

Efficient AI Modeling

John Li, Jacob Lin, and Ray Kuo
Digital Center and AI Center

創新

Innovation

品質

Quality

虛心

Open Mind

力行

Execution

Inventec



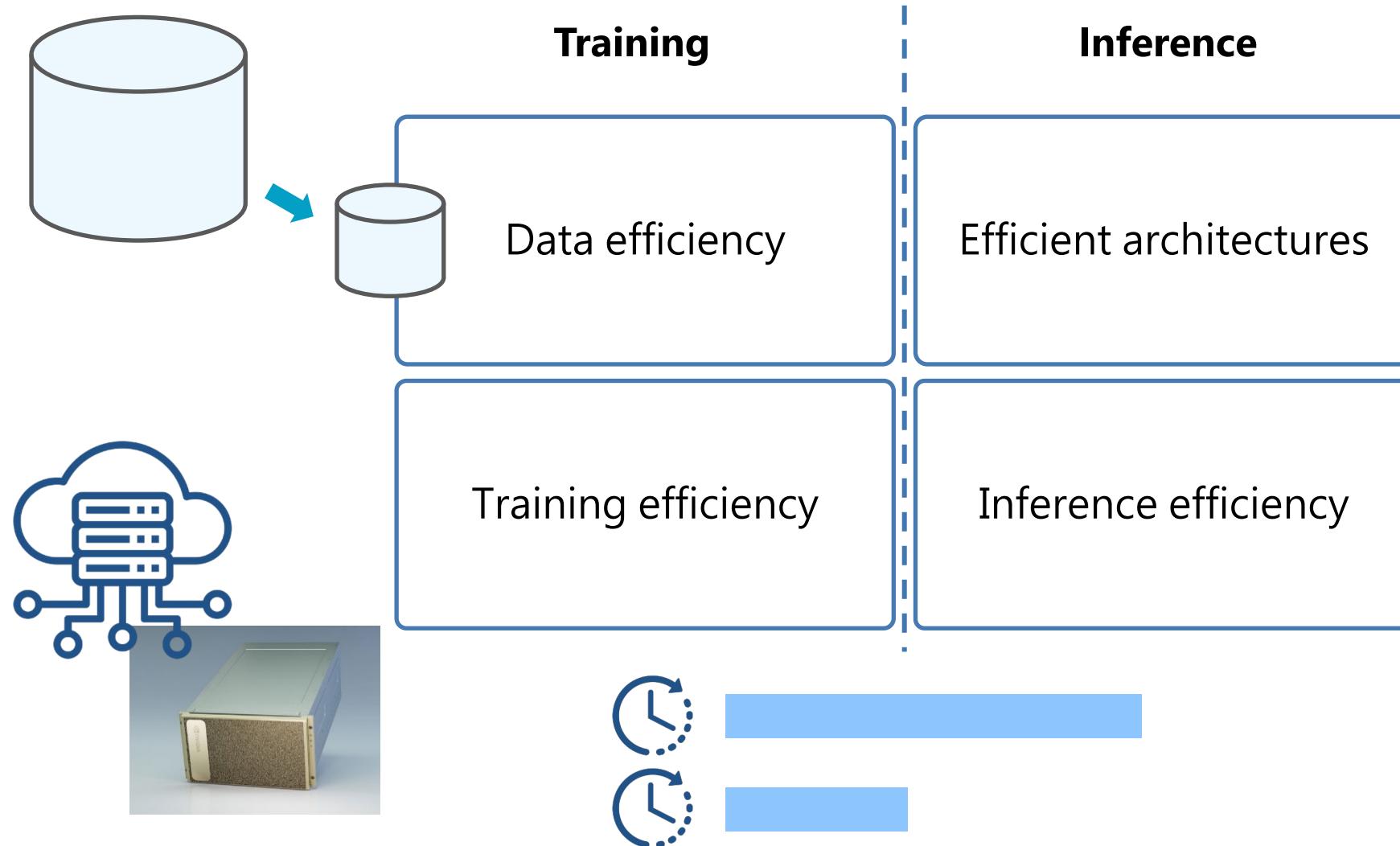
Outline

- Overview
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

Outline

- **Overview**
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

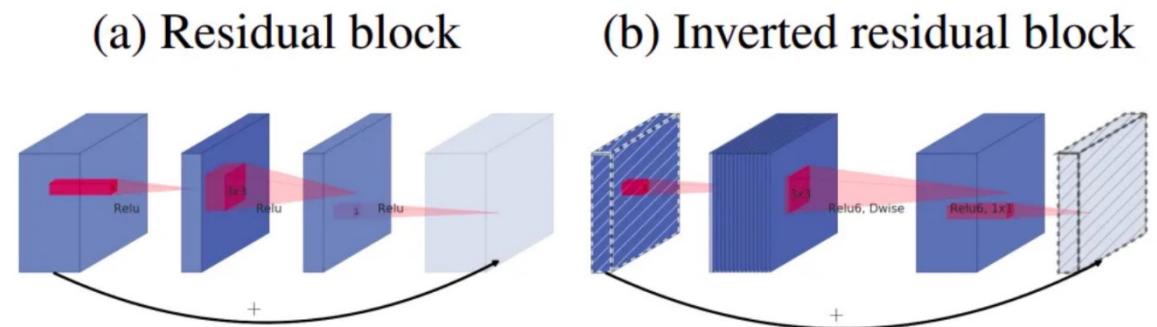
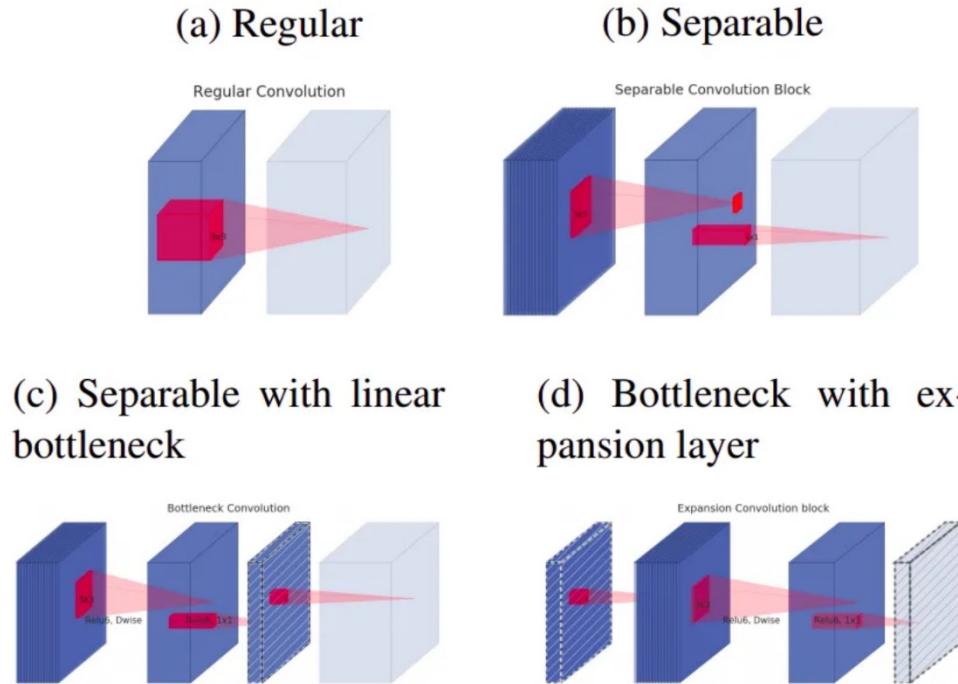
Efficient AI



Model Design

MobileNetV2 and MobileNetV3

- V2: Bottleneck Residual Block
- V3: Squeeze-and-Excitation, NAS and NetAdapt



Sandler, Mark, et al. "MobileNetv2: Inverted residuals and linear bottlenecks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.

