# Challenges in Designing and Evaluating Neural Processing Units

Jun-Seok Park

S.LSI, Samsung Electronics

# Outline

- ☐ **Introduction**
  - ■ **Challenges in Designing NPU**
- ☐ NPU architecture overview
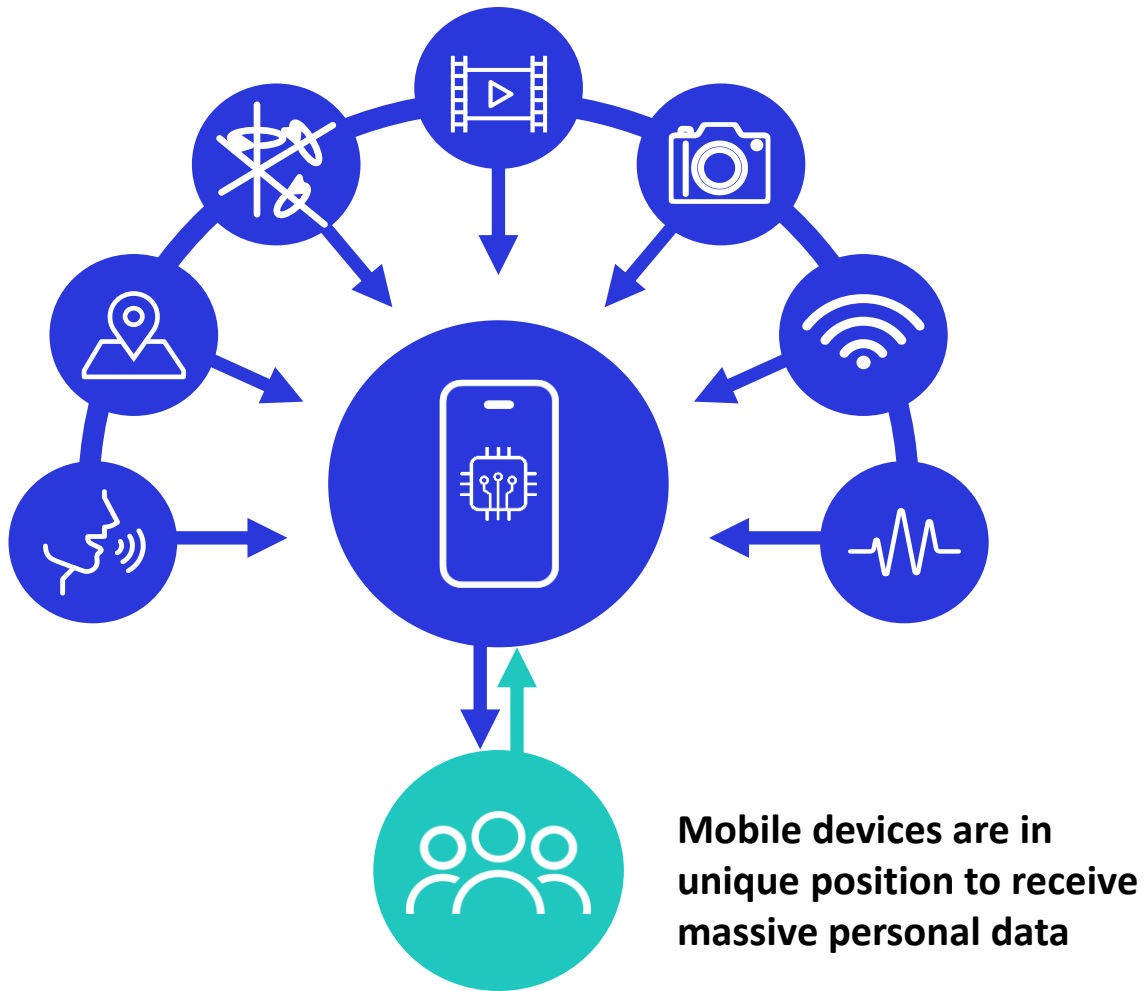- ☐ Evolution of NPU over the years
  - ■ Compute density
  - ■ HW utilization improvement
  - ■ Power efficiency & heterogeneous NPU
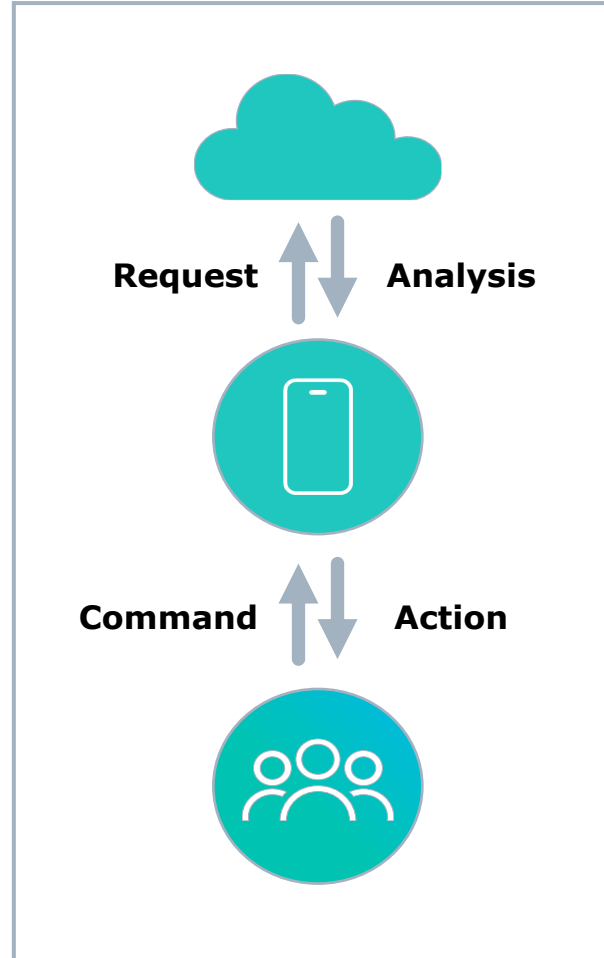- ☐ Challenges in Evaluating NPU
  - ■ Domain-specific architectures (Heterogeneous NPUs, CIM/PIM)
  - ■ Fast-Evolving Neural Networks
  - ■ Rapid scaling-up of NPUs outpacing CPUs or DRAMs
  - ■ Expoliting HW parallelism (Tensor-Vector, Memory pre-fetch)
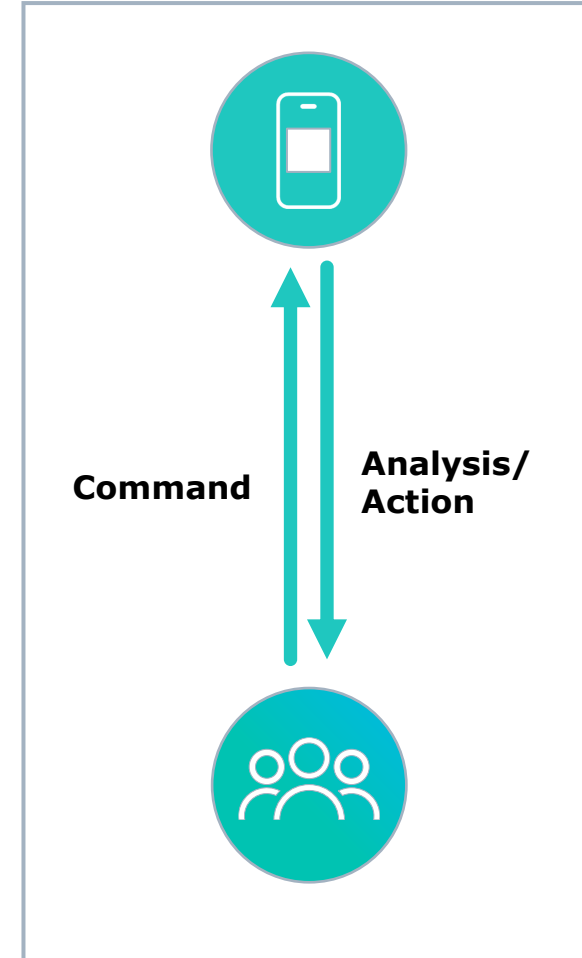- ☐ Conclusion

# On-device AI



Mobile devices are in unique position to receive massive personal data

**AI using Cloud Servers**

Request ↕ Analysis

Command ↕ Action

**On-Device AI**

Command / Analysis/ Action
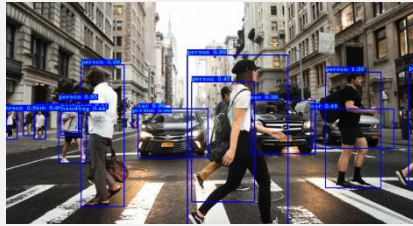
# Endless Applications



**Object detection**

**Voice recognition**

**Live text**

**Reflection removal**

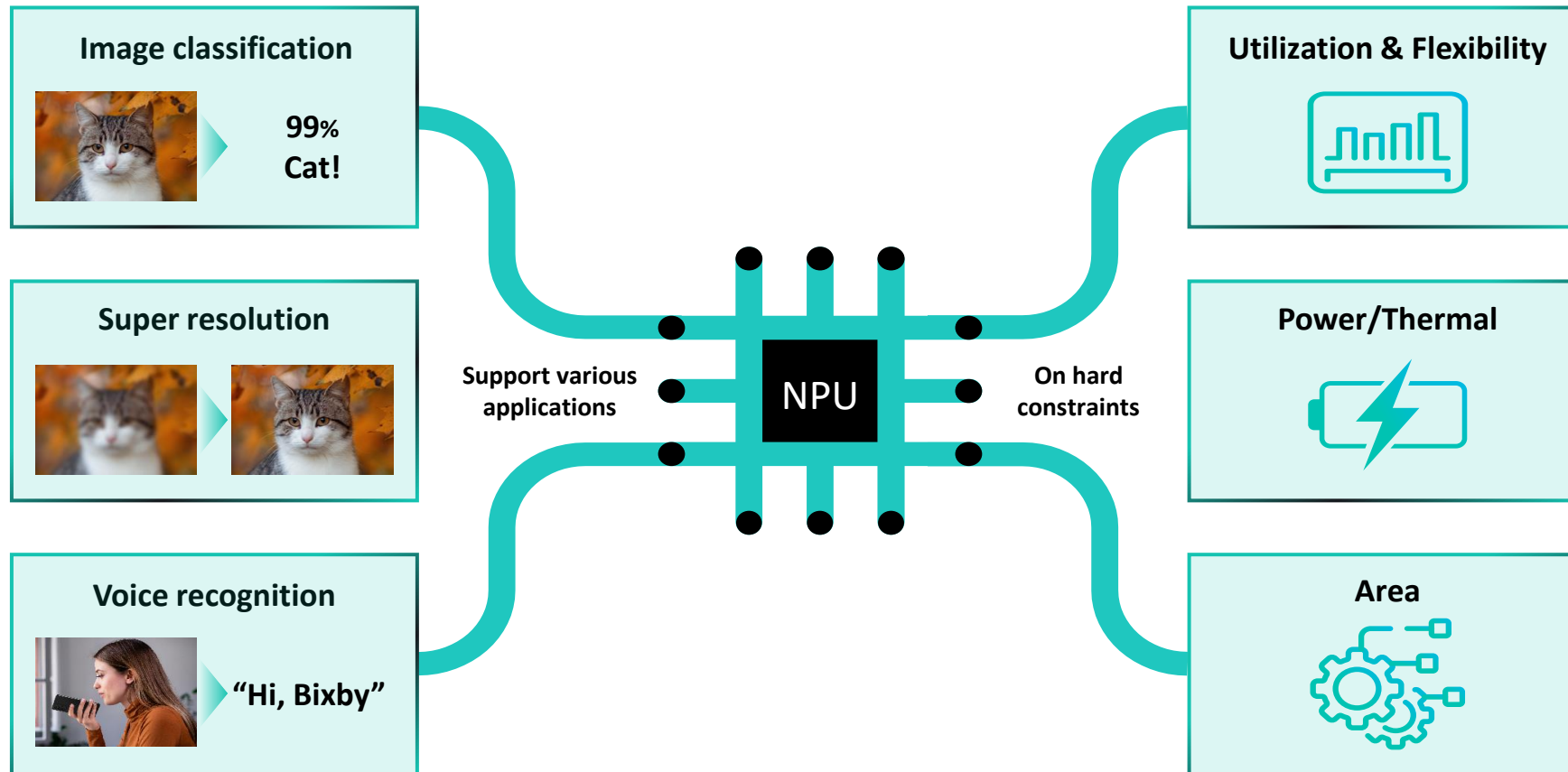**Object removal**

Original    Result

**Photo remastering**

Original    Remastered

**Scene detection**

**Face landmark detection**

**Bokeh**

# NPU on a Mobile AP



NPU

**Image classification**

99% Cat!

