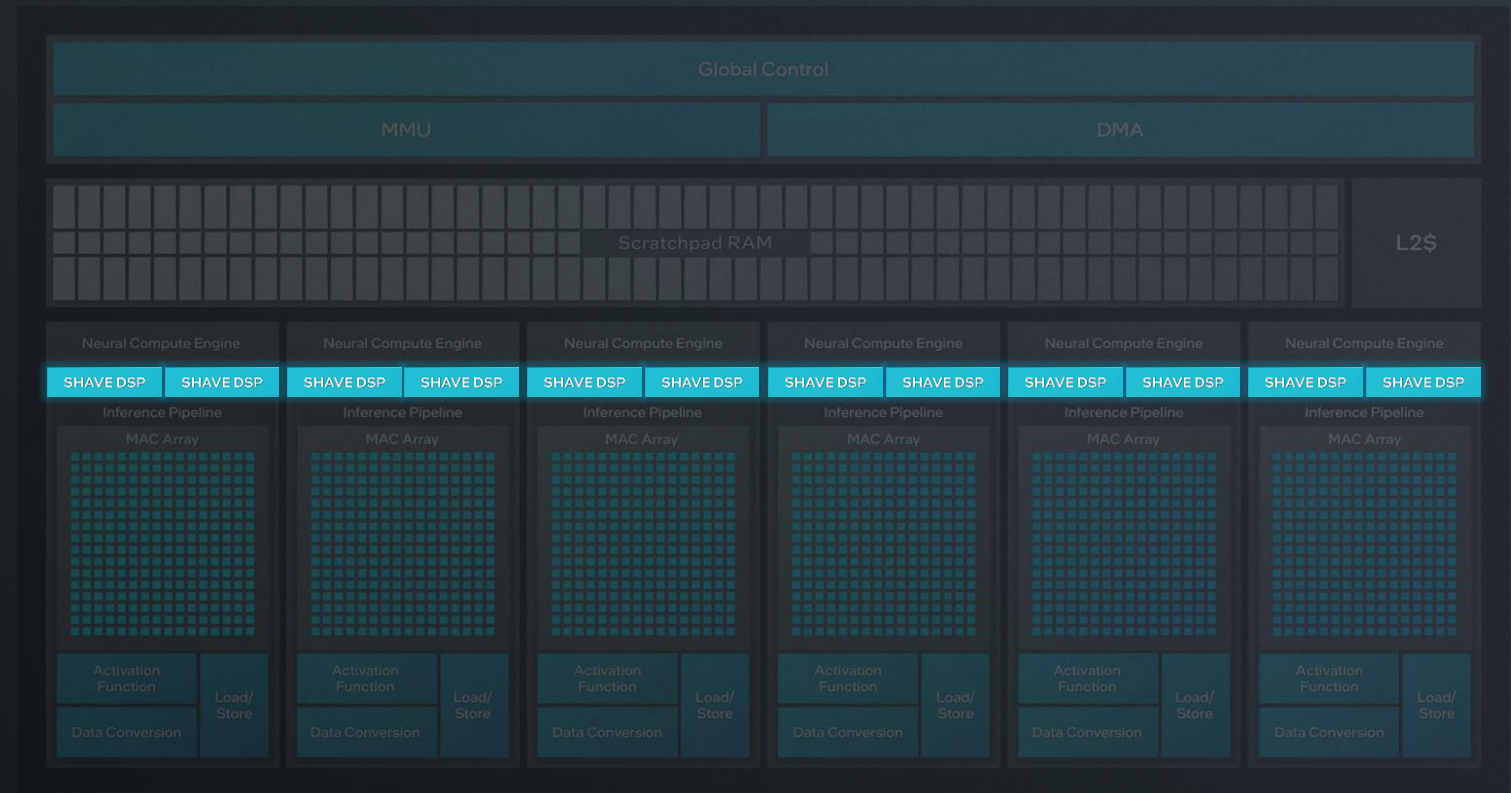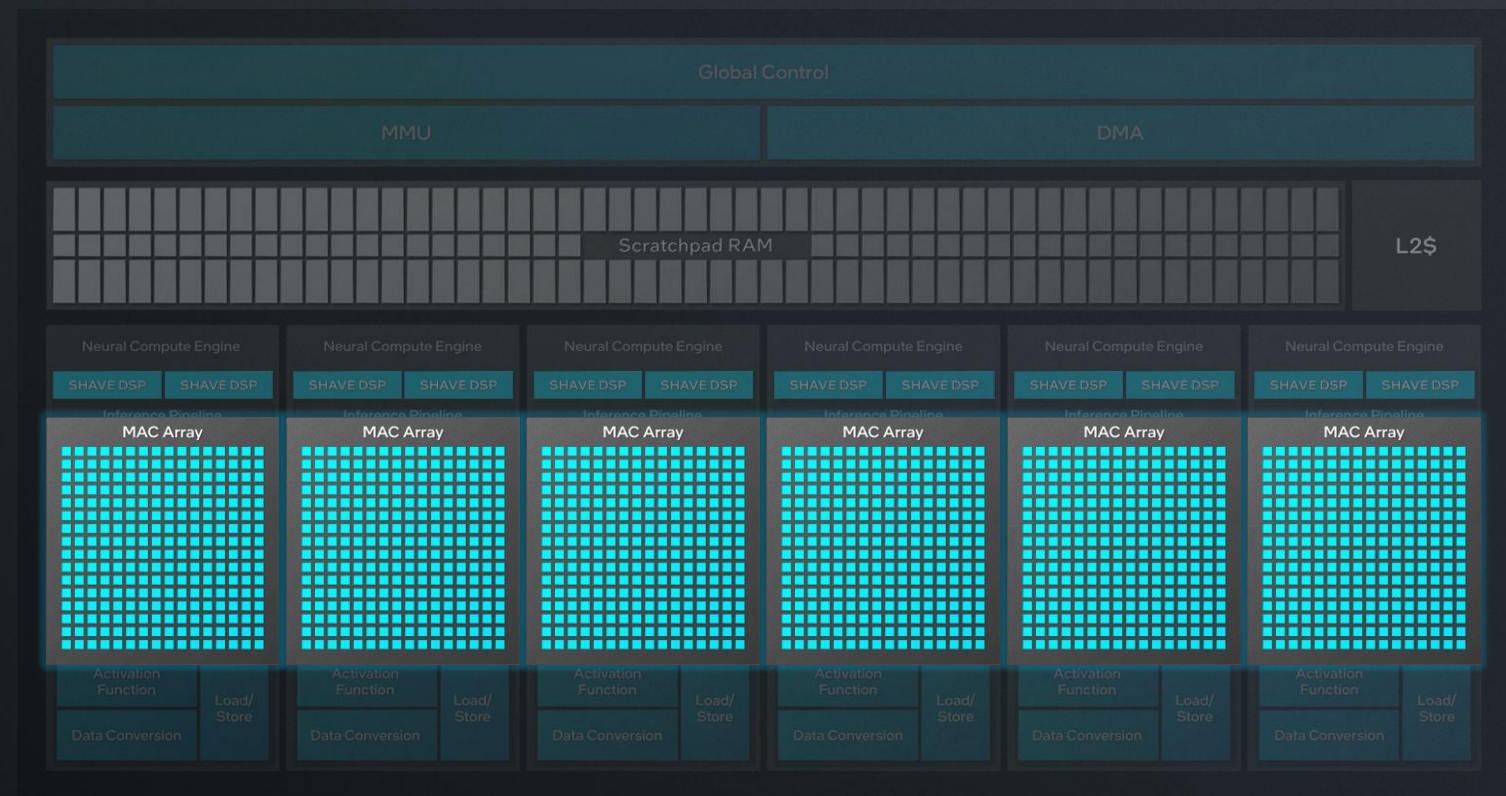