Slide credit: John

Inventec

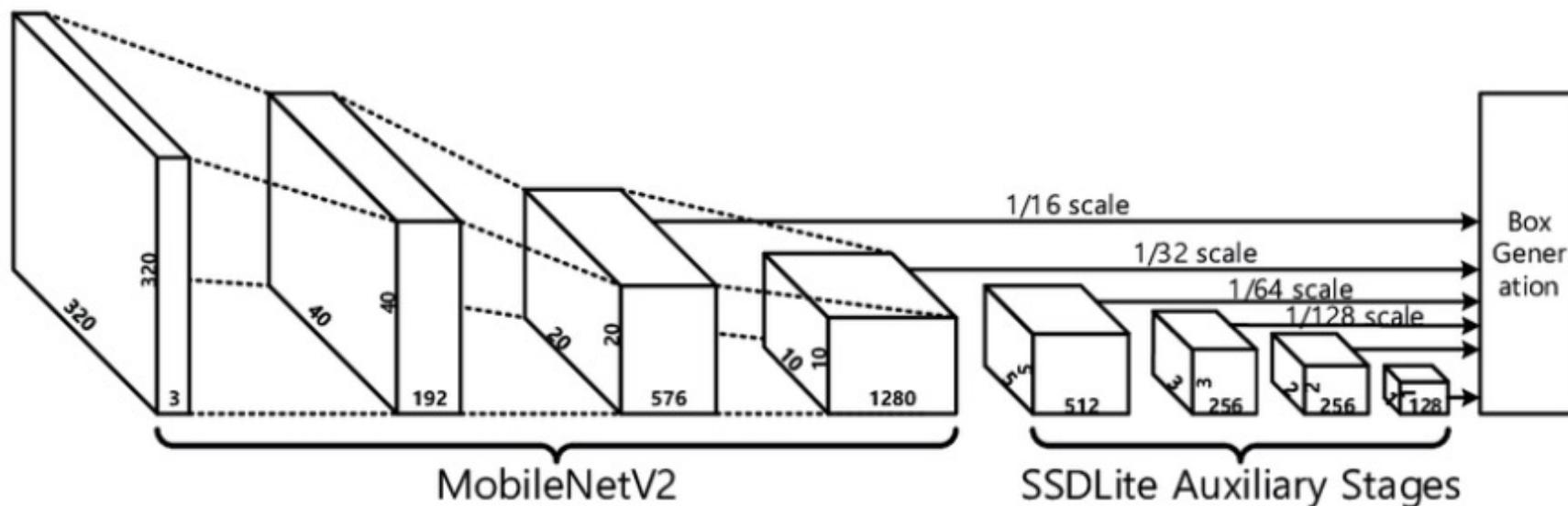
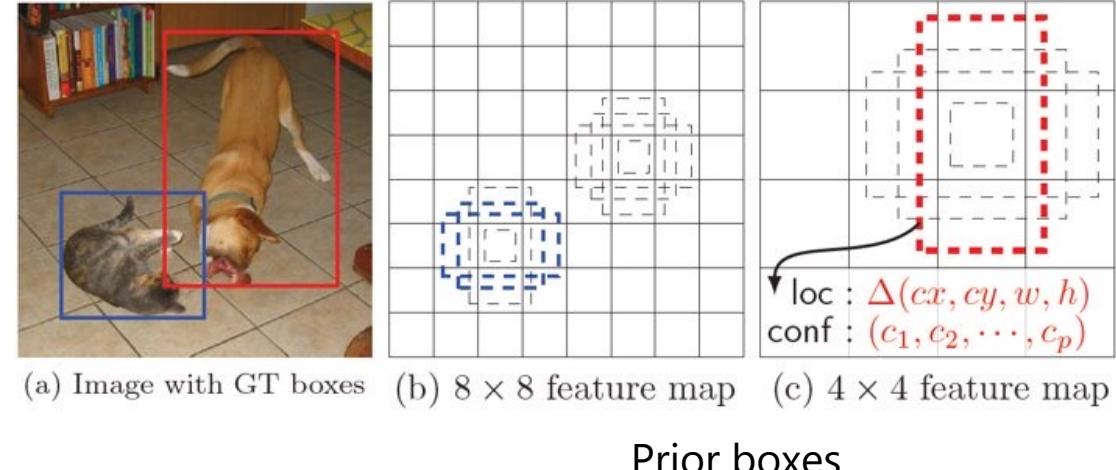
Inventec Confidential



Model Design

Face detection

- Network: MobileNet + SSD (Single Shot Detector)



<https://github.com/Linzaer/Ultra-Light-Fast-Generic-Face-Detector-1MB>

Heeong-Ju Kang, "SSDLiteX: Enhancing SSDLite for Small Object Detection," *Applied Sciences*, 2023

Slide credit: Jacob

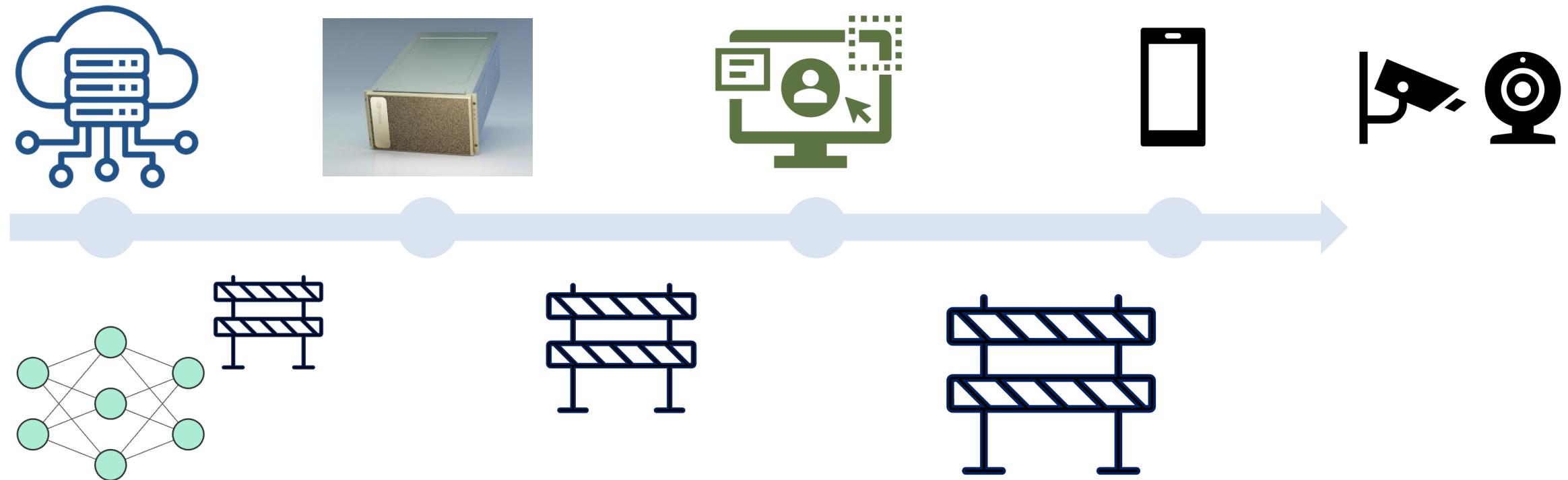


Inventec Confidential

Gap of Efficiency Expectation

New challenges

- Edge-of-the-edge

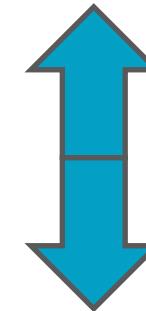


Efficient AI

AI on the edge

- Define the use case
- Select the right hardware
- Develop the AI model
- Optimize the model for edge deployment
- Test the system
- Deploy the system to production

Face Detection



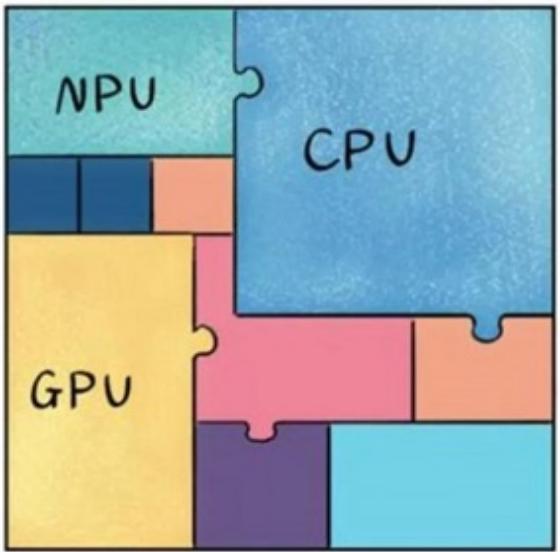
AI on Chips

Outline

- Overview
- **AI on Chip (AloC) team work flow**
- Quantization
- Testing System
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