**Super resolution**

**Voice recognition**

"Hi, Bixby"

Support various applications

On hard constraints

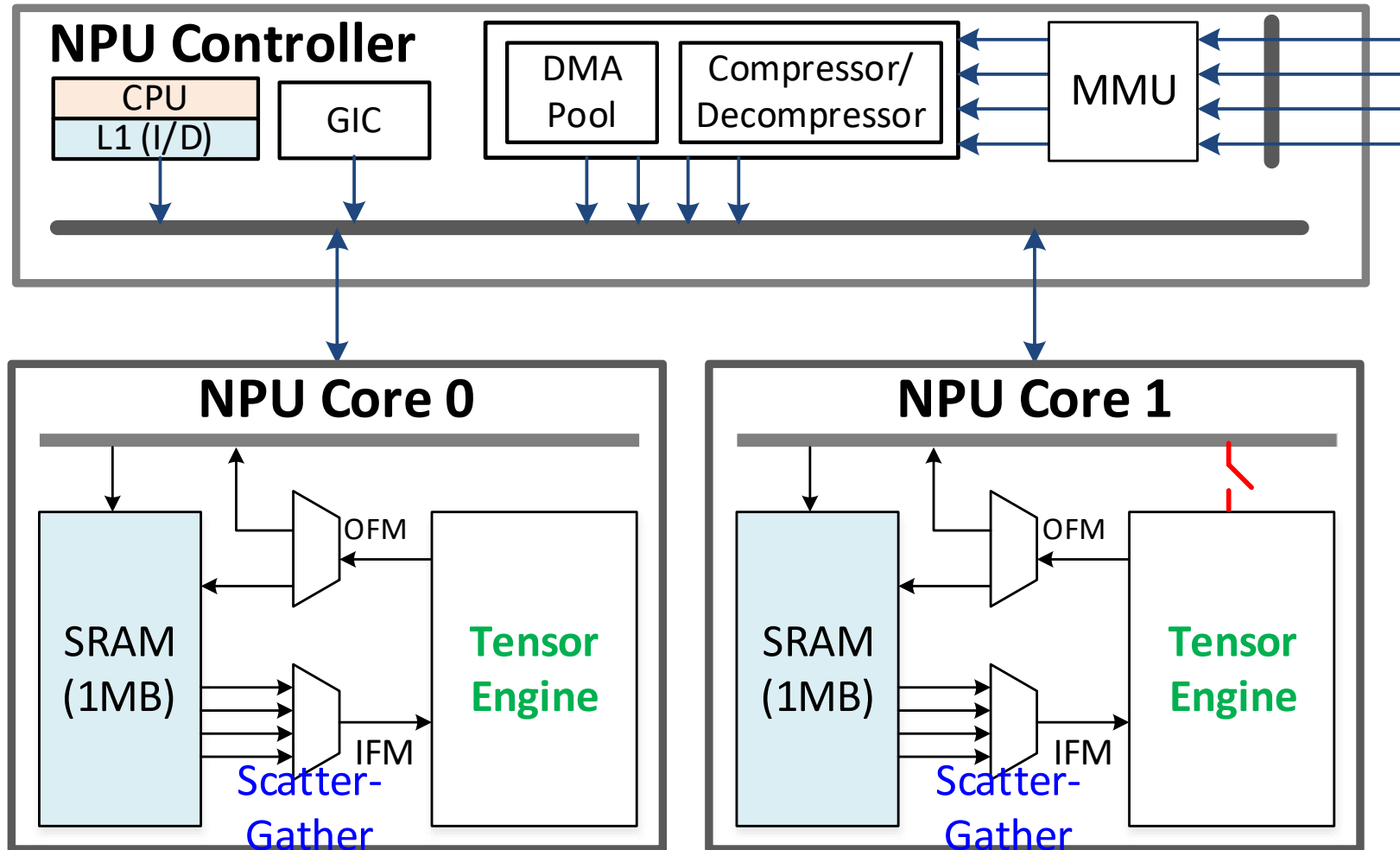**Utilization & Flexibility**

**Power/Thermal**

**Area**

☐ NPU on a mobile AP should be area and energy efficient.

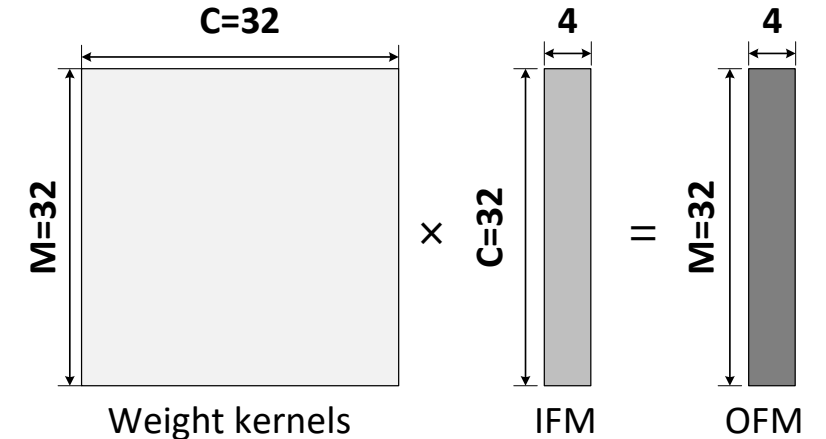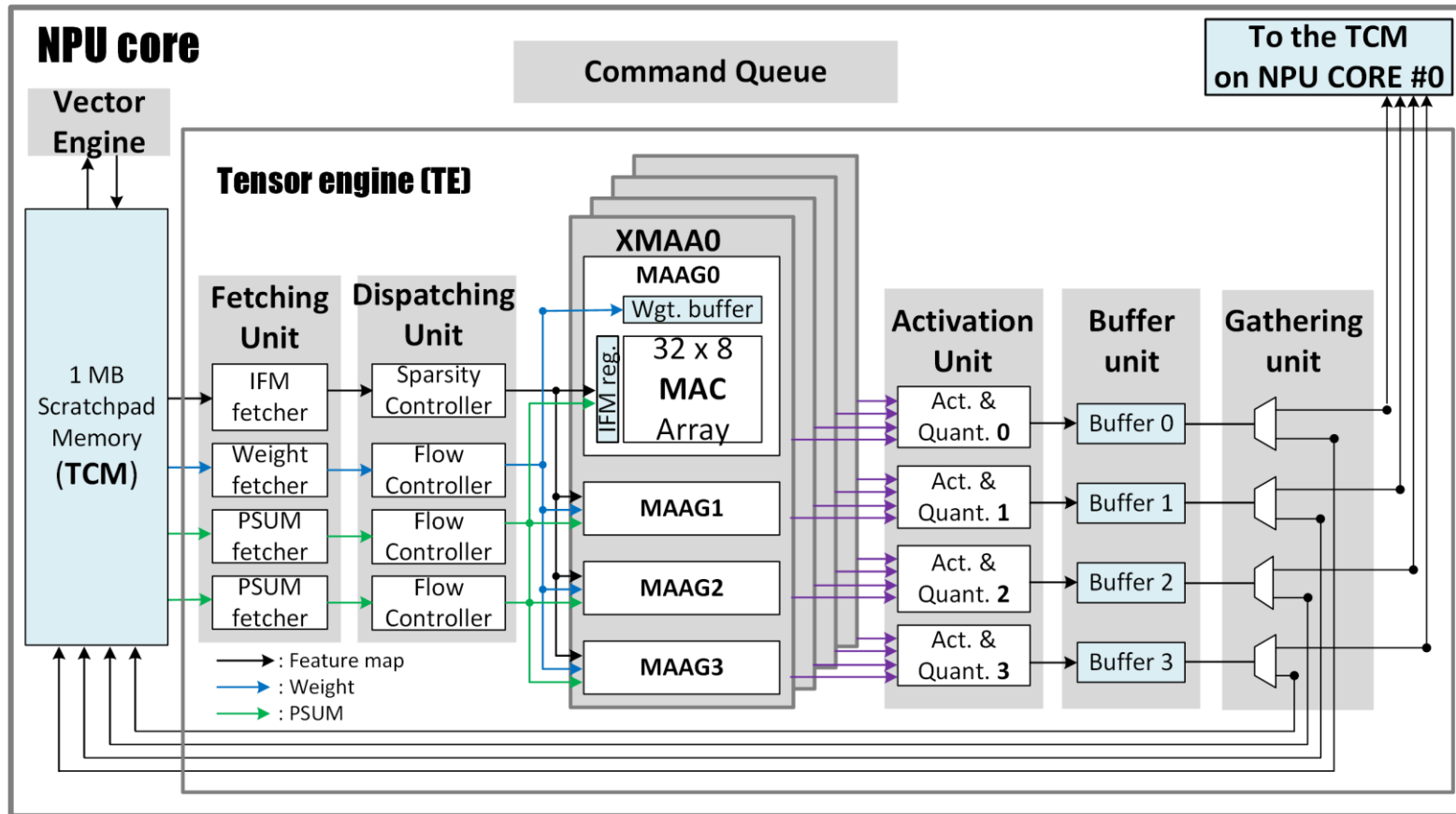☐ Plus, it should support a comprehensive range of NN.

# Outline

- ☐ Introduction
  - ■ Challenges in Designing NPU
- ☐ **NPU architecture overview**
- ☐ Evolution of NPU over the years
  - ■ Compute density
  - ■ HW utilization improvement
  - ■ Power efficiency & heterogeneous NPU
- ☐ Challenges in Evaluating NPU
  - ■ Domain-specific architectures (Heterogeneous NPUs, CIM/PIM)
  - ■ Fast-Evolving Neural Networks
  - ■ Rapid scaling-up of NPUs outpacing CPUs or DRAMs
  - ■ Expoliting HW parallelism (Tensor-Vector, Memory pre-fetch)
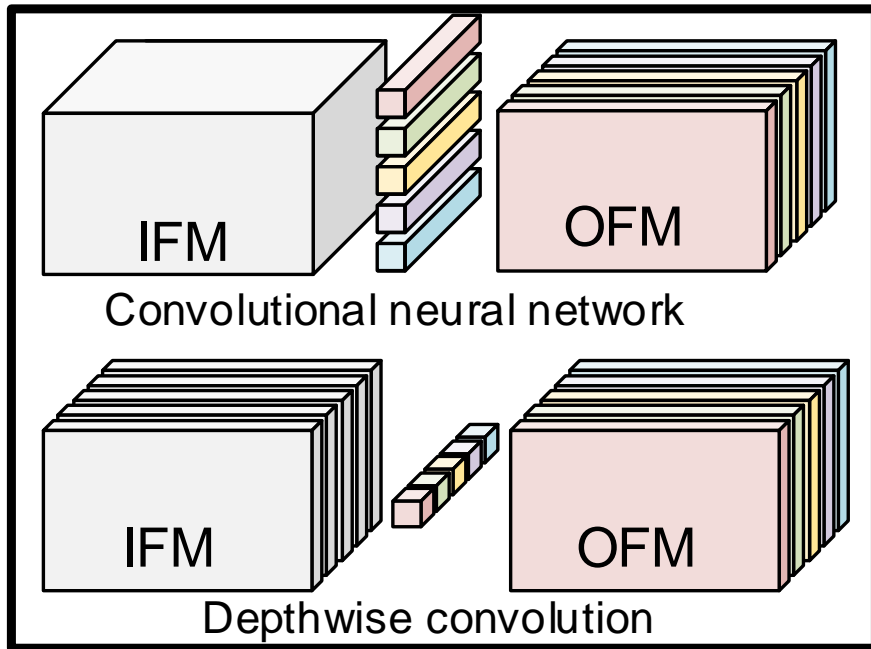- ☐ Conclusion

# NPU architecture

# Tensor Engine (TE)



Weight kernels × IFM = OFM

$C=32$, $M=32$, $4$, $4$

**One NPU core**
**: 4K** MACs / core
= **32** (C) x **32** (M) x 4 (**X**)
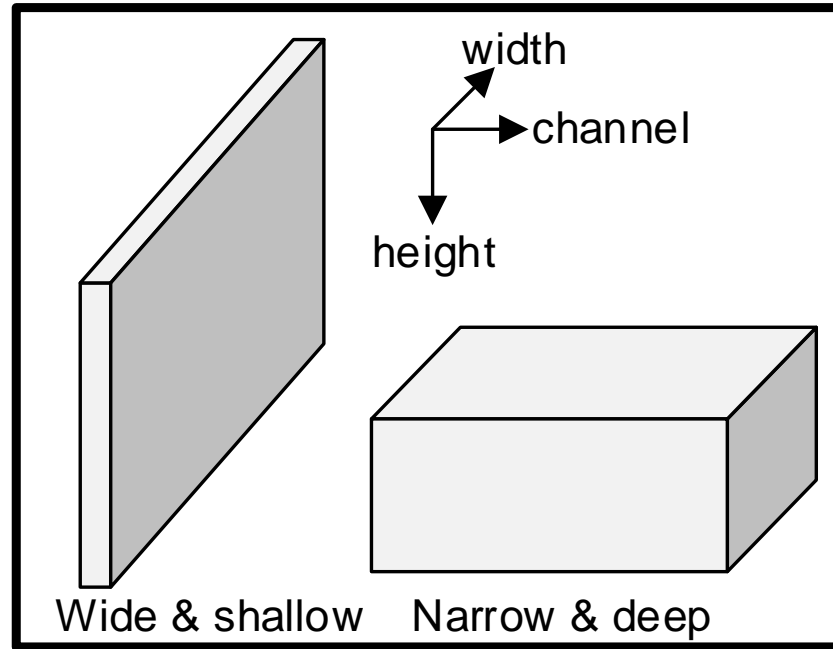
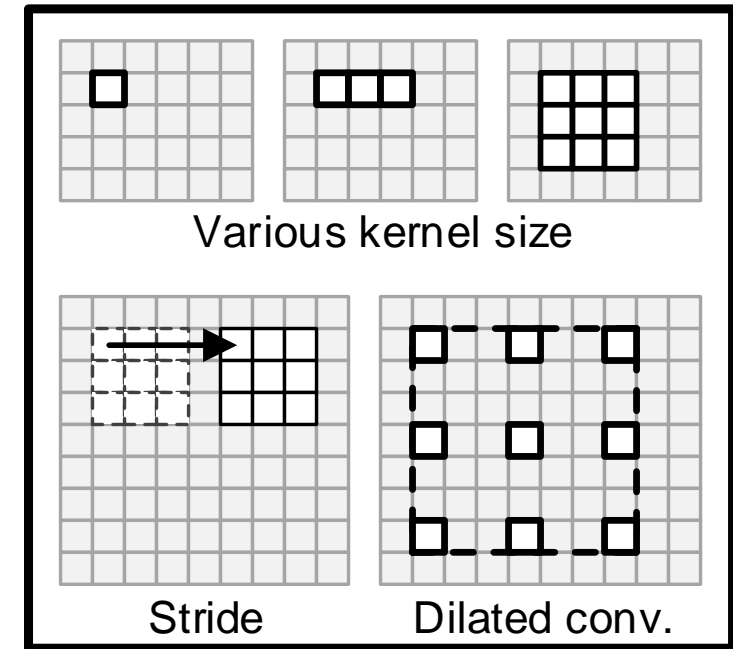☐ TE executes four 32-elements dot-products of 32-dim vectors per cycles