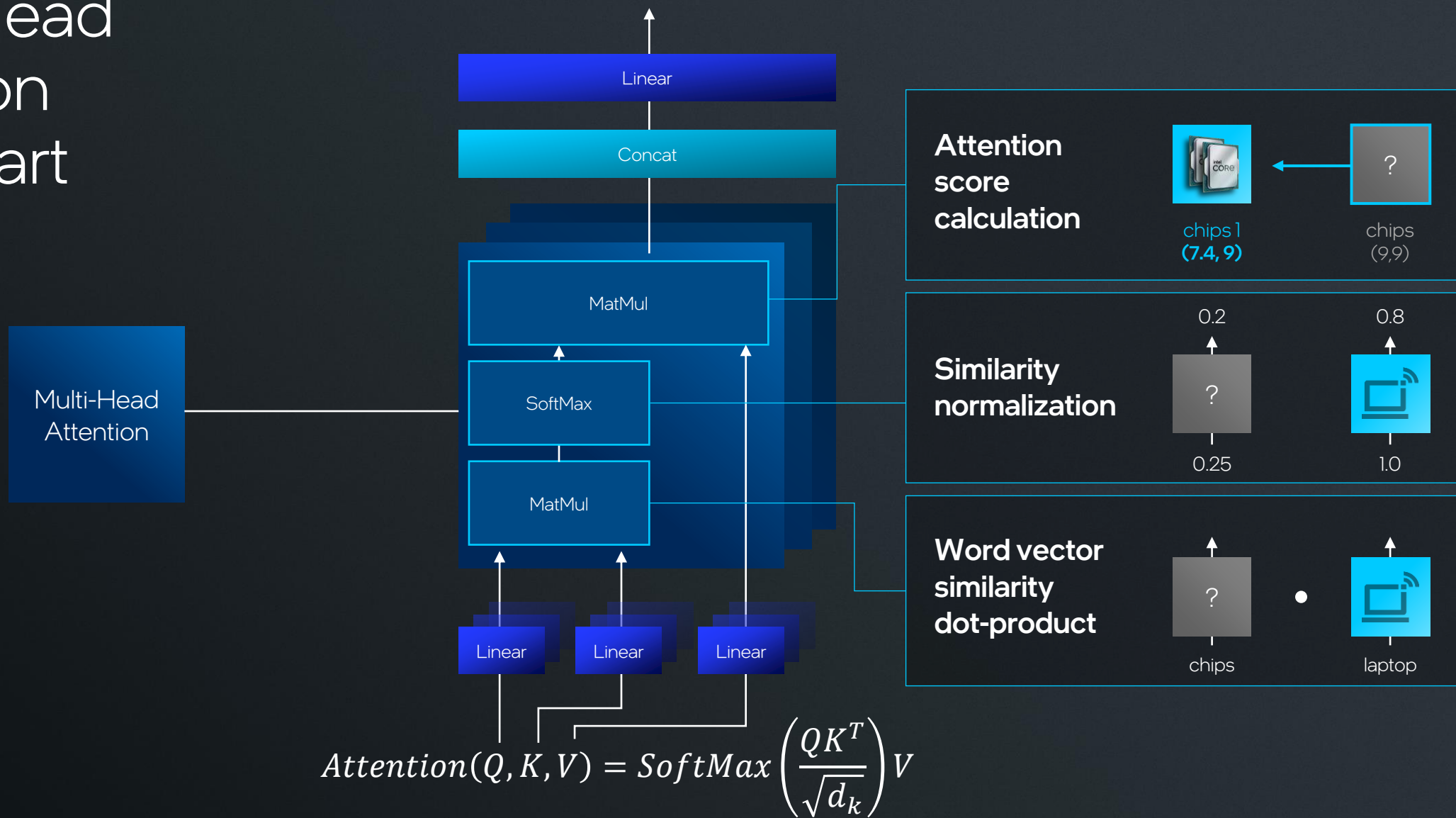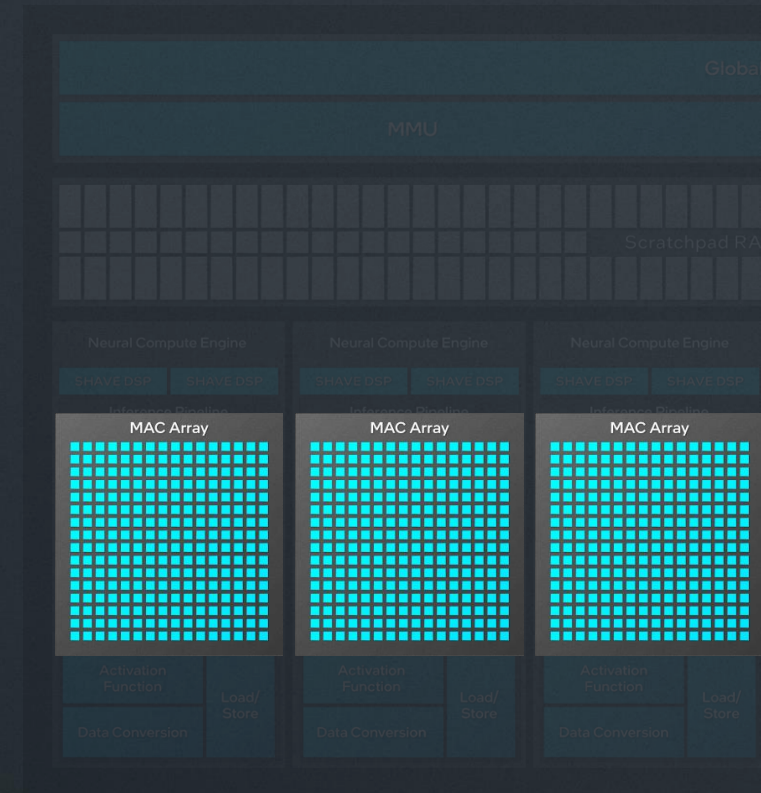