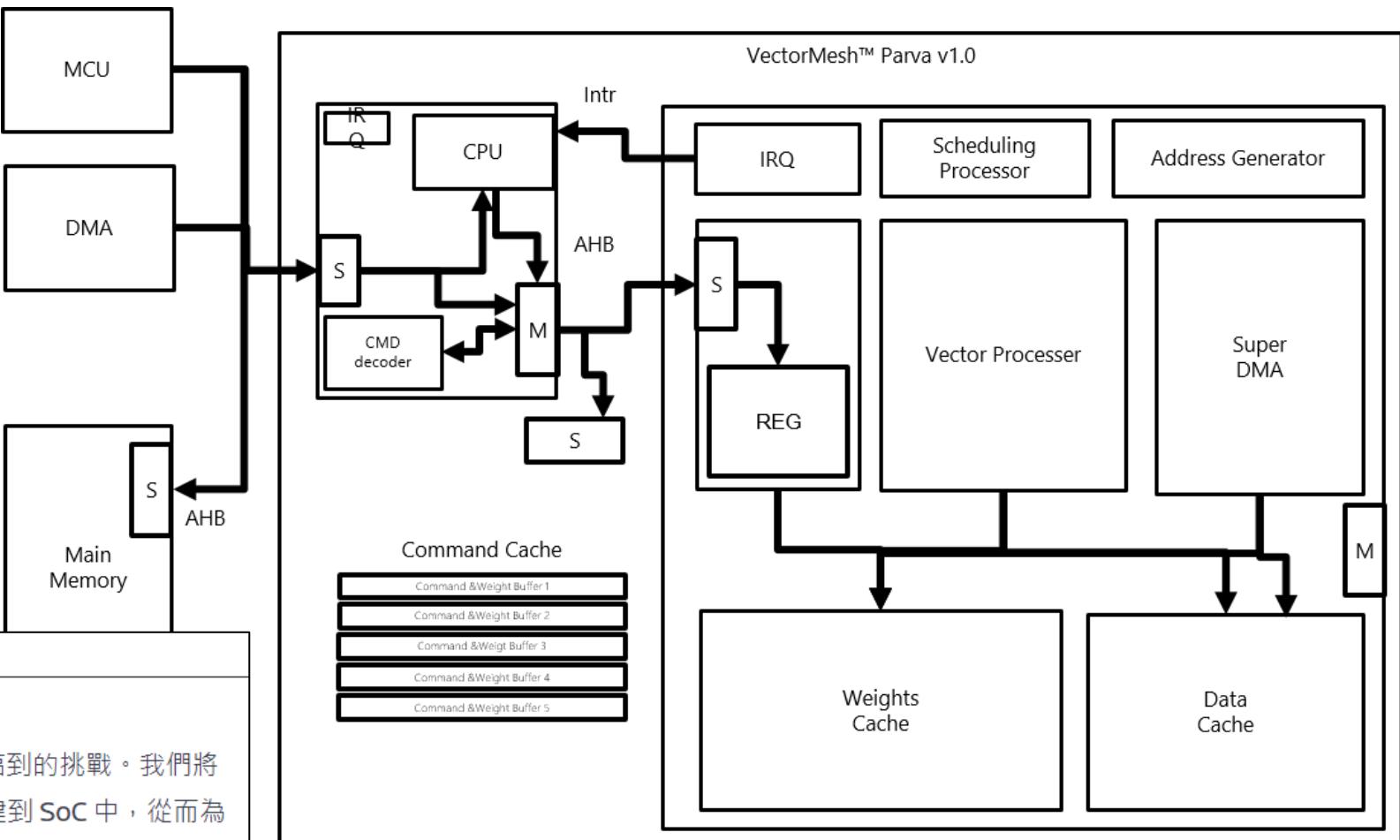
What is NPU IP?

System on the chip (SoC)



聯發科技 NeuroPilot

我們以聯發科技 NeuroPilot 對應在終端人工智慧所面臨到的挑戰。我們將 CPU、GPU 和 APU (AI 處理單元) 等異構運算功能內建到 SoC 中，從而為人工智慧功能和應用提供了高性能和低功耗。針對 SoC 中的這些特定處理單元，開發人員可以讓聯發科技 NeuroPoint SDK 智慧地為他們處理所分配到的任務。



<https://www MEDIATEK.tw/innovations/artificial-intelligence>

Slide credit: Jacob

Inventec Confidential

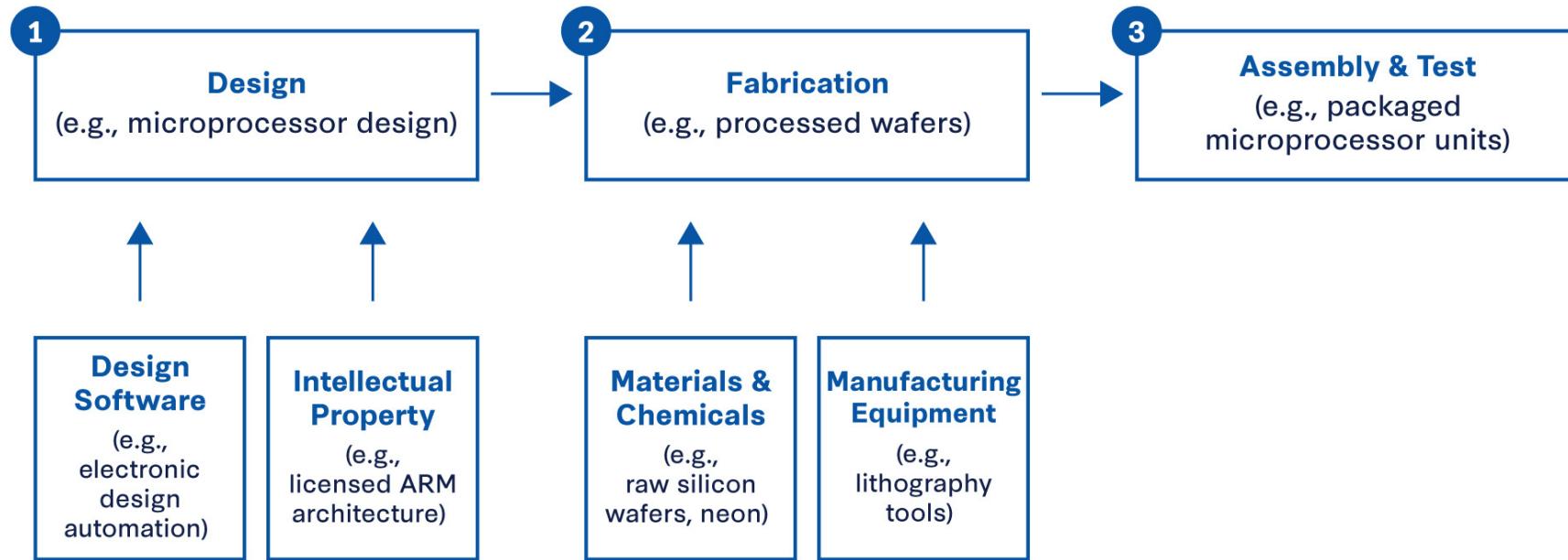


Why NPU IP?

System on the chip (SoC)

- SIP: Silicon Intellectual Property (矽智財)
- Faster, Better, Cheaper

Simplified Depiction of the Semiconductor Value Chain



<https://www.csis.org/analysis/mapping-semiconductor-supply-chain-critical-role-indo-pacific-region>