# Motivation - Flexibility



**Types of layer**

IFM — Convolutional neural network — OFM

IFM — Depthwise convolution — OFM

**Feature map shape**

width, channel, height

Wide & shallow    Narrow & deep

**Kernels**

Various kernel size

Stride    Dilated conv.
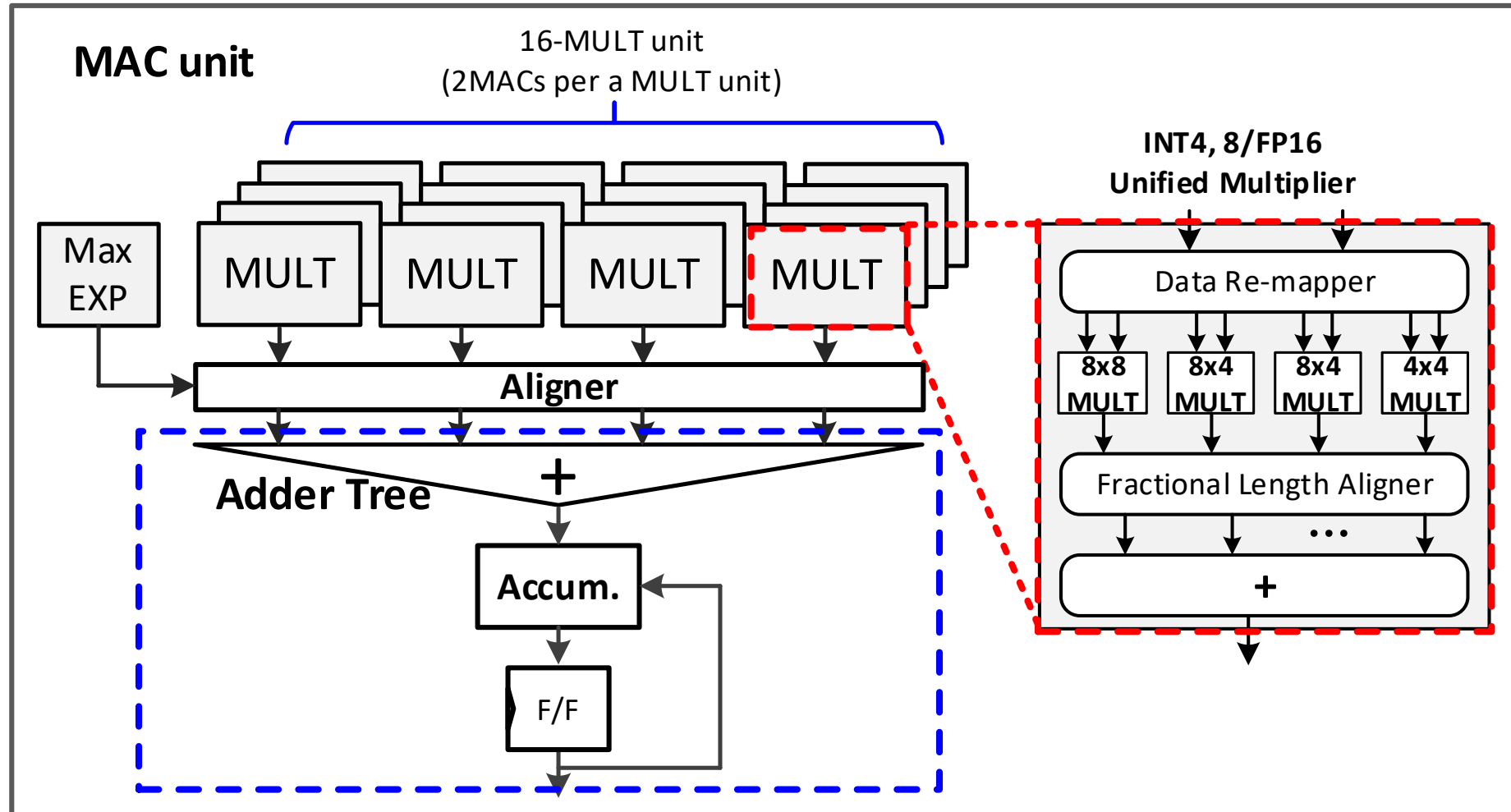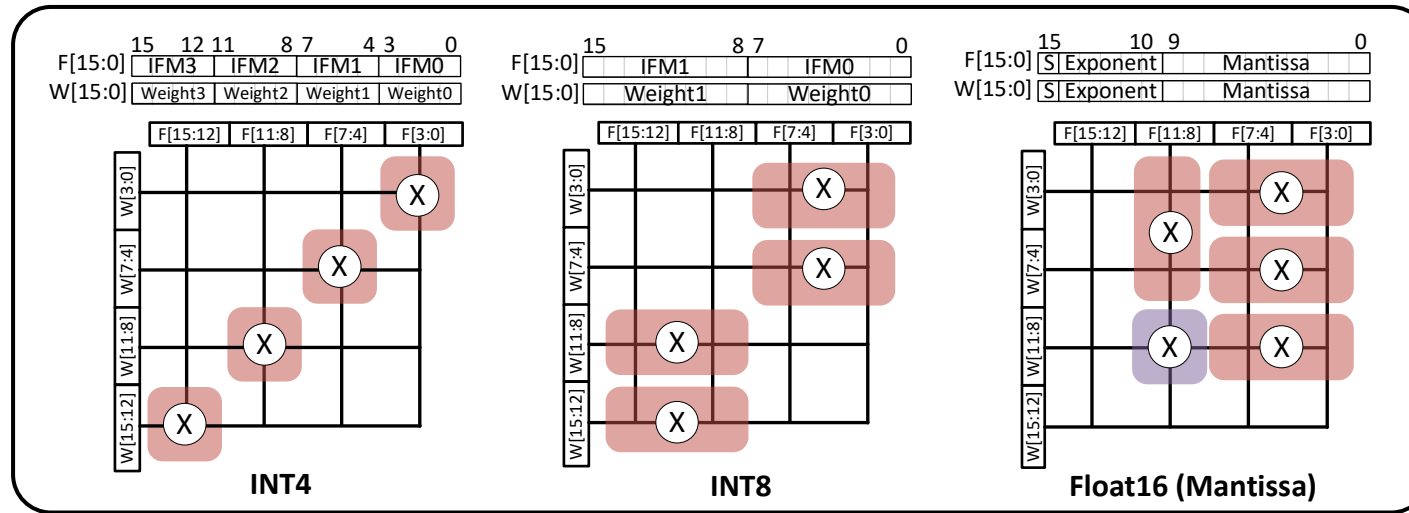
☐ NPU needs to support a comprehensive range of NN

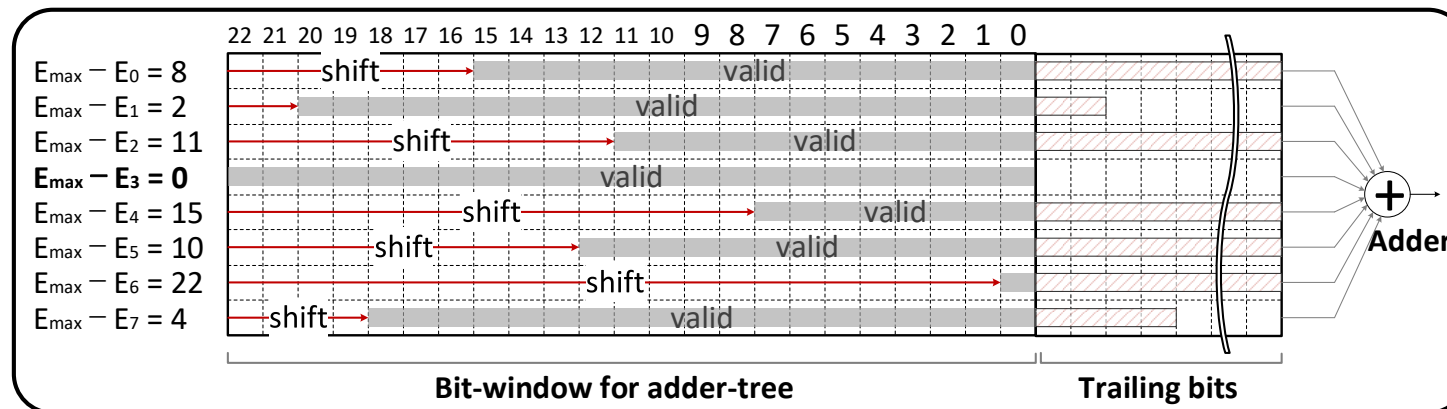   ■ Diverse kernel sizes, dilation rates and strides

# Unified Multi-precision MAC

# Multi-precision MAC



(a) Reconfigurable datapath in the unified multiplier



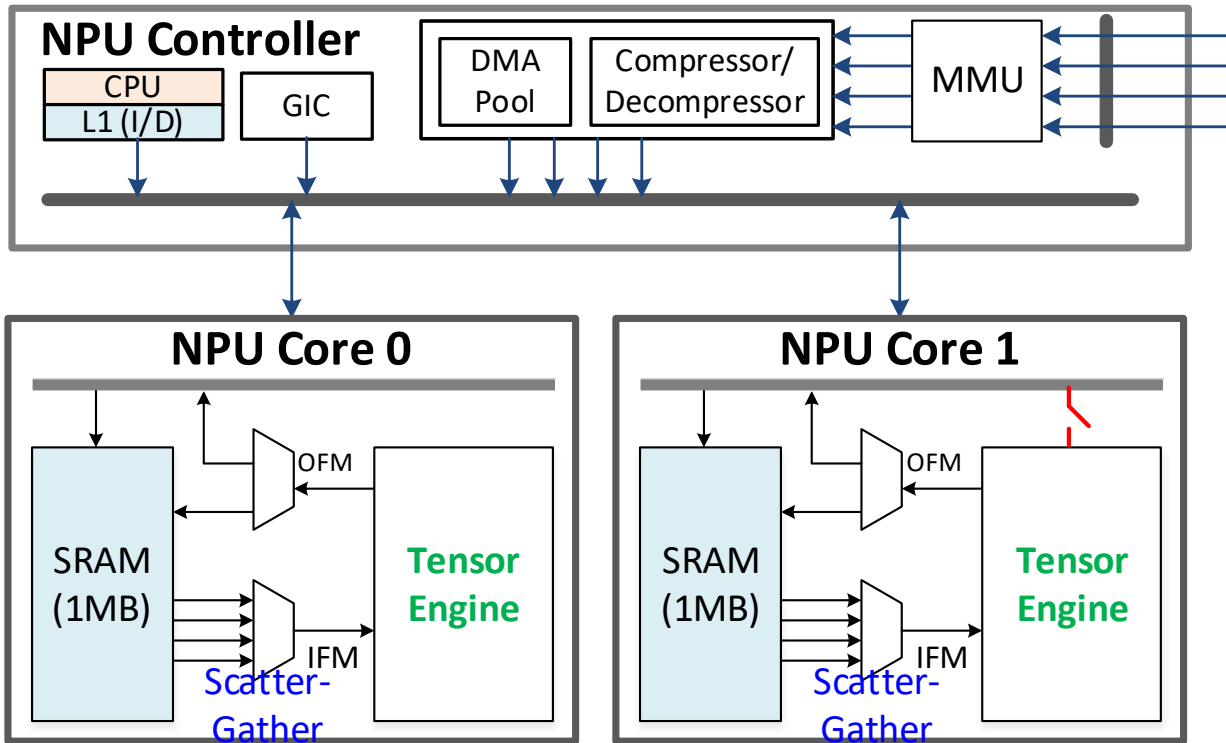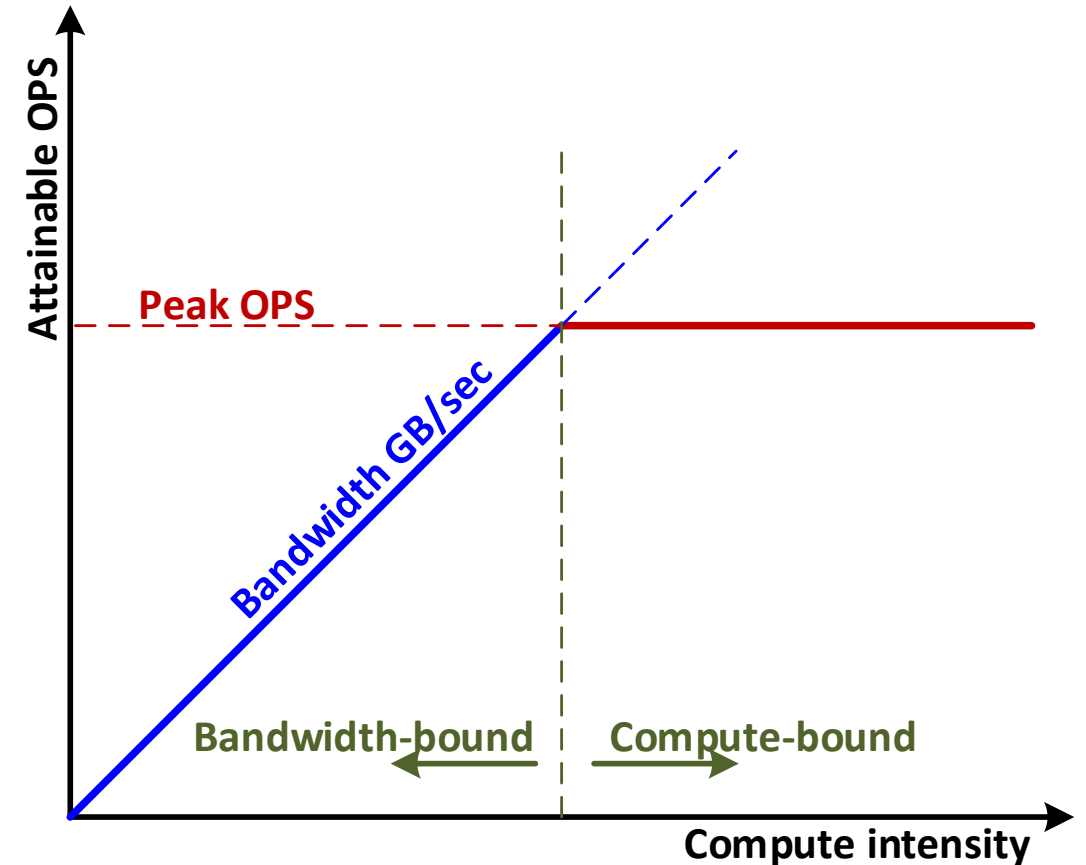(b) Alignment of mantissa in the fused dot-product for FP16