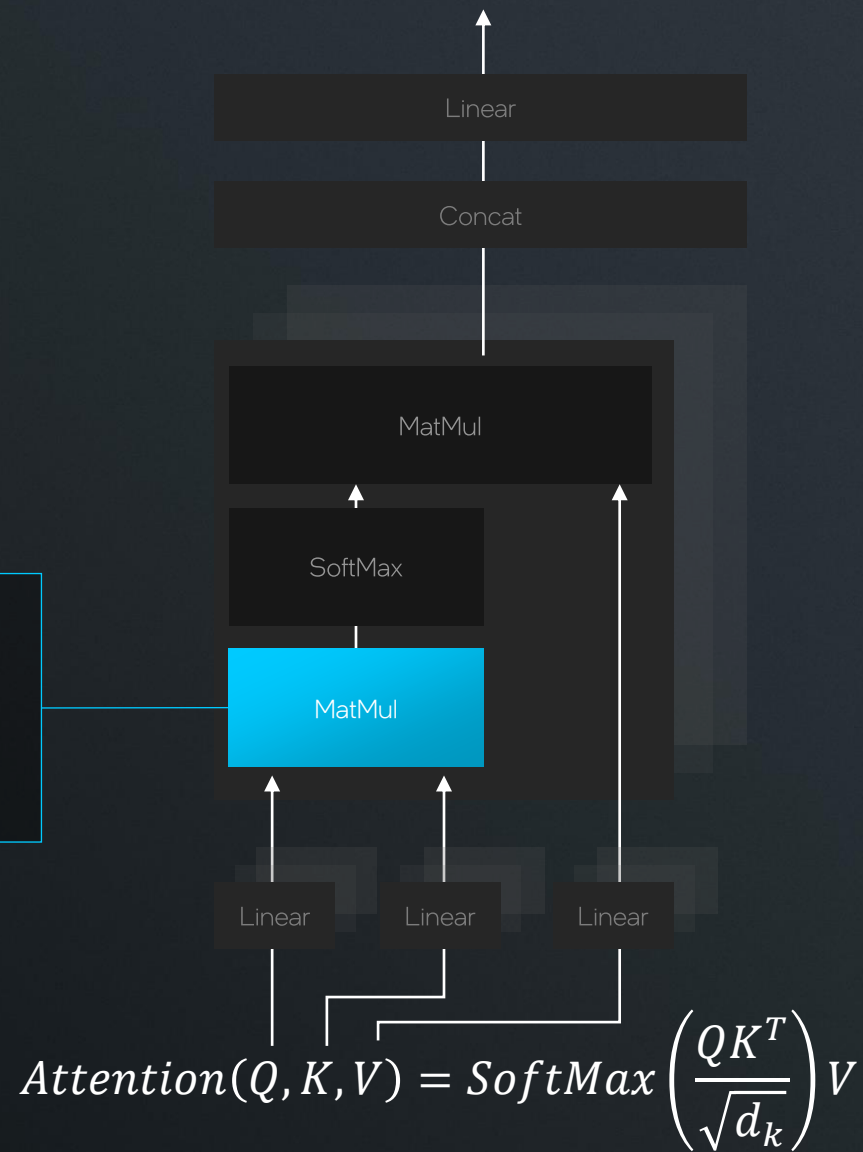
# Transformer Architecture on Intel's NPU
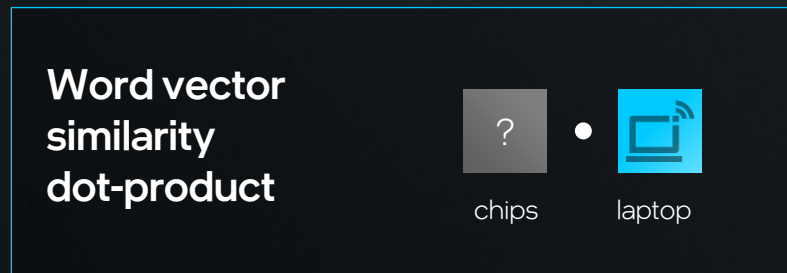