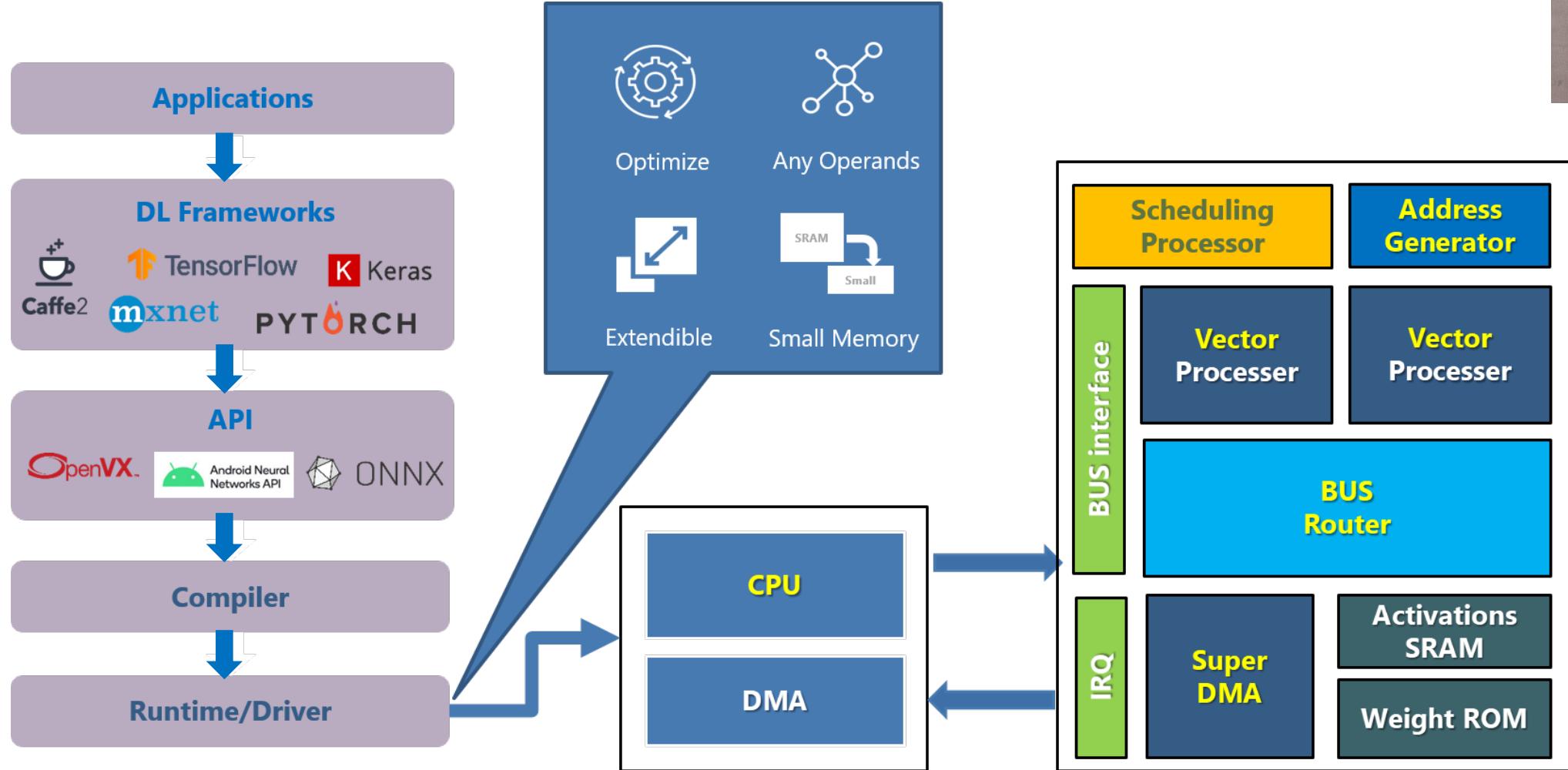
Slide credit: Jacob



Inventec Confidential

One-stop Service?

From AI Model Training to NPU IP Design



Opposite: Separate solutions for
AI model, compiler, NPU
design, ... etc



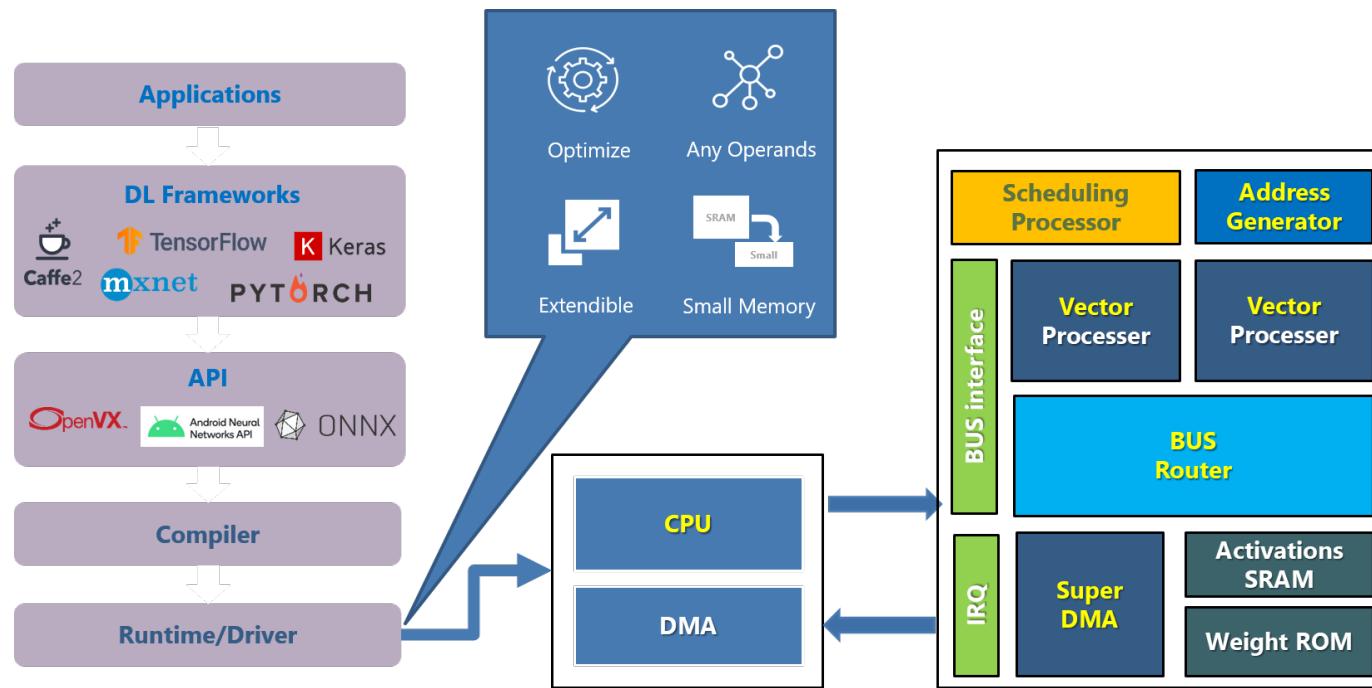
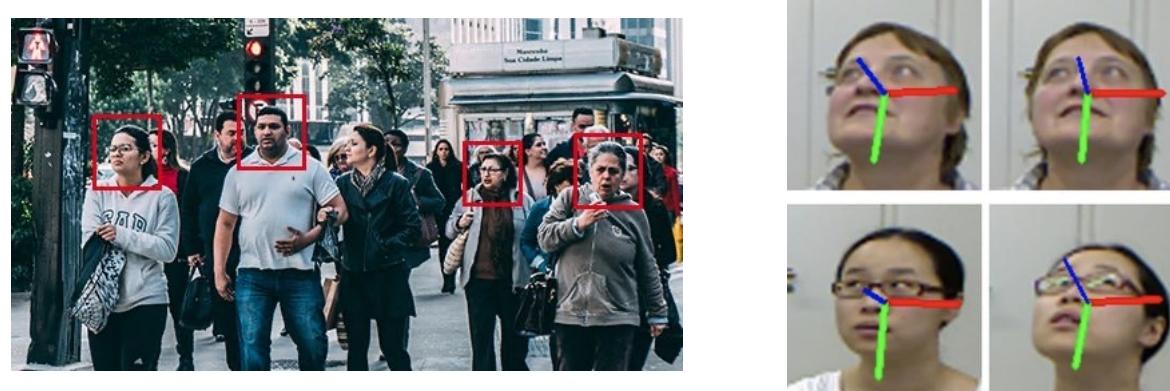
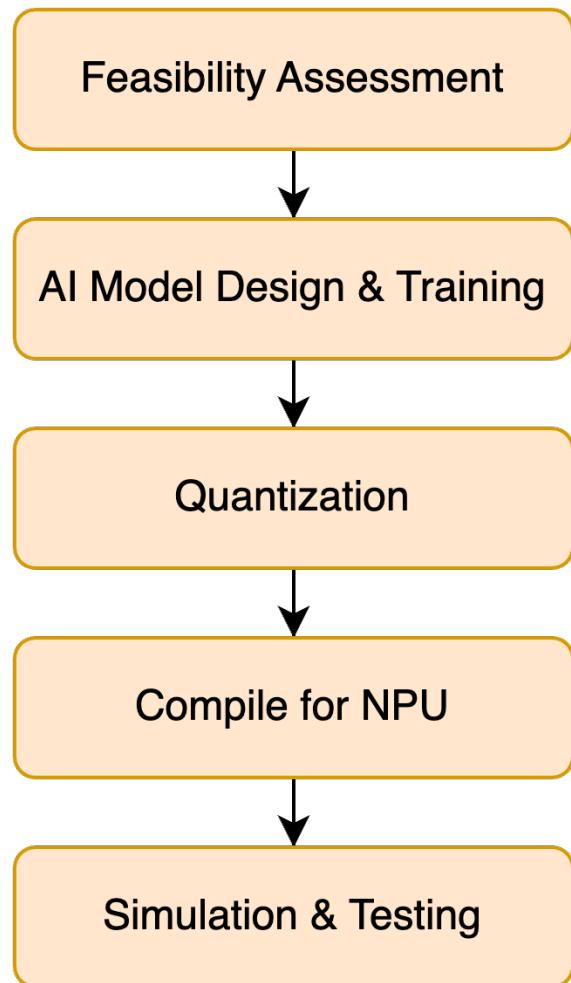
<https://www.quwenxuei.com/classical/21785.html>