Challenges in Designing and Evaluating Neural
Processing Units

# Outline

- ☐ Introduction
  - ■ Challenges in Designing NPU
- ☐ NPU architecture overview
- ☐ **Evolution of NPU over the years**
  - ■ **Compute density**
  - ■ **HW utilization improvement**
  - ■ **Power efficiency & heterogeneous NPU**
- ☐ Challenges in Evaluating NPU
  - ■ Domain-specific architectures (Heterogeneous NPUs, CIM/PIM)
  - ■ Fast-Evolving Neural Networks
  - ■ Rapid scaling-up of NPUs outpacing CPUs or DRAMs
  - ■ Expoliting HW parallelism (Tensor-Vector, Memory pre-fetch)
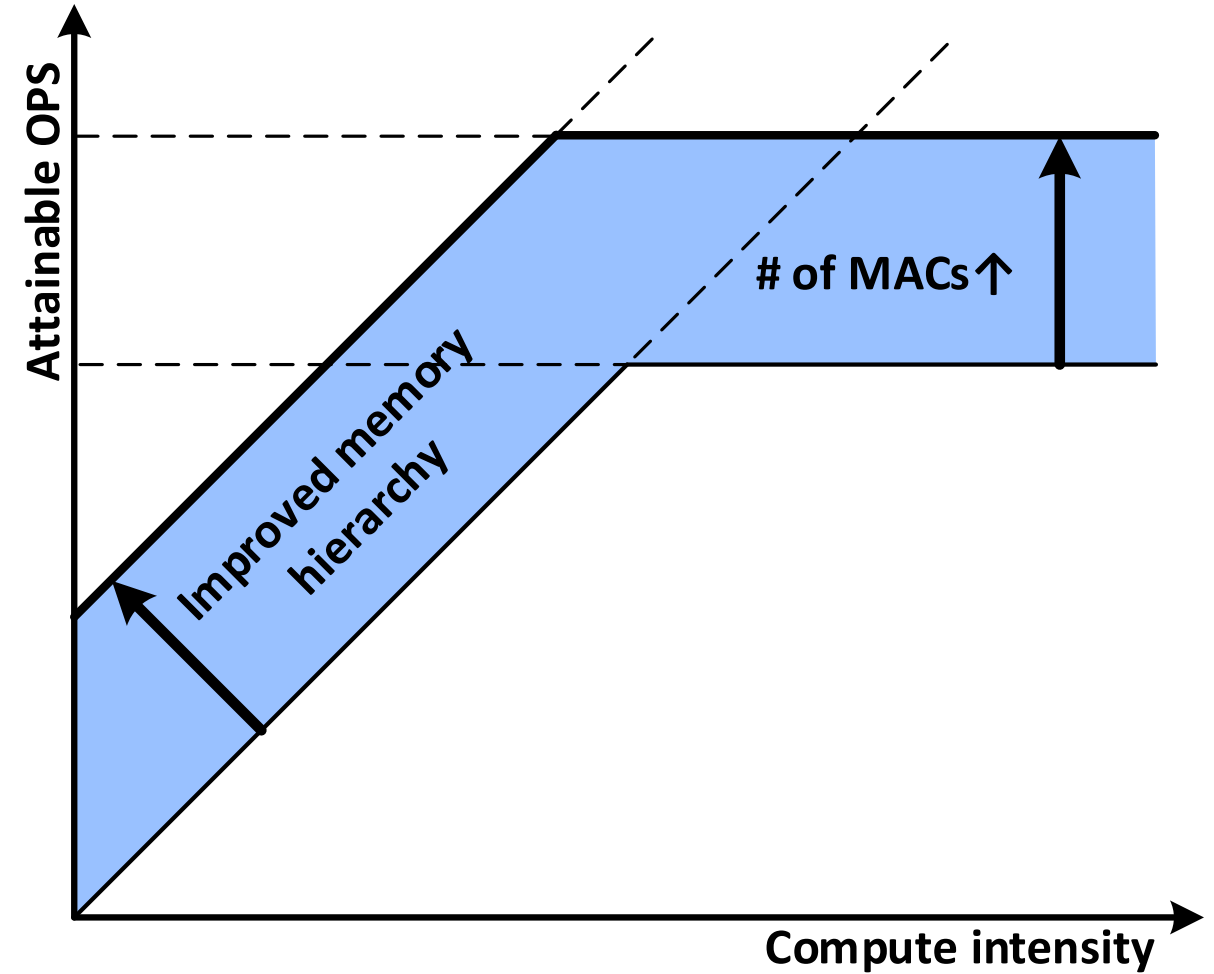- ☐ Conclusion

# NPU architecture

☐ A block diagram of an NPU

☐ Roofline of an NPU architecture

**NPU Controller**

CPU
L1 (I/D)
GIC
DMA Pool
Compressor/Decompressor
MMU

**NPU Core 0**

SRAM (1MB)
OFM
Tensor Engine
IFM
Scatter-Gather

**NPU Core 1**

SRAM (1MB)
OFM
Tensor Engine
IFM
Scatter-Gather

Attainable OPS

Peak OPS

Bandwidth GB/sec

Bandwidth-bound     Compute-bound

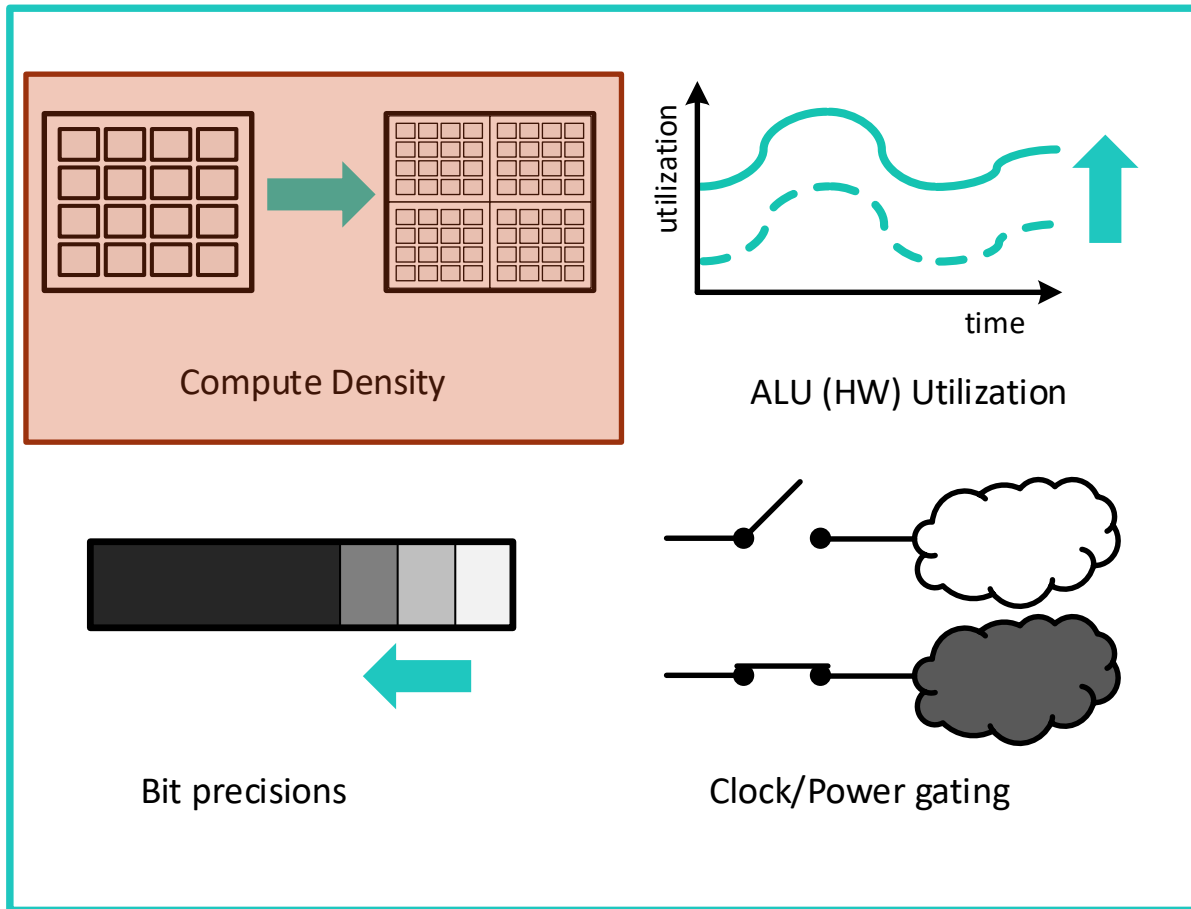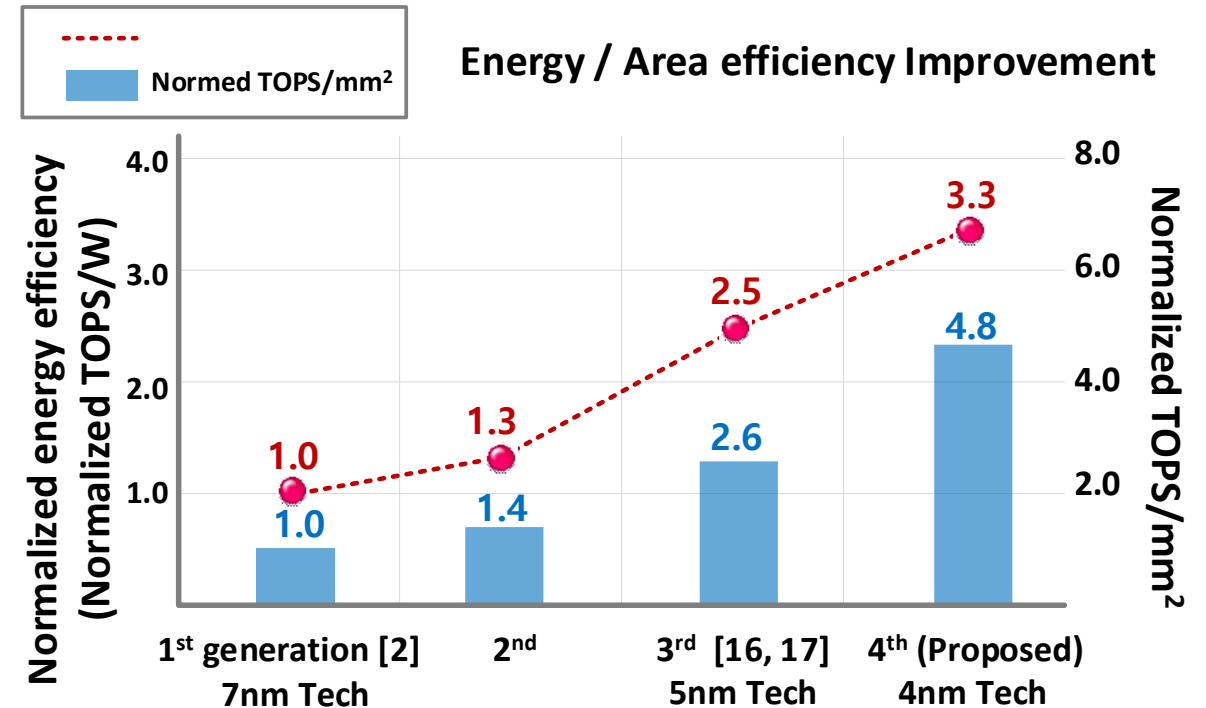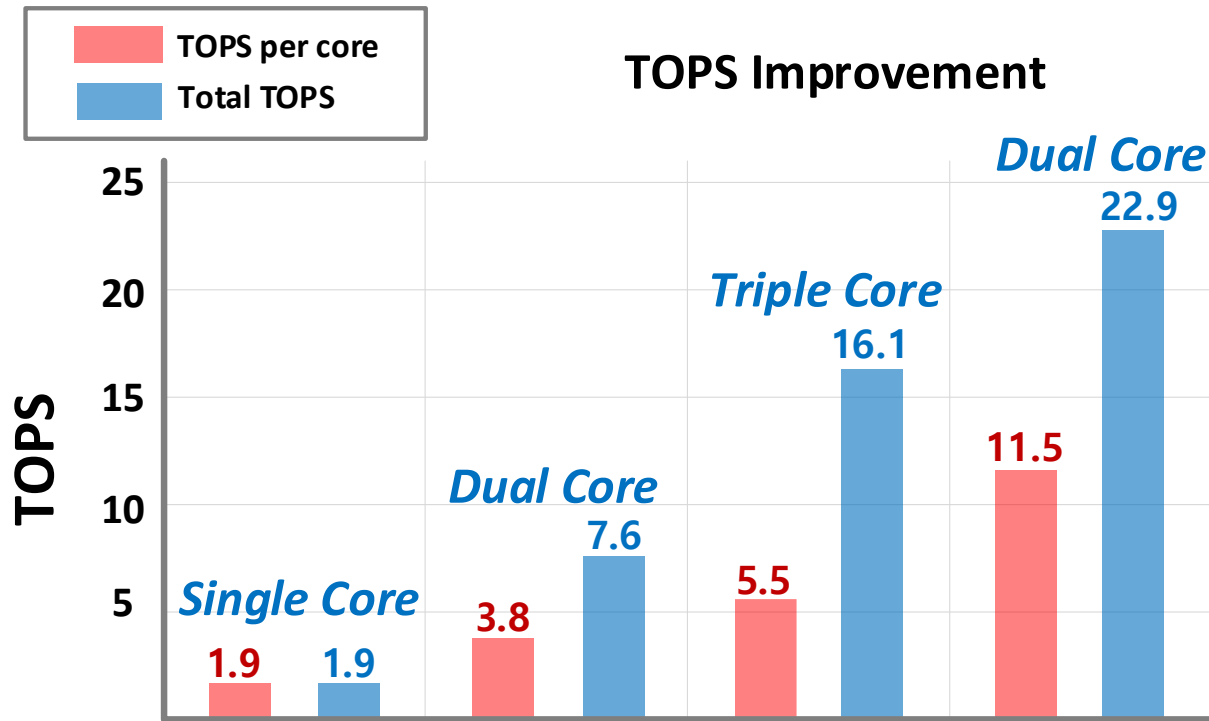Compute intensity

[7] S Williams, "Roofline: an insightful visual performance model for multicore architectures", 2009.