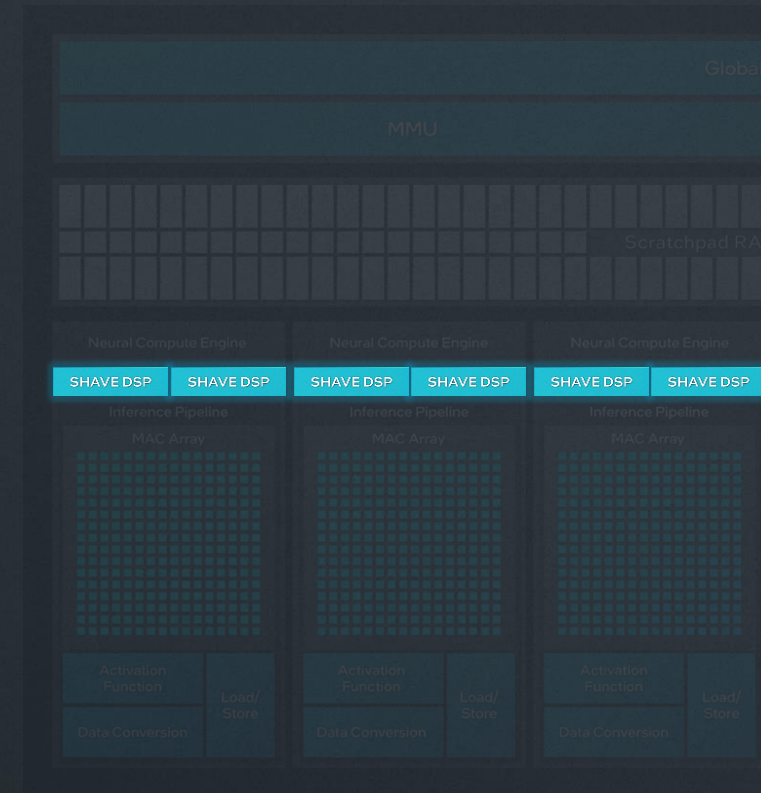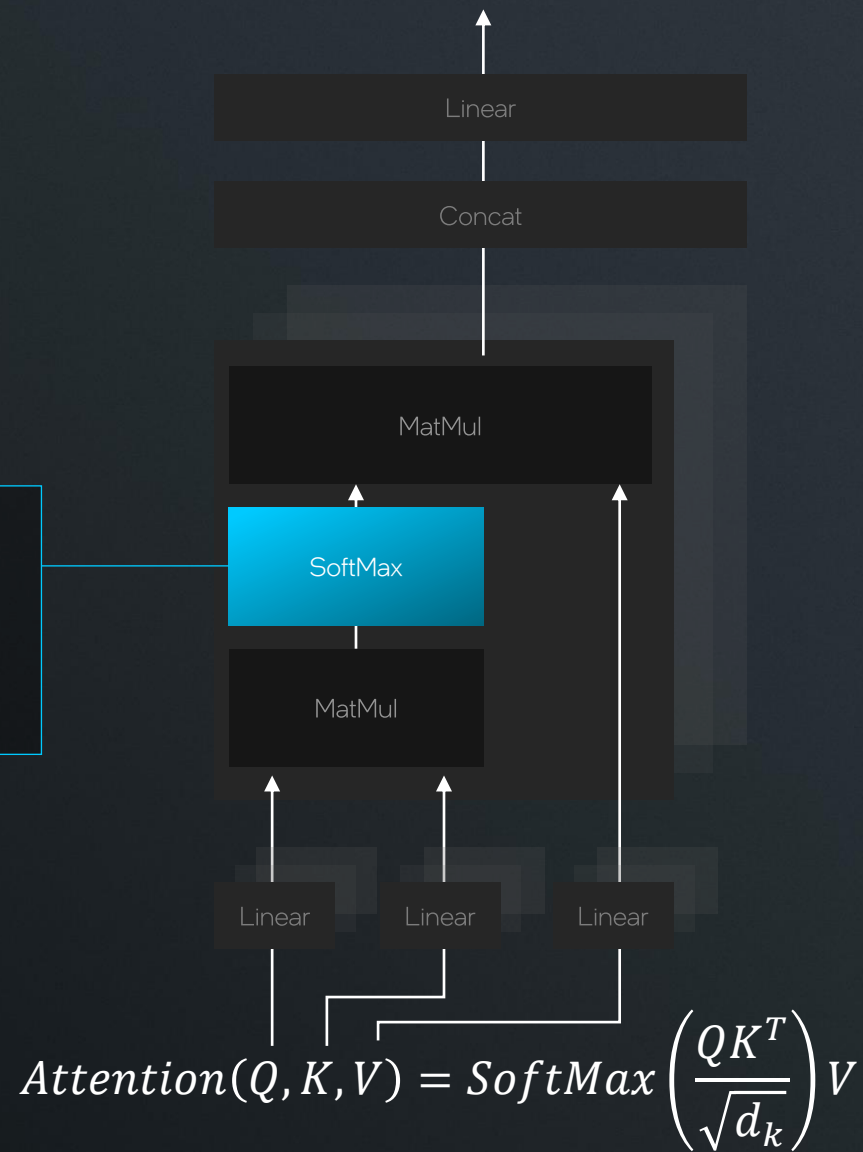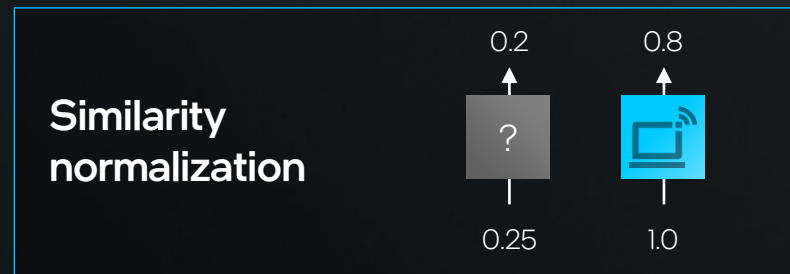Transformer Architecture on Intel's NPU