Slide credit: Jacob

Inventec Confidential



AIoC Team Workflow

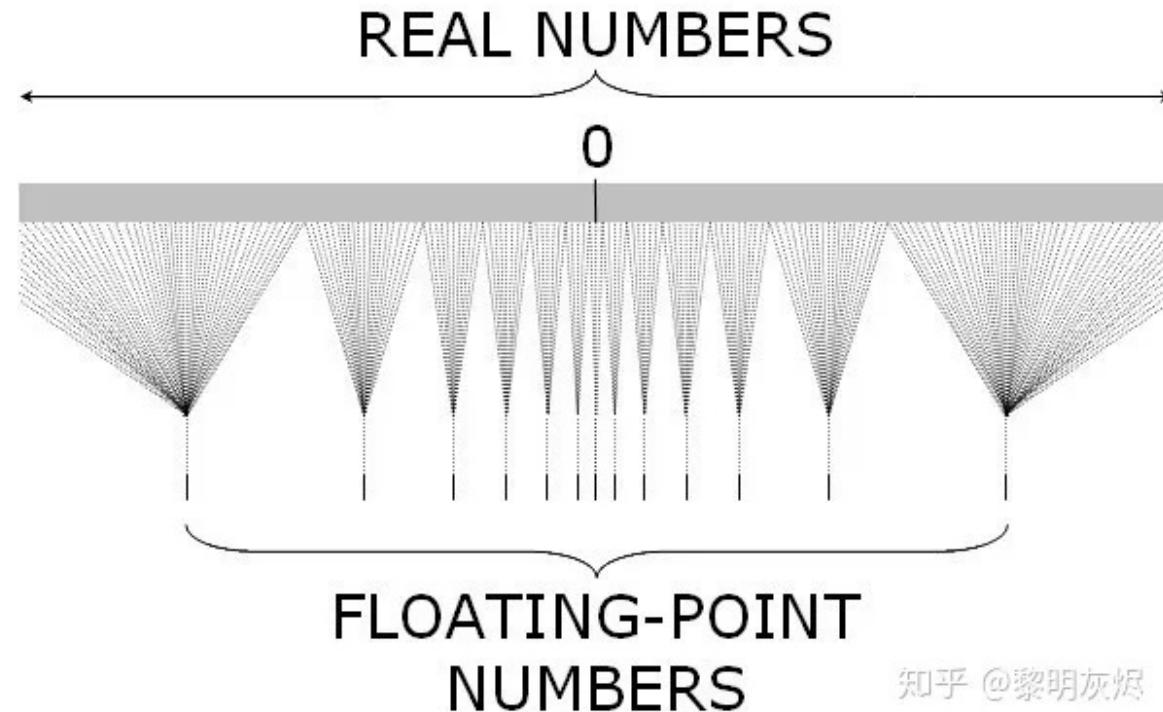
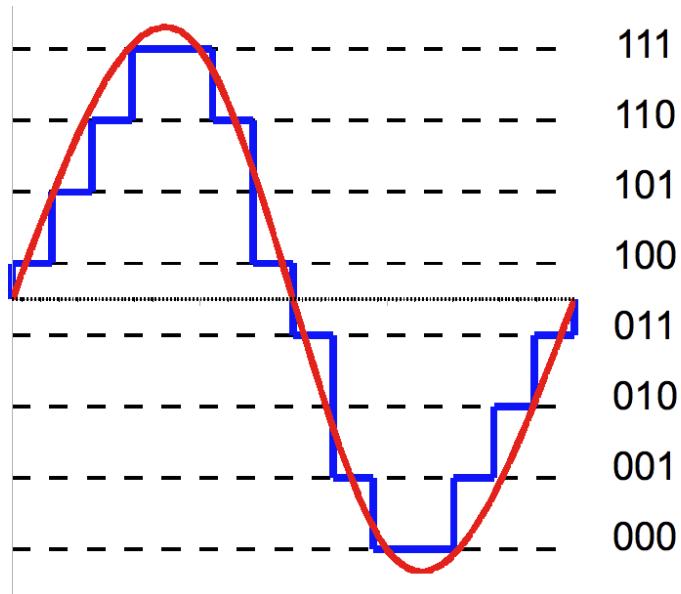


Outline

- Overview
- AI on Chip (AloC) team work flow
- **Quantization**
- Testing System
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

What is Quantization?

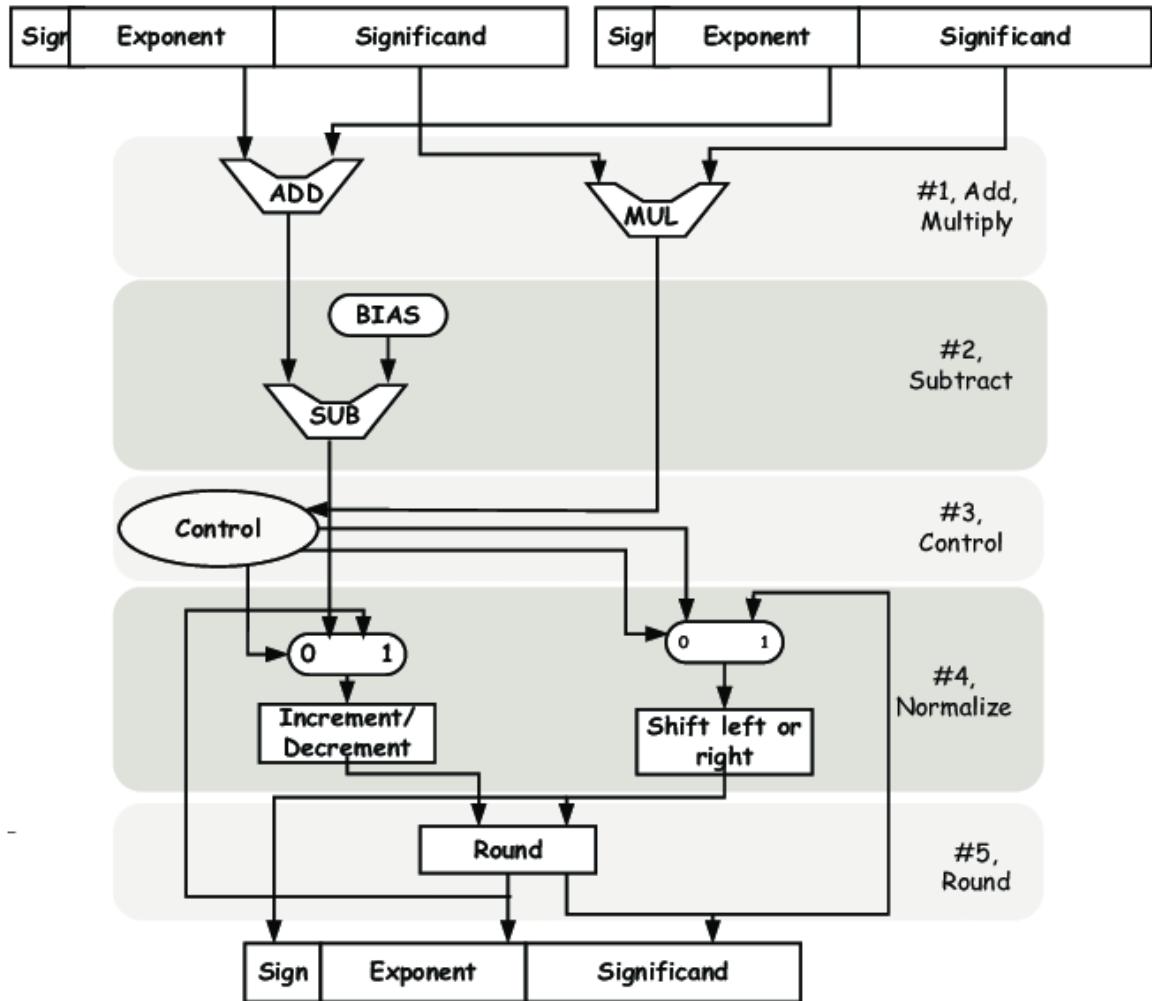
- Process of turning continuous values into a finite set of discrete values.
- Computer Science & DSP: Used to reduce data precision.
 - Resource Efficiency
 - Faster Computation
 - **Hardware Implementation**