# NPU Enhancement: Compute Density



Compute Density

ALU (HW) Utilization

Bit precisions

Clock/Power gating

Attainable OPS

Improved memory hierarchy

# of MACs↑

Compute intensity

# Increase the compute density



**TOPS Improvement**

Legend: TOPS per core, Total TOPS

- Single Core: 1.9 / 1.9
- Dual Core: 3.8 / 7.6
- Triple Core: 5.5 / 16.1
- Dual Core: 11.5 / 22.9

TOPS

**Energy / Area efficiency Improvement**

Legend: Normed TOPS/mm$^2$

Normalized energy efficiency (Normalized TOPS/W) — Normalized TOPS/mm$^2$

- 1st generation [2] 7nm Tech: 1.0 / 1.0
- 2nd: 1.3 / 1.4
- 3rd [16, 17] 5nm Tech: 2.5 / 2.6
- 4th (Proposed) 4nm Tech: 3.3 / 4.8

# Feature Map Lossless Compressor



**Quad-Tree Header**

: **Non-zero Features**

: **Zero Features**

**Feature-map Groups**

**Compressed feature map**

Meta-data

Stream length | 010 | 0110 | 1110 | 1111 | 1111 | ⋯ | 0111 | xxx | xxx | ⋯ | xxx |

Truncated nonzero bitwidth

Quad-tree Header

Nonzero Features
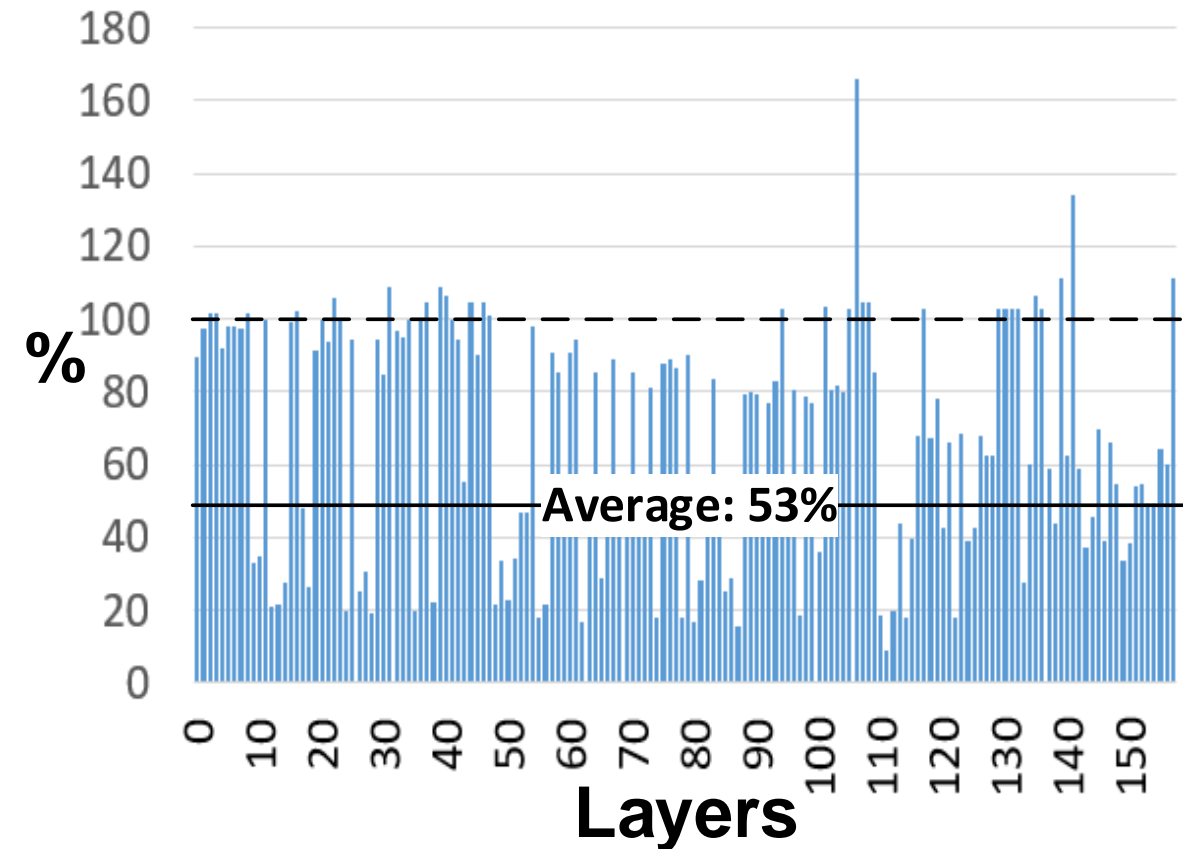
# Percentage of Compressed FM Size
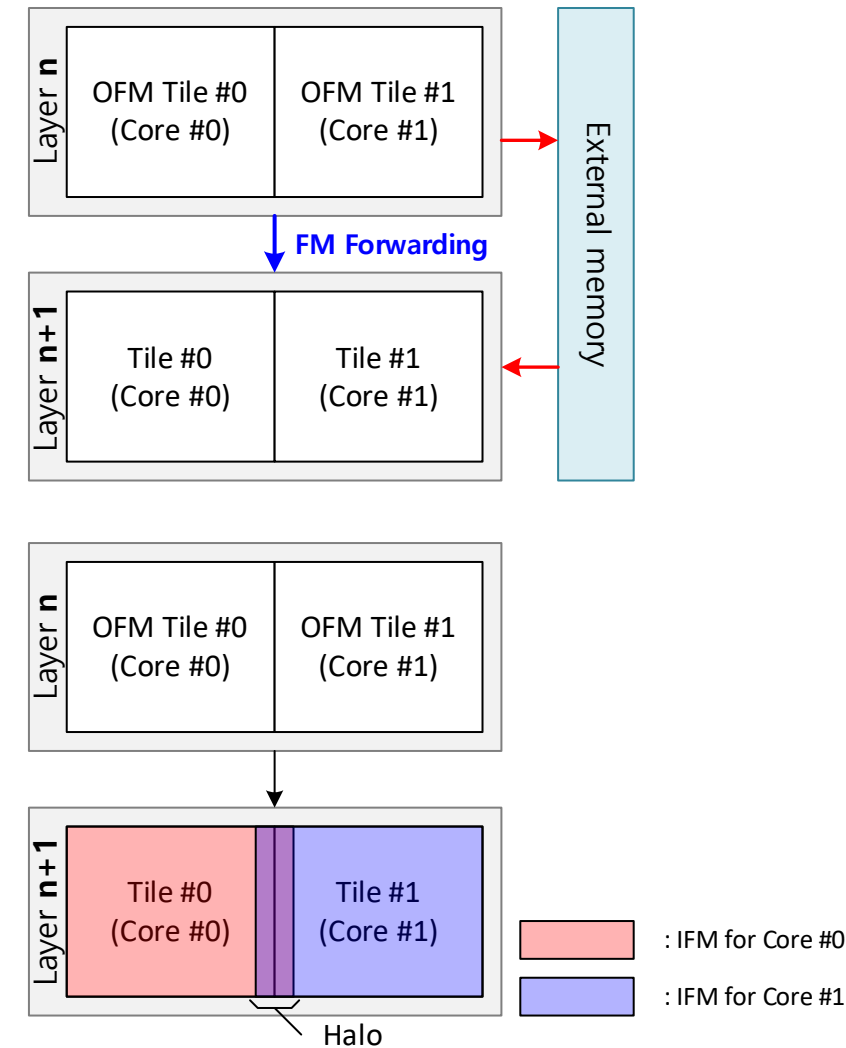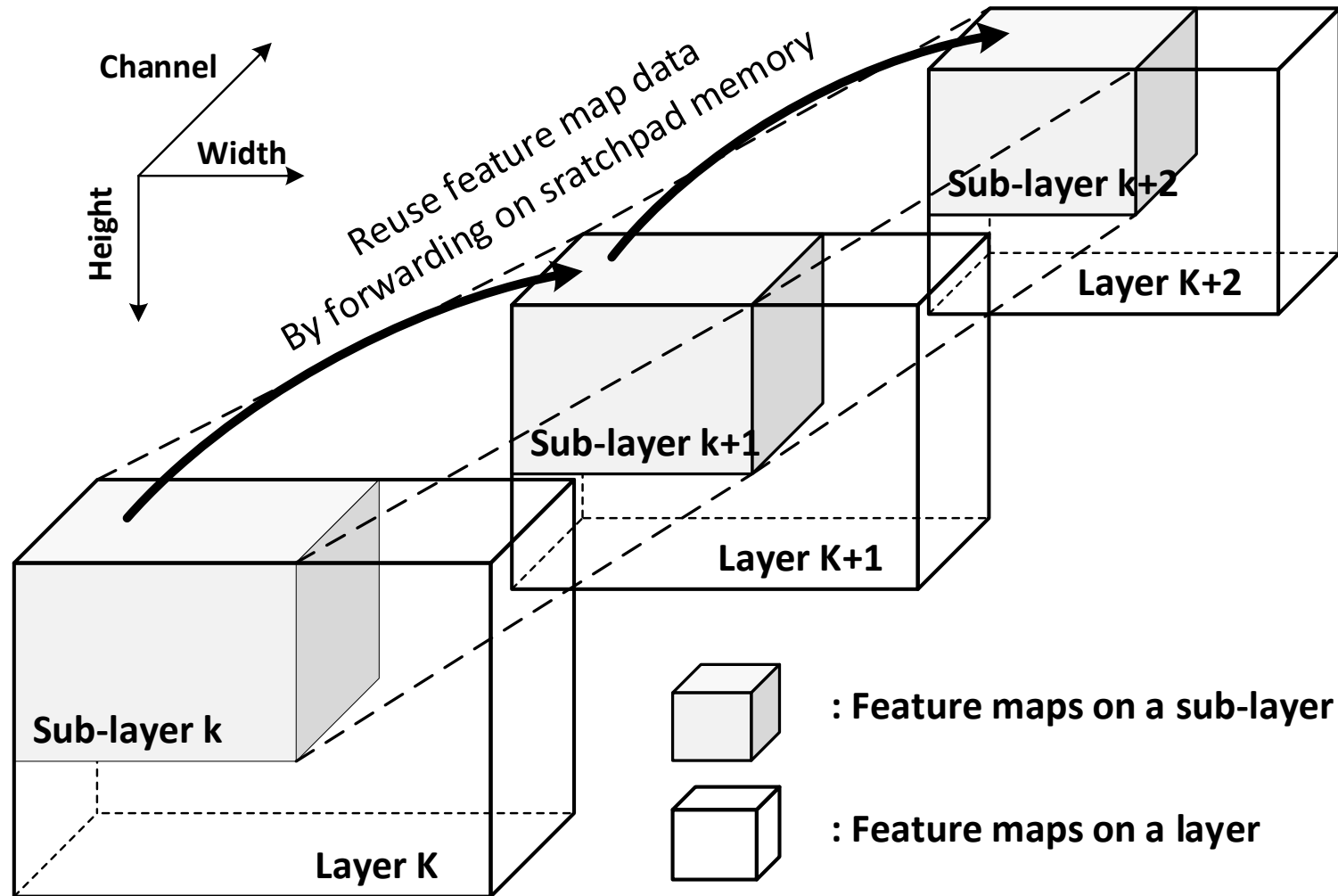