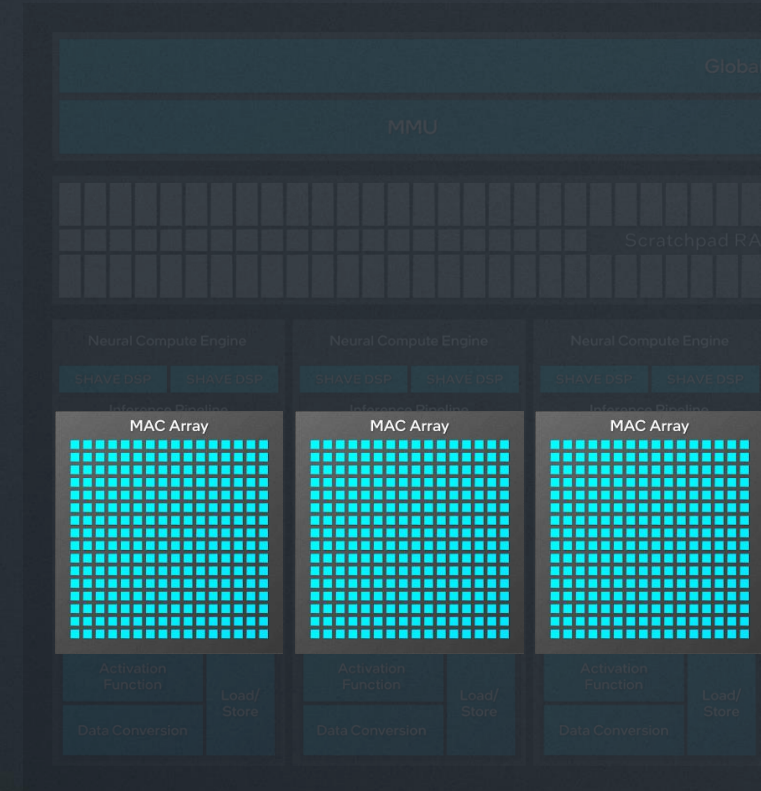# Transformer Architecture on Intel's NPU