知乎 @黎明灰烬

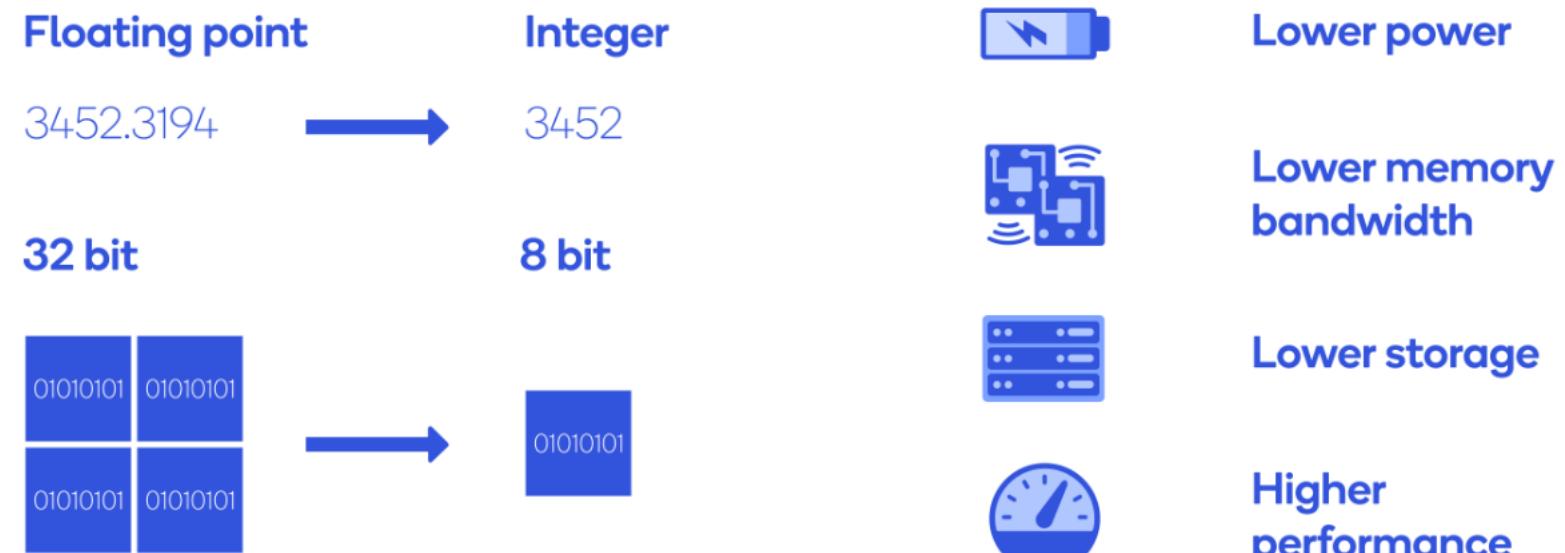
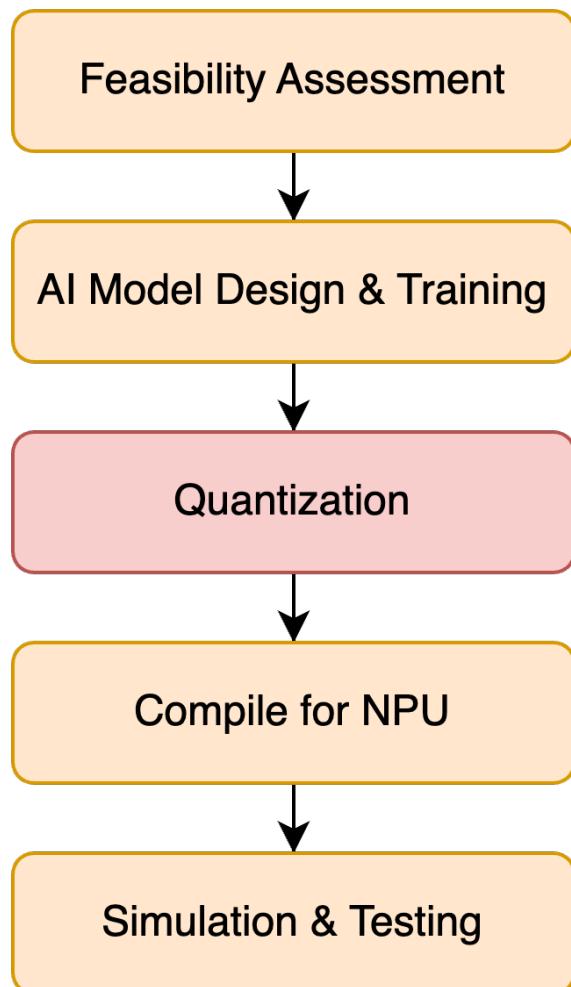
Hardware Implementation

Floating point vs fixed point



https://www.researchgate.net/figure/Floating-Point-Multiplier-Architecture_fig2_224698133

Quantization for Neural Network



<https://www.qualcomm.com/news/onq/2019/03/heres-why-quantization-matters-ai>

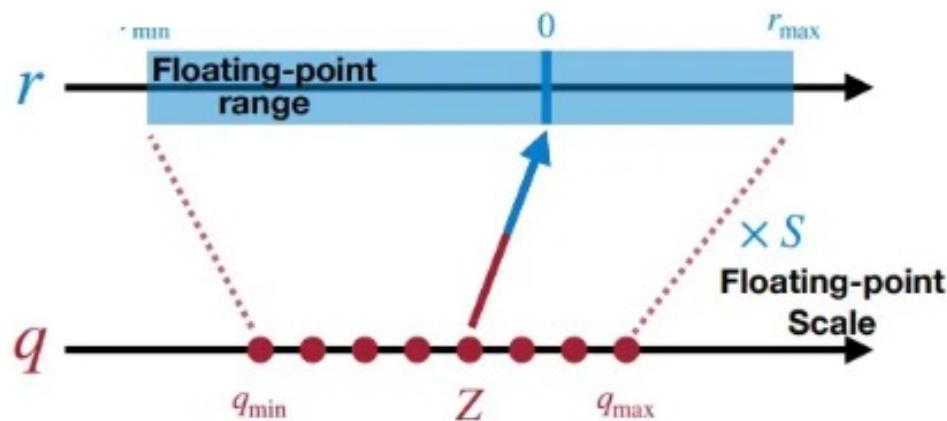
Quantization for Neural Network

Asymmetric vs. Symmetric

- r: floating point
- q: fixed point after quantization
- S: scale
- Z: zero point

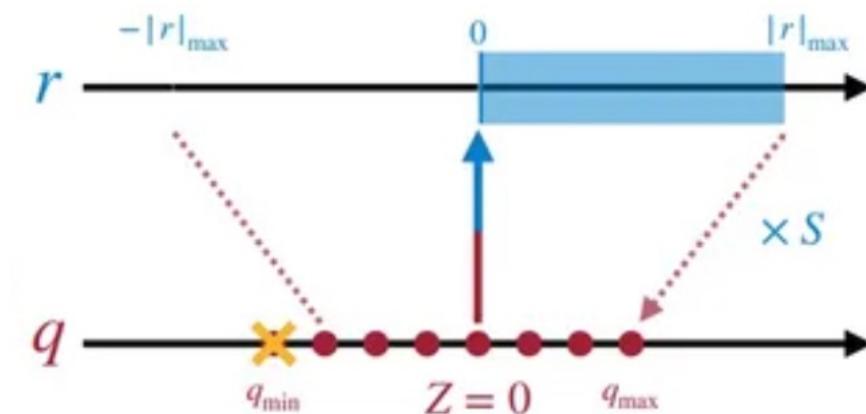
$$S = \frac{r_{max} - r_{min}}{q_{max} - q_{min}}$$

$$Z = round(q_{max} - \frac{r_{max}}{S})$$



$$r = S(q - Z)$$

$$q = round\left(\frac{r}{S} + Z\right)$$



Quantization for Neural Network

Quantization of matrix operations

- 2D Convolution can be converted to matrix multiplication
- $\mathbf{R}_3 = \mathbf{R}_1 \times \mathbf{R}_2$

$$r = S(q - Z)$$

$$r_3^{i,k} = \sum_{j=1}^N r_1^{i,j} r_2^{j,k}$$

$$S_3(q_3^{i,k} - Z_3) = \sum_{j=1}^N S_1(q_1^{i,j} - Z_1) S_2(q_2^{j,k} - Z_2)$$

$$q_3^{i,k} = \frac{S_1 S_2}{S_3} \sum_{j=1}^N (q_1^{i,j} - Z_1)(q_2^{j,k} - Z_2) + Z_3$$

Convolution
↓
Matrix Multiply

$$\text{Filter} \quad \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \text{Input fmap} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} = \text{Output fmap} \quad \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\text{Toeplitz Matrix (w/redundant data)} \quad \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 4 & 5 \\ 2 & 3 & 5 & 6 \\ 4 & 5 & 7 & 8 \\ 5 & 6 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$



Quantization for Neural Network

- M is usually between 0 and 1 (statistical results)

$$q_3^{i,k} = \boxed{\frac{S_1 S_2}{S_3}} \sum_{j=1}^N (q_1^{i,j} - Z_1)(q_2^{j,k} - Z_2) + Z_3$$

$$\boxed{M} = \frac{S_1 S_2}{S_3}$$

$$q_3^{i,k} = M \sum_{j=1}^N (q_1^{i,j} - Z_1)(q_2^{j,k} - Z_2) + Z_3 = \boxed{M} P + Z_3$$

$$\boxed{M} = 2^{-n} M_0$$

$$q_3^{i,k} = 2^{-n} M_0 P + Z_3$$

Convolution



Matrix Multiply

$$r = S(q - Z)$$

$$q = \text{round}\left(\frac{r}{S} + Z\right)$$

Filter	$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$	*	Input fmap	$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$	=	$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$
--------	--	---	------------	---	---	--