(a) Inception V3

(b) DeepLab V3

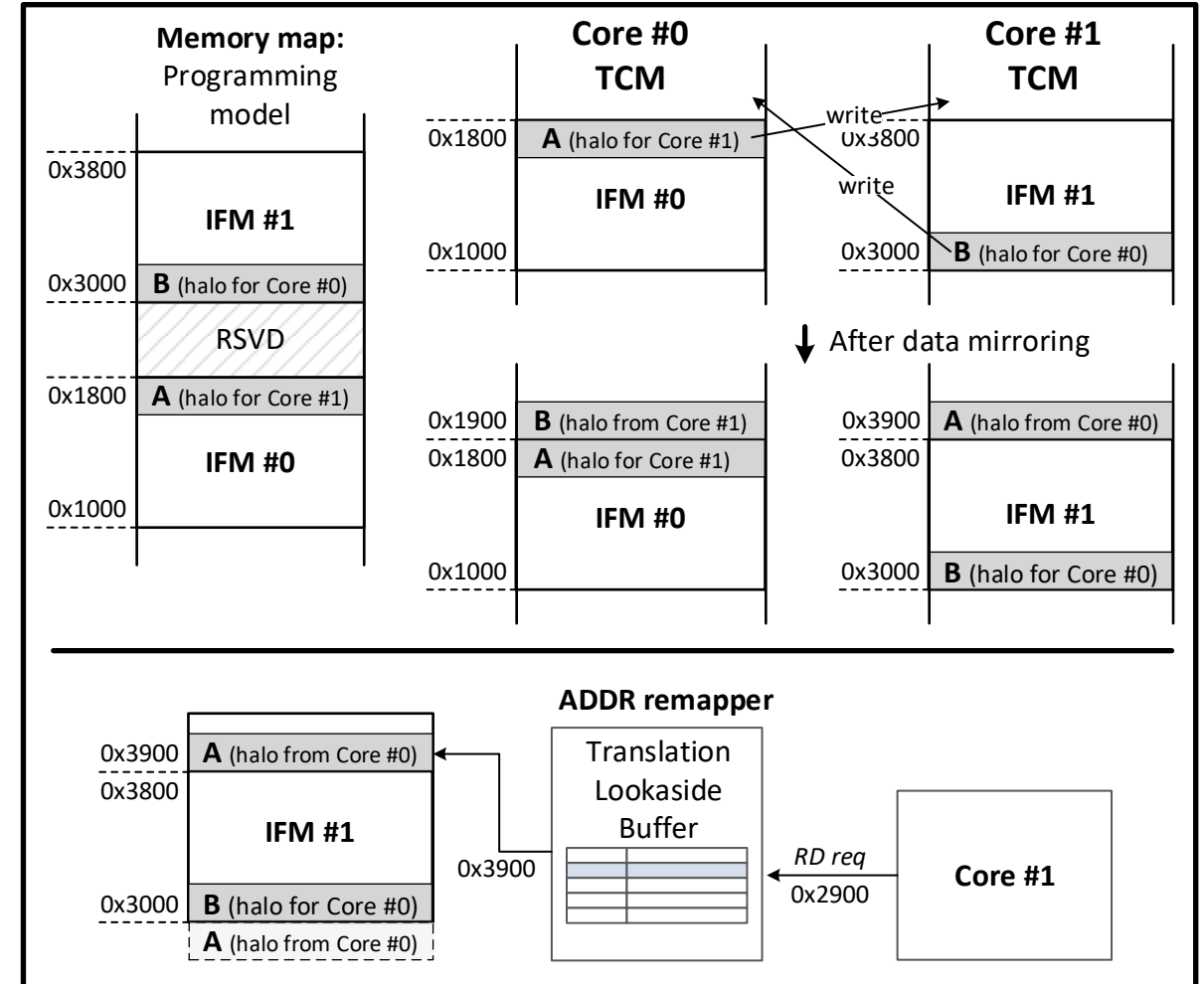# Feature-map forwarding



: Feature maps on a sub-layer

: Feature maps on a layer

Channel

Width

Height

Reuse feature map data
By forwarding on sratchpad memory

Sub-layer k
Layer K

Sub-layer k+1
Layer K+1

Sub-layer k+2
Layer K+2

Layer n: OFM Tile #0 (Core #0), OFM Tile #1 (Core #1)

FM Forwarding

Layer n+1: Tile #0 (Core #0), Tile #1 (Core #1)

External memory

Layer n: OFM Tile #0 (Core #0), OFM Tile #1 (Core #1)

Layer n+1: Tile #0 (Core #0), Tile #1 (Core #1)
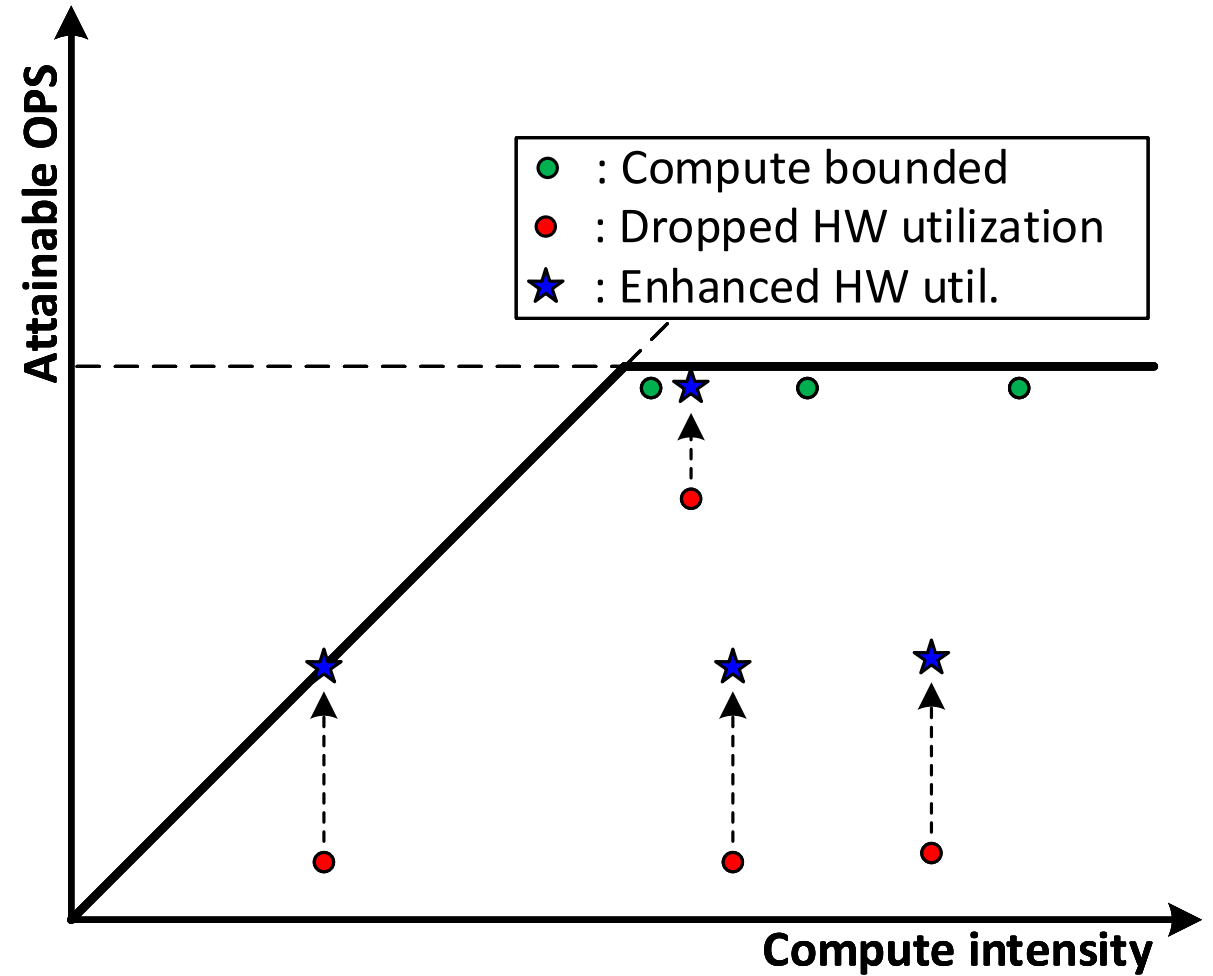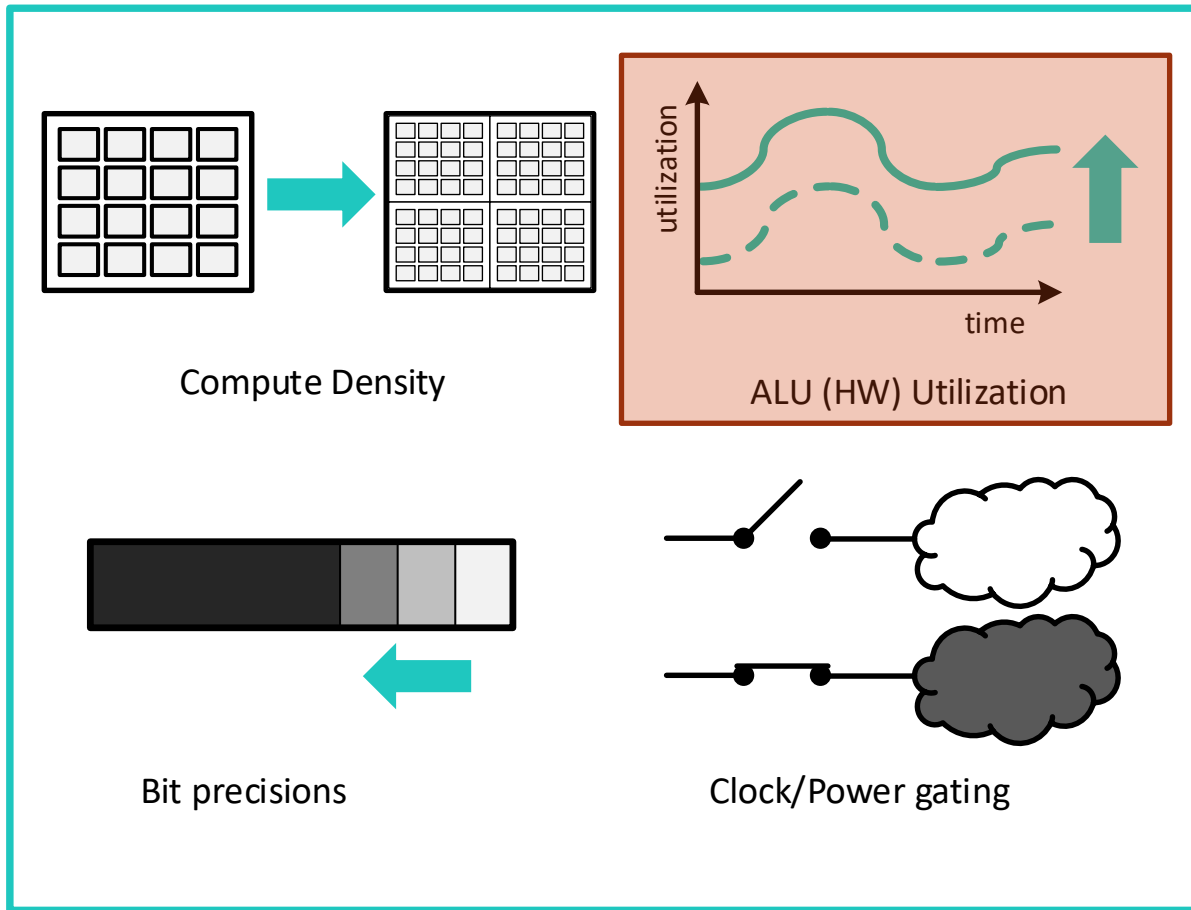
Halo

: IFM for Core #0

: IFM for Core #1
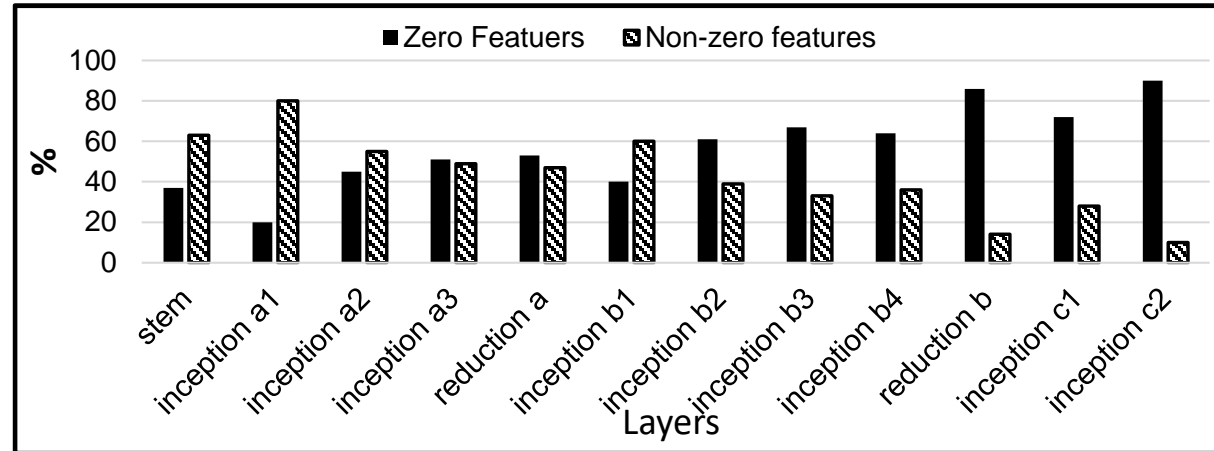
# Low-latency mode



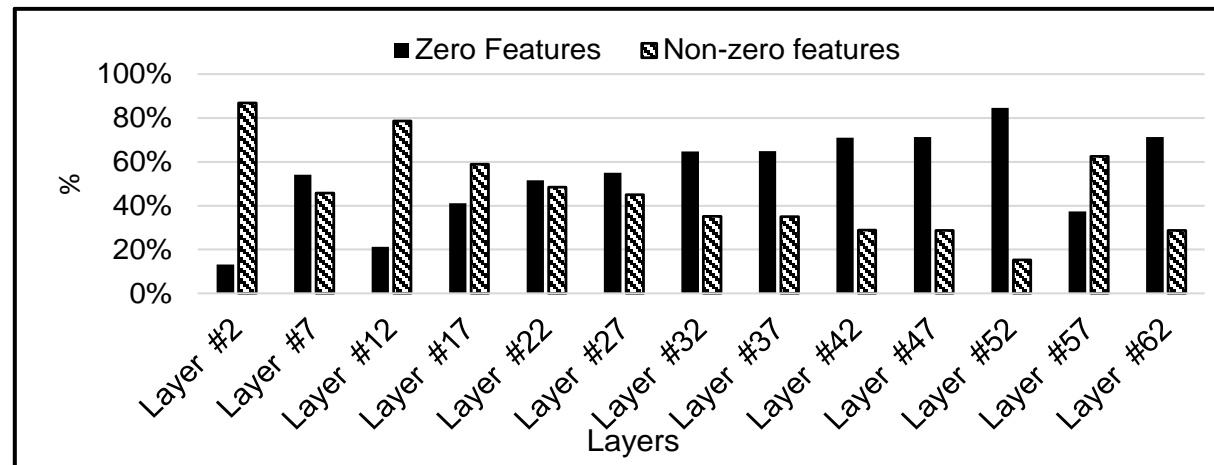- Two NPU cores process a network together to minimize latency.