# Multi-Head Attention Flowchart

Linear

Concat

MatMul

SoftMax

MatMul

Multi-Head Attention

Linear

Linear

Linear

$$Attention(Q, K, V) = SoftMax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

**Attention score calculation**

chips 1
**(7.4, 9)**

chips
(9,9)

**Similarity normalization**

0.2

0.8

?

0.25

1.0

**Word vector similarity dot-product**

?

•

chips

laptop

intel. TECH. tour.TW

# Multi- Head Attention Flowchart



Linear

Concat

MatMul

SoftMax

MatMul

Similarity normalization

0.2     0.8

?

0.25     1.0

Linear     Linear     Linear

$$Attention(Q, K, V) = SoftMax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

intel. TECH tour.TW

# Stable Diffusion Architecture



"Cute kitten with a pink bow" → Text prompt understanding (Text encoder) → U-Net diffusion (U-Net+ | U-Net-) → Image decoder (VAE) → Output

# Accelerating Multi-Head Attention

## Performance on U-Net



up to **9x faster** attention calculation

0ms 50ms 100ms 150ms 200ms 250ms 300ms 350ms

MAC ARRAY
SHAVE DSP
SHAVE DSP

intel. TECH tour.TW

*See details in backup

# Accelerating Multi-Head Attention

Performance on U-Net

Stable Diffusion
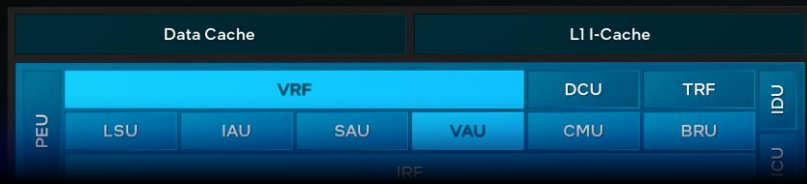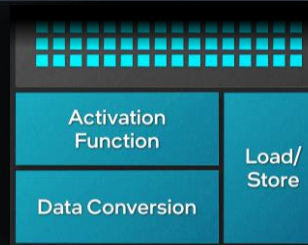
Demo