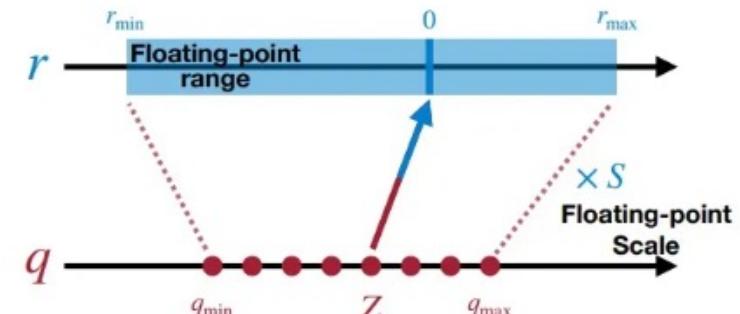
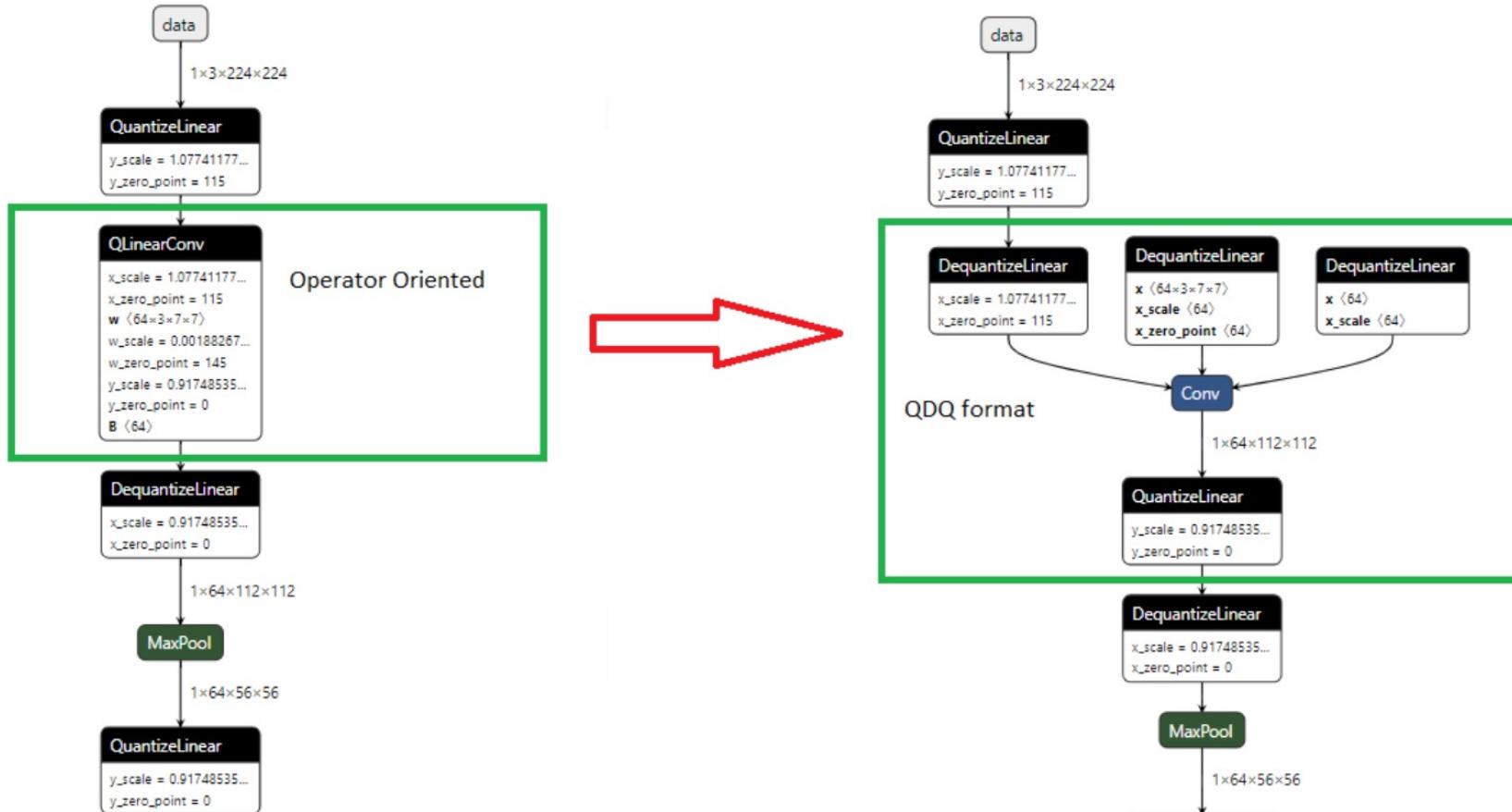
Toeplitz Matrix
(w/redundant data)

$\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$	\times	$\begin{bmatrix} 1 & 2 & 4 & 5 \\ 2 & 3 & 5 & 6 \\ 4 & 5 & 7 & 8 \\ 5 & 6 & 8 & 9 \end{bmatrix}$	=	$\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$
---	----------	--	---	---



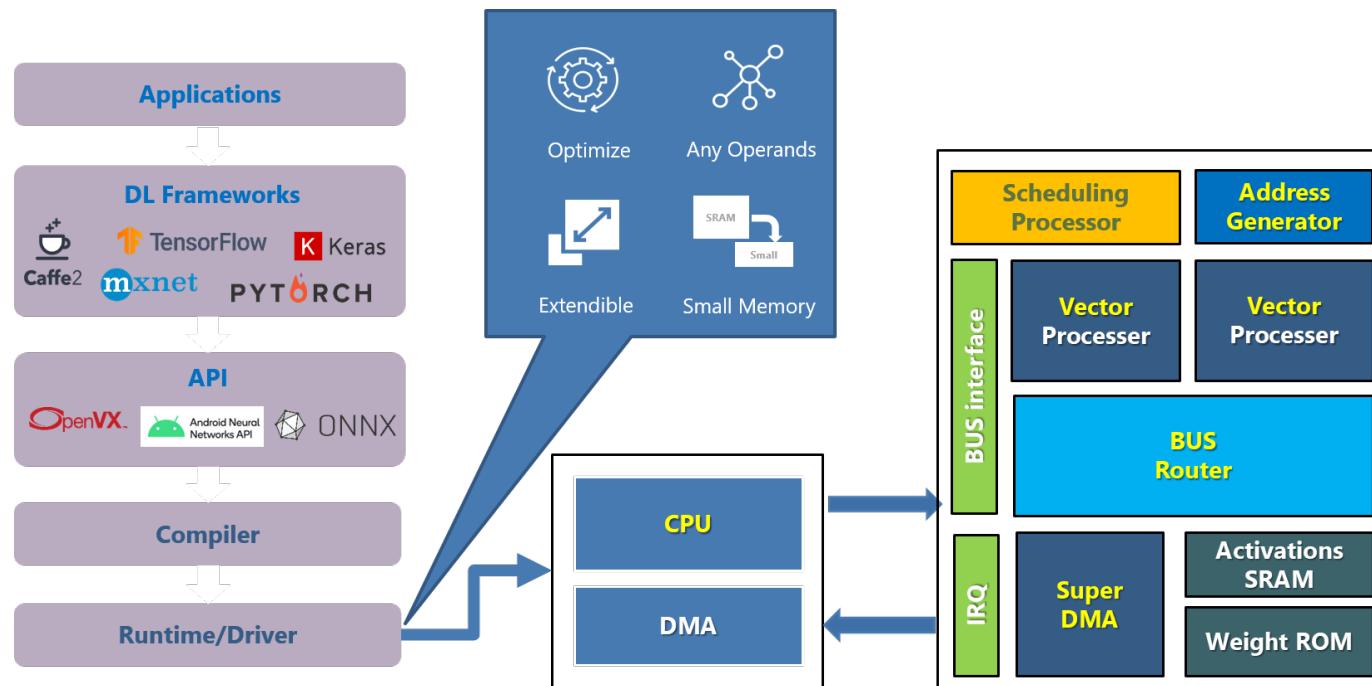
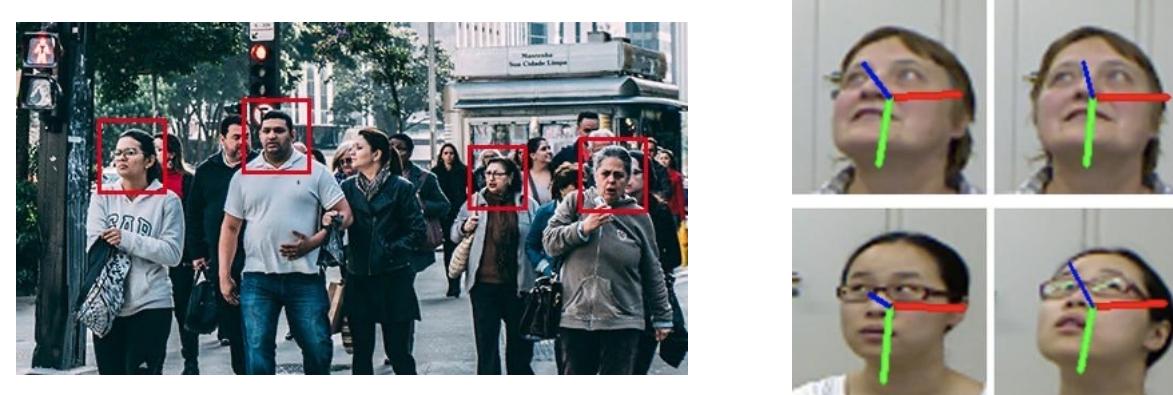
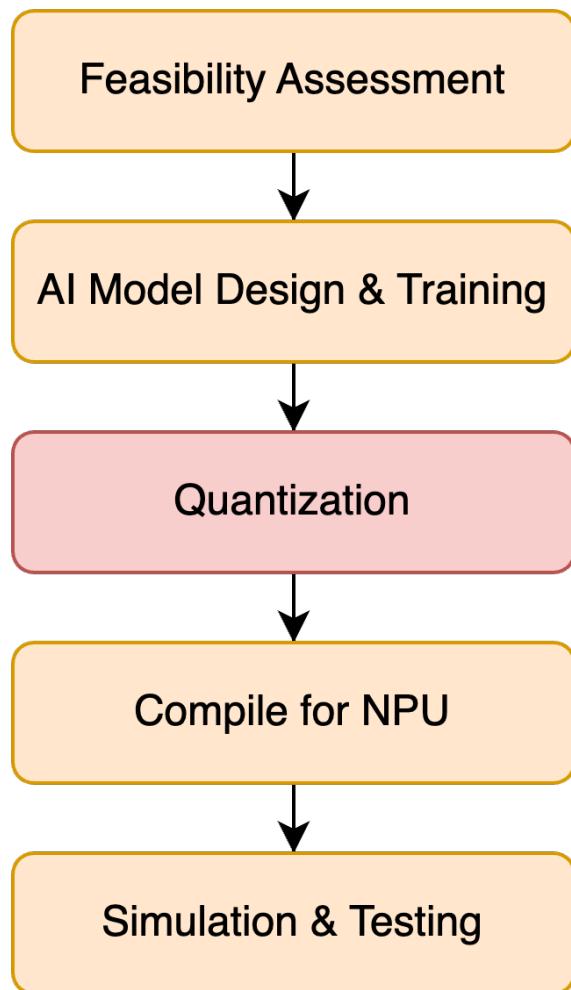
Onnx Quantization