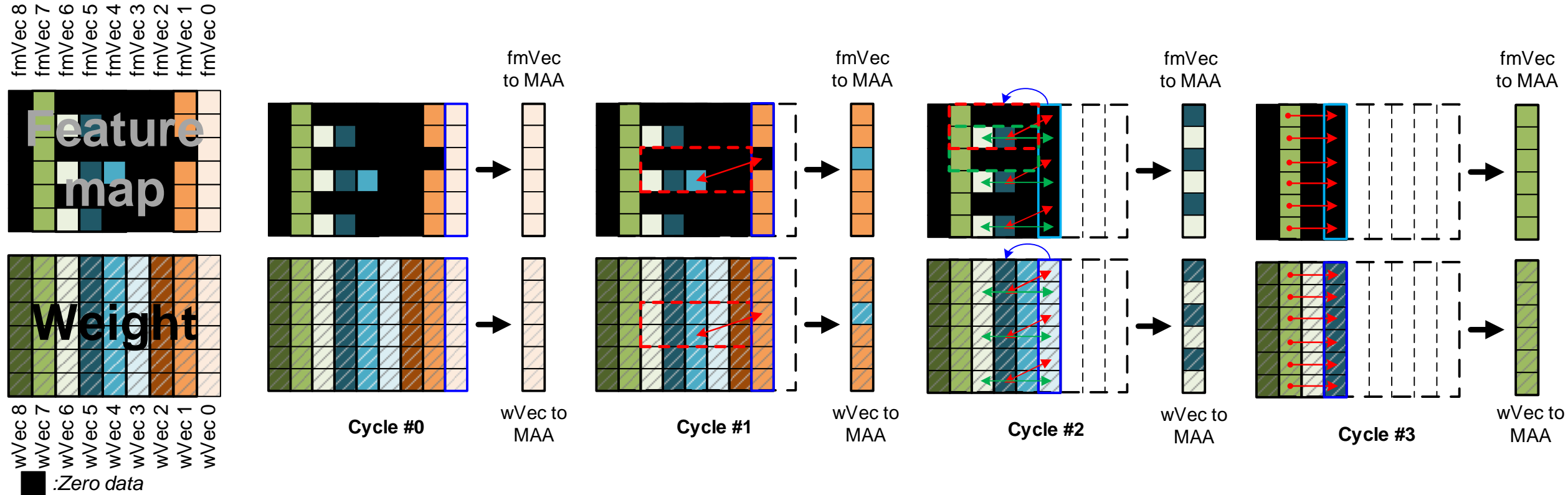# NPU Enhancement: HW Utilization

# Zeroness on Feature map



**Feature map Distribution for Neural Layers**



**Feature map Distribution for Neural Layers on DeepLabV3**

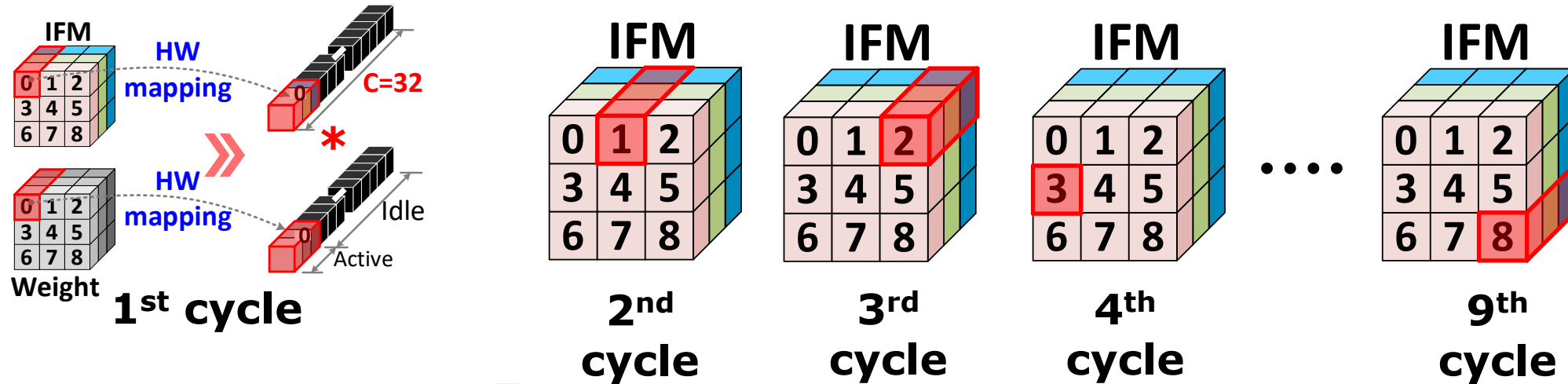Challenges in Designing and Evaluating Neural
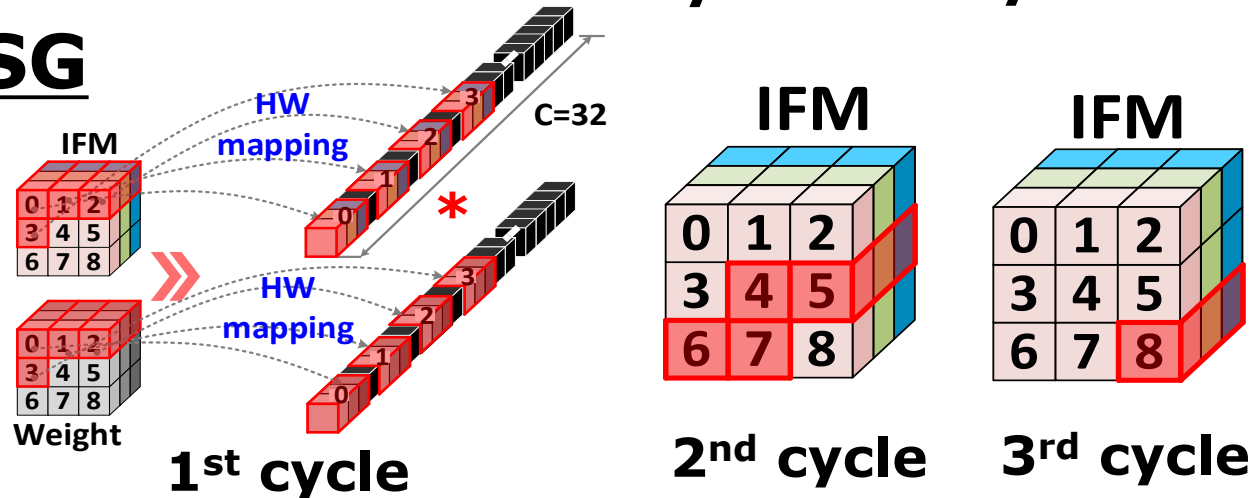Processing Units

# Feature-map aware zero skipping



□ Motivation: A lot of zeros on feature-maps, No effect on the final result.

□ Find and move non-zero in the search window

  ■ Corresponding weight data is loaded.

# Example of Scatter-Gather (SG)

## 3x3 convolution in a layer having 3 input channel



**1st cycle**

**2nd cycle**

**3rd cycle**

**4th cycle**

**9th cycle**

## With SG



**1st cycle**

**2nd cycle**

**3rd cycle**
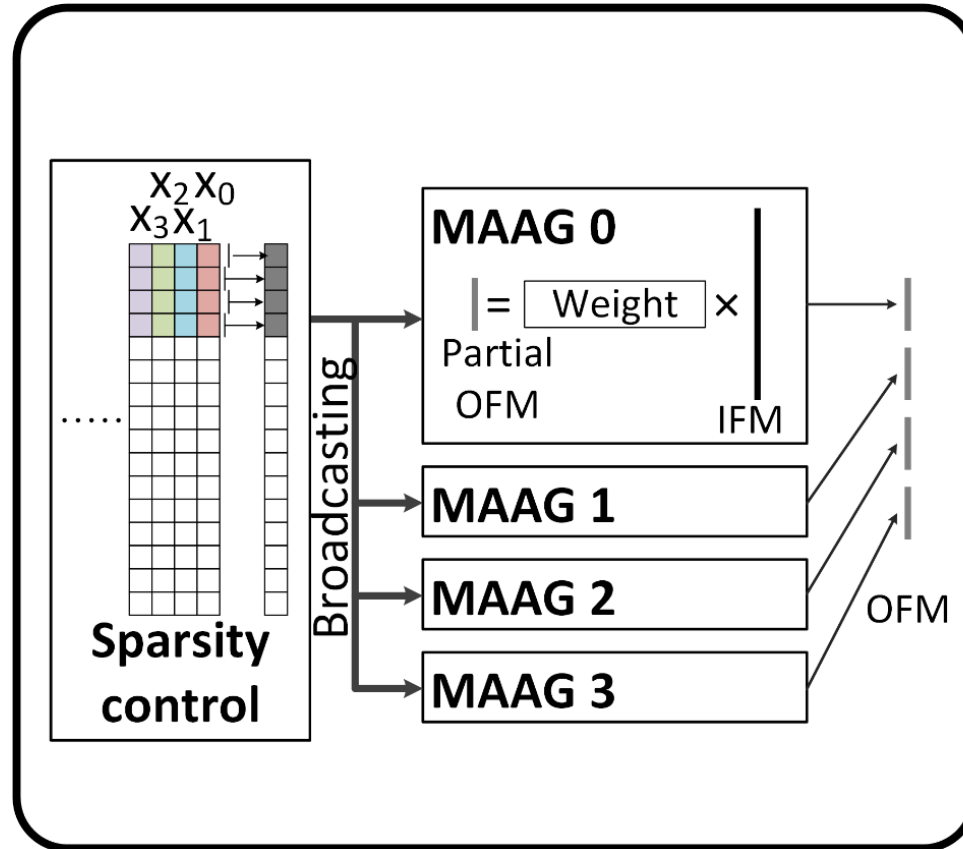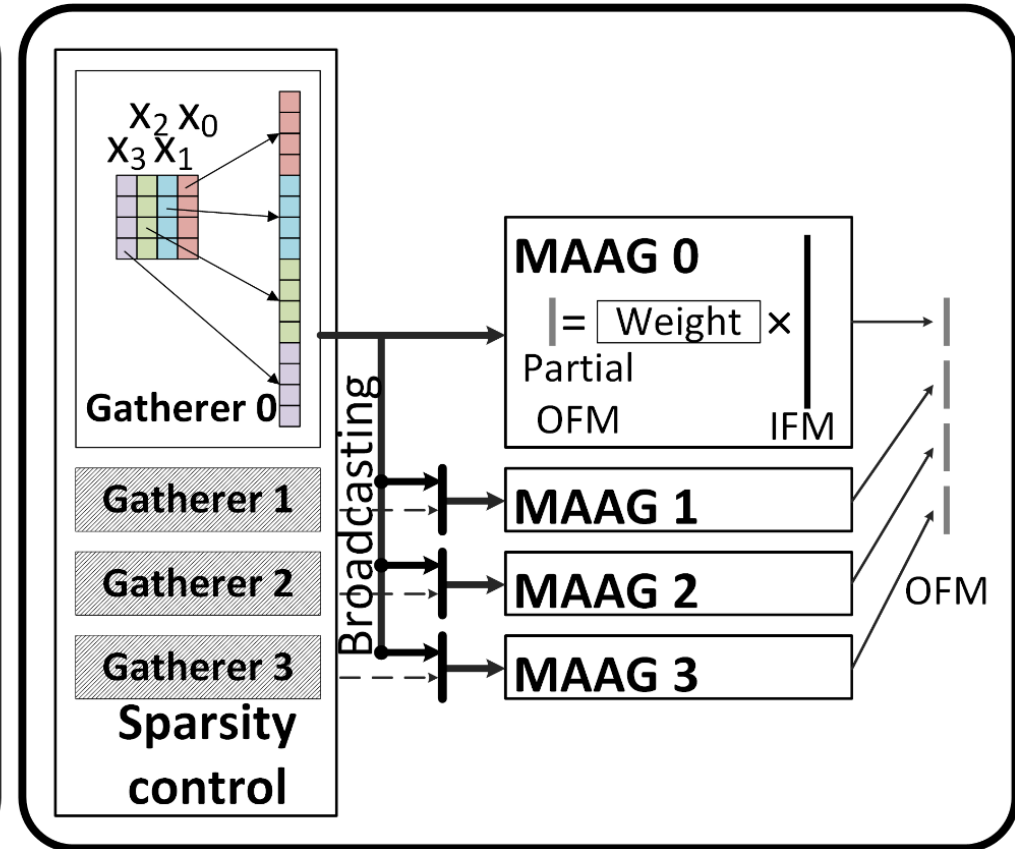
- **Utilization without S/G**
  - 9.4% (= 3/32)
- **Utilization with S/G**
  - **37.5% (= 12/32)**