# Next Gen NPU 4

Largest integrated and dedicated AI accelerator for the AI PC

**12** Enhanced SHAVE DSPs

Accelerating LLM & transformer operations



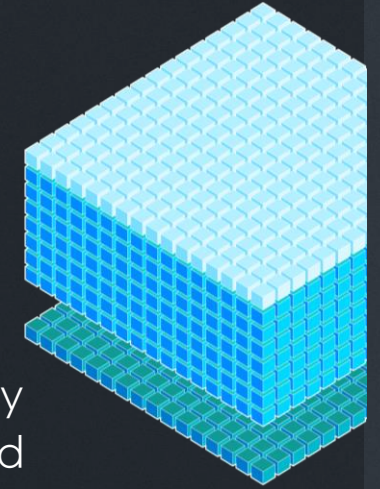Native activation function & data conversion support
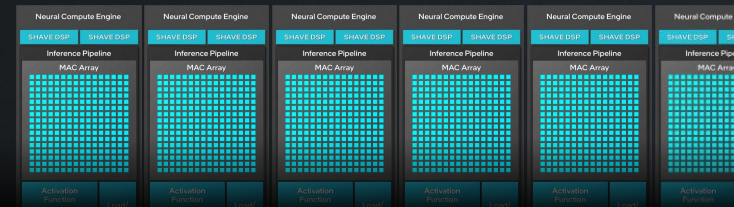


Up to **48** TOPS

**2x** Bandwidth

Efficiency optimized MAC array

DMA

Embedding tokenization used for LLMs

**6** Neural compute engines



intel

# Notices & Disclaimers

intel® TECH tour.TW

# APPENDIX

| Claim # & Statement | Slide # & Title/Details |
|---|---|
| | SLIDE 22: Increased Efficiency & Increased Performance |
| 2x performance at ISO power vs. Meteor Lake | Testing by Intel as of January 2024. Based on VPU-EM simulation. Power data is generated from the simulation tool based on power data that has been extracted from circuit simulation tools. This simulation, which is a ~100% utilization int8 network, is expected to correlate well with silicon. |
| 4x peak performance | 4x peak performance is based on TOPS increase from MTL (11 TOPS) to LNL (48 TOPS). |
| | SLIDE 34: NPU4 Shave DSP |
| 4x Vector compute | Based on 4x vector width increase vs. NPU3 . NPU3 has 8 FP16 Vector ops/clock, NPU4 has 32 |
| 12x overall vector performance | Vector performance = 3x tiles and 4x the vector width (vs. NPU3 ) |
| | SLIDE 38: NPU 4 Performance |
| 12x vector performance | Vector performance = 3x tiles and 4x the vector width (vs. NPU3 ) |
| 4x TOPS | TOPS calculation is # of tiles * fmax frequency * ops clock<br>Meteor Lake is up to 11.5 TOPS, Lunar Lake is up to 48 TOPS;<br>Meteor Lake TOPS = (2 tiles * 1.4GHz * 4096 ops/clock)/1000<br>Lunar Lake TOPS = (6 tiles * 1.95GHz * 4096 ops/clock)/1000 |
| 2x IP bandwidth | IP Bandwidth: Meteor Lake is 64GB/s; Lunar Lake is 136 GB/s. |
| | SLIDE 55: Stable Diffusion v1.5 |
| Lunar Lake vs. Meteor Lake performance, power and efficiency ratio | Testing by Intel as of May 2024. Data based on Lunar Lake reference validation platform vs. Intel® Core™ Ultra 7 155H 32GB LPDDR5-6400Mhz (Meteor Lake). Calculated using open source GIMP with NPU plug in. Text Encoder, & Unet +/- are running on the NPU. VAE  is running on the built-in GPU. |