- Convert floating point model to QDQ format and perform **calibration**

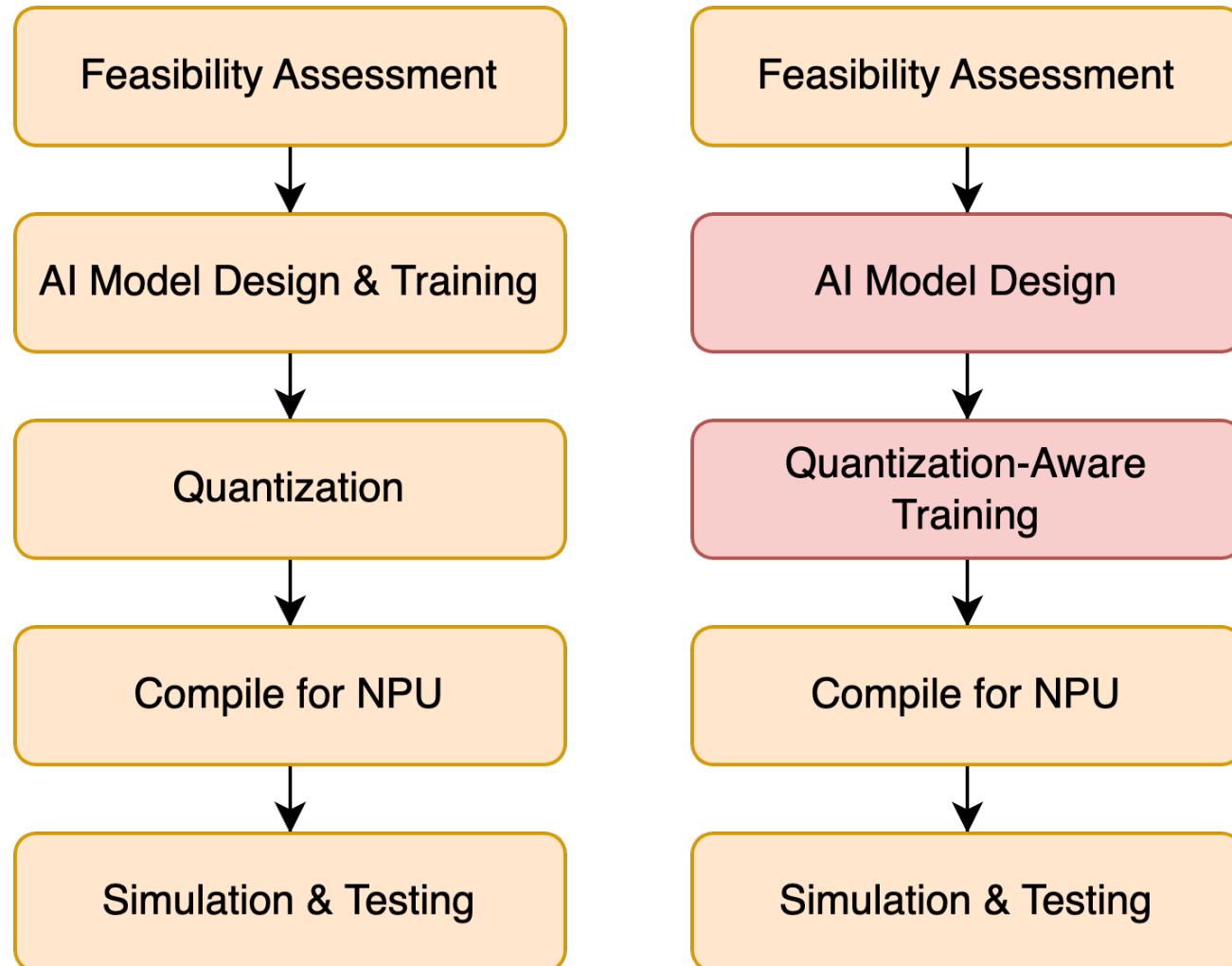


<https://onnxruntime.ai/docs/performance/model-optimizations/quantization.html>

Quantization



Quantization Aware Training



Outline

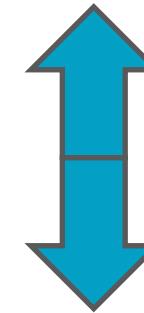
- Overview
- AI on Chip (AloC) team work flow
- Quantization
- **Testing System**
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

Efficient AI

AI on the edge

- Define the use case
- Select the right hardware
- Develop the AI model
- Optimize the model for edge deployment
- Test the system
- Deploy the system to production

Face Detection



AI on Chips

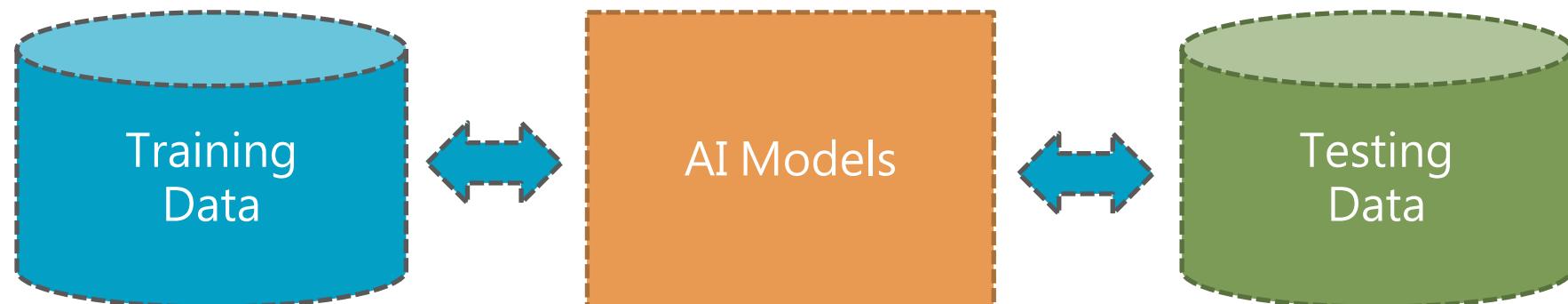
What is Learning in AI?

A program that writes itself... using framework and data