# Enhance HW Util with Scatter-Gather (SG)



(a) Tensor engine without Scatter-Gather
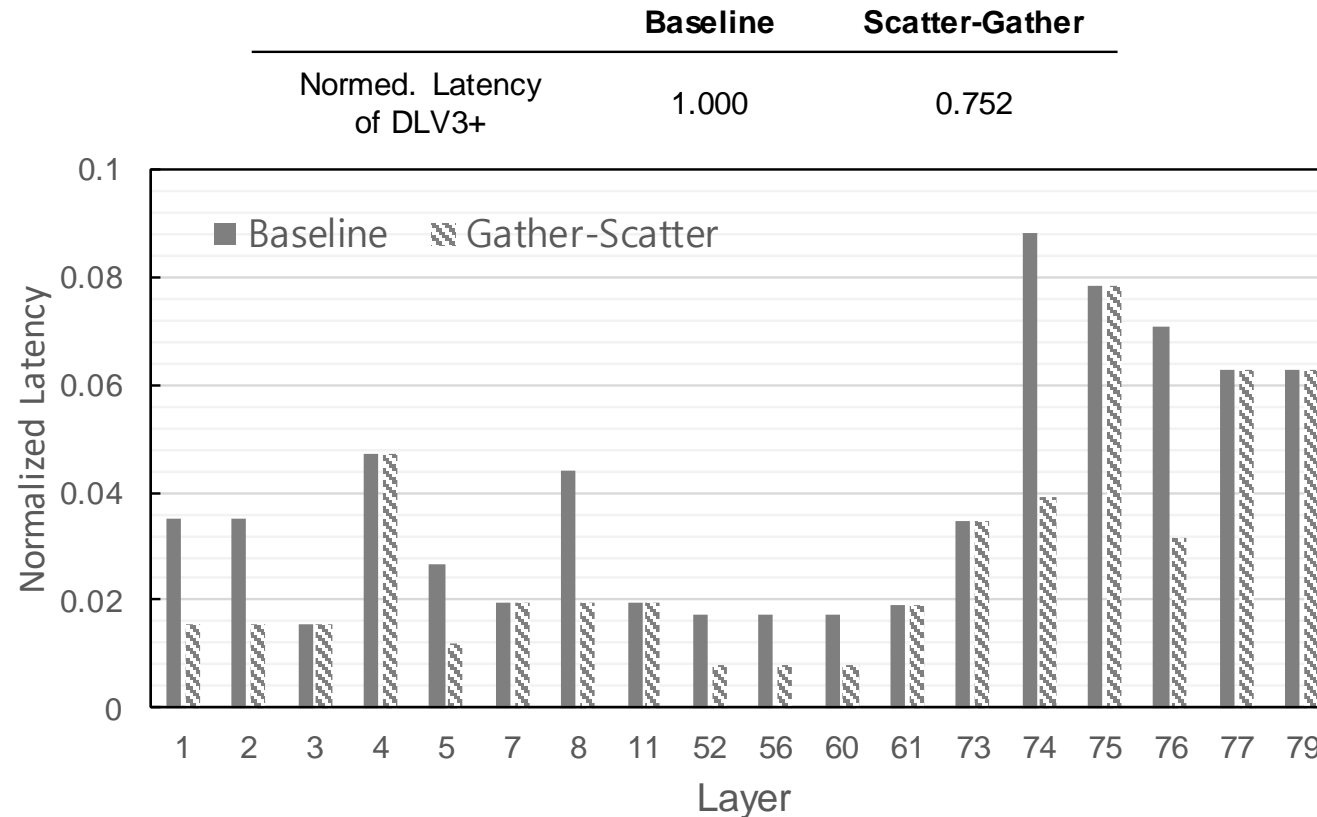
(b) Tensor Engine with Scatter-Gather for shallow input channel

# Latency enhancement by Scatter-Gather

| | Baseline | Scatter-Gather |
|---|---|---|
| Normed. Latency of DLV3+ | 1.000 | 0.752 |



☐SG reads 4 consecutive input vectors and stacks them into a vector.

- ■ SG enhances utilization up to 4 times.
- ■ Compute utilization is improved from 37% to 50% for DeepLabV3 (FP16)

# NPU Enhancement: Power Efficiency



ALU (HW) Utilization

Compute Density

Bit precisions

Clock/Power gating

Attainable OPS

Due to Power/Thermal

Improved energy efficiency

Limitation

Compute intensity

# Enhance Effective Energy Efficiency

☐ Quantization: Power vs. Accuracy Trade off

**Formats for mixed precisions**







☐ MAC datapath optimization



**(b) Accumulator-based dotproduct engine**



**(a) Adder-tree based dotproduct engine**

# Outline

- ☐ Introduction
  - ■ Challenges in Designing NPU
- ☐ NPU architecture overview
- ☐ Evolution of NPU over the years
  - ■ Compute density
  - ■ HW utilization improvement
  - ■ Power efficiency & heterogeneous NPU
- ☐ **Challenges in Evaluating NPU**
  - ■ **Domain-specific architectures (Heterogeneous NPUs, CIM/PIM)**
  - ■ **Fast-Evolving Neural Networks**
  - ■ **Rapid scaling-up of NPUs outpacing CPUs or DRAMs**
  - ■ **Expoliting HW parallelism (Tensor-Vector, Memory pre-fetch)**
- ☐ Conclusion

Challenges in Designing and Evaluating Neural Processing Units

# Multiple-generations of NPU HW



[3]

| Process | 8nm CMOS technology (Samsung) |
|---|---|
| Area | 5.5mm$^2$ |
| Voltage | 0.5-to-0.8V |
| Frequency | 67-to-933-MHz |
| Best Peak Performance | 1910/6937* TOPS (* 75% weight zeros) |
| Best Energy Efficiency | 11.5 (8b) @ 0.5V |

[4]

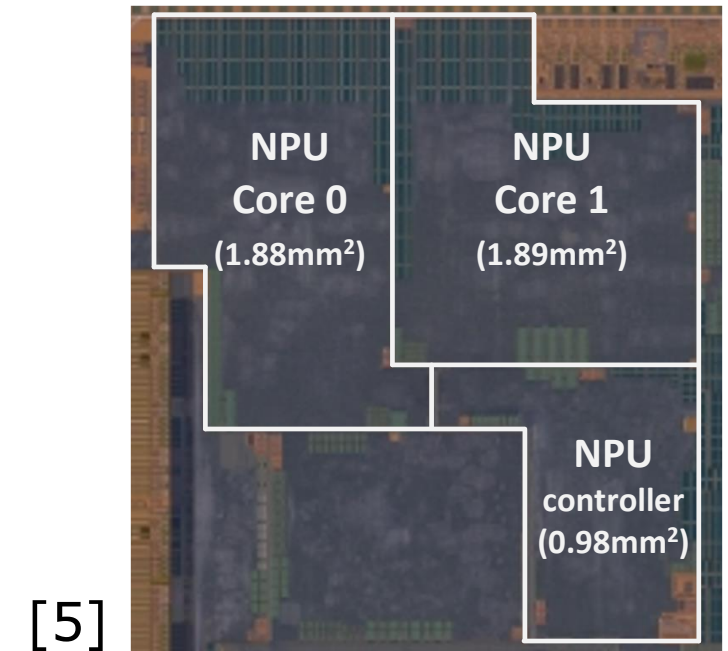| Process | 5nm CMOS technology (Samsung) |
|---|---|
| Area | 5.46mm$^2$ |
| Voltage | 0.55-to-0.9V |
| Frequency | 332-to-1196-MHz |
| Best Peak Performance | **623 inferences/s** @ 0.9V (Inception V3) |
| Best Energy Efficiency | **13.6 TOPS/W** @ 0.6V (Inception V3) |

[5]

| Process | 4nm CMOS technology (Samsung) |
|---|---|
| Area | 4.74mm$^2$ |
| Voltage | 0.55-to-1.0V |
| Frequency | 332-to-1196-MHz |
| Best Peak Performance | **3433 inferences/s** @ 1.0V (MobileNetTPU_INT8) |
| Best Energy Efficiency | **11.59 TOPS/W** @ 0.55V (MobileNetTPU_INT8) |

Challenges in Designing and Evaluating Neural Processing Units

# Examples of Physical Limitations: Temp.

□ Thermal management

- Throttling: Intentional performance reduction for temp. control
- Tripping: forced shutdown to protect SOC

No more than 40℃ not to cause skin burn

[2] DT Team, "Design Methodologies for Advanced Mobile SOCs", Samsung Foundry

# Fast-Evolving Neural Networks

☐ Neural Networks: Fast-Paced Evolution and the Challenge of a Moving Target



ZFNet, Dropout — 13

ResNet, U-net — 15

Transformer, mobileNet, DeepLab — 17

EfficientNet, YOLO v3 — 19.5

ChatGPT, AlphaFold, StableDiffusion — 21.5

AlexNet — 12

GAN, GRU, GoogLeNet — 14

AlphaGo, WaveNet — 16

GPT, BERT — 18

ViT, NeRF — 20

DALL-E, Swin Trans. — 22

☐ The computational configuration greatly varies among different neural networks



[1] JaewanChoi et al, "Accelerating Transformer Networks through Recomposing SoftmaxLayers", Arxiv, 2023.

# Computational Characteristics of NNs



☐ Divided into two Groups in terms of parameters # and operation speed

# Compute-in-Memory

☐ Energy efficient CIM cells



A 28nm 16.9-300TOPS/W Computing-in-Memory Processor Supporting Floating-Point NN Inference/Training with Intensive-CIM Sparse-Digital Architecture

☐ System-level benefits

# Specialization and Heterogeneous NPUs

# Rapid scaling-up of NPU outpacing CPU/DRAM



□ SW and DMA overheads such as data transfer and host time take up a big portion of the total latency in the small size of neural networks.

# Exploiting HW parallelism

☐ It is required scheduling with HW components in NPU to exploit the HW parallelism as considering the load balancing between them.

- ■ The components inside the NPU become more diverse
- ■ There can be more types of NPU engines due to domain-specific architecture.

# Challenges in Evaluation

☐ Various optimization points considered at the design stage influence the performance of the NPU. HW-SW Co-optimization is also important.

- ■ SW
  - ☐ Precision
  - ☐ Tiling & scheduling techniques
  - ☐ Thermal/Power Management
- ■ Architecture
  - ☐ Domain-specific architecture
  - ☐ CIM/PIM

- ■ SoC
  - ☐ Technologies
  - ☐ Memory Hierarchy ( + System Bus, LLC)
  - ☐ CPU time to set up the application

☐ Benchmarks & Guidelines for a fair comparison

- ■ V. J. Reddi et al. "MLPerf; Mobile Inference Benchmark: An Industry-Standard Open-Source Machine Learning Benchmark for On-Device AI", MLSys, 2022.
- ■ Geoffrey Burr et al., "Fair and Comprehensive Benchmarking of Machine Learning Processing Chips", IEEE Design & Test, 2022

# Outline

- ☐ Introduction
  - ◾ Challenges in Designing NPU
- ☐ NPU architecture overview
- ☐ Evolution of NPU over the years
  - ◾ Compute density
  - ◾ HW utilization improvement
  - ◾ Power efficiency & heterogeneous NPU
- ☐ Challenges in Evaluating NPU
  - ◾ Domain-specific architectures (Heterogeneous NPUs, CIM/PIM)
  - ◾ Fast-Evolving Neural Networks
  - ◾ Rapid scaling-up of NPUs outpacing CPUs or DRAMs
  - ◾ Expoliting HW parallelism (Tensor-Vector, Memory pre-fetch)
- ☐ **Conclusion**

# Conclusion

☐ NPUs in a mobile AP are enabling efficient NN inference

☐ NPUs keep scaling the performance by enhancing the HW utilization and putting more MACs

☐ NPUs face physical limitations which make many kinds of design challenges

☐ Various optimization points considered at the design stage influence the performance of the NPU such as HW-SW co-optimization

☐ We need a methodology that not only evaluates the NPU but also encompasses the evaluation of the entire system

# Reference

- [1] Jaewan Choi et al, "Accelerating Transformer Networks through Recomposing SoftmaxLayers", Arxiv, 2023.
- [2] DT Team, "Design Methodologies for Advanced Mobile SOCs", Samsung Foundry.
- [3] J Song et al., "An 11 5TOPS/W1024 MAC Butterfly Structure Dual Core Sparsity-Aware Neural Processing Unit in 8nm Flagship Mobile SoC", ISSCC, 2019.
- [4] J-S Park et al., "A 6K-MAC Feature-map-sparsity-aware Neural Processing Unit in 5nm Flagship Mobile SoC", ISSCC, 2021.
- [5] J-S Park et al., "A Multi-Mode 8K-MAC HW-Utilization-Aware Neural Processing Unit with a Unified Multi-Precision Datapath in 4nm Flagship Mobile SoC", ISSCC, 2022.
- [6] JW Jang, "Sparsity-aware and re-configurable NPU architecture for Samsung flagship mobile SoC", ISCA, 2021.
- [7] S Williams, "Roofline: an insightful visual performance model for multicore architectures", Communications of the ACM, 2009.
- [8] Colby R. Bnabury et al., "Benchmarking Tiny ML Systems", arxiv, 2021.
- [9] TingxingTim Dong et al., "RenderSR: A lightweight Super-Resolution Model for Mobile Gaming Upscaling", CVPR, 2022.
- [10] Eli Schwartz et al., "DeepISP: Learning End-to-End Image Processing Pipeline", ICLR 2018.

# Reference

- [11] Tobias Boceck et al., "Force Touch Detection on Capacitive Sensors using Deep Neural Networks", MobileHCI, 2019.
- [12] Jialing Li et al., "Deep Learning-based Massive CSI Feedback", ICOCN, 2019.
- [13] V. J. Reddi, et al. "MLPerf mobile inference benchmark: An industry-standard open-source machine learning benchmark for on-device AI." Proceedings of Machine Learning and Systems, 2023.
- [14] C-H. Lin et al., "A 3.4-to-13.3TOPS/W 3.6TOPS Dual-Core Deep-Learning Accelerator for Versatile AI Applications in 7nm 5G Smartphone SoC," ISSCC, 2020.
- [15] Y. Jiao et al, "A 12nm Programmable Convolution-Efficient Neural-Processing-Unit Chip
- Achieving 825TOPS," ISSCC, 2020.
- [16] A. Agrawal et al., "7nm 4-Core AI Chip with 25.6TFLOPS Hybrid FP8 Training, 102.4TOPS INT4 Inference and Workload-Aware Throttling," ISSCC, 2021.
- [17] J-S Park et al., "A Multi-Mode 8K-MAC HW-Utilization-Aware Neural Processing Unit with a Unified Multi-Precision Datapath in 4nm Flagship Mobile SoC", JSCC, 2022.

# Thank you for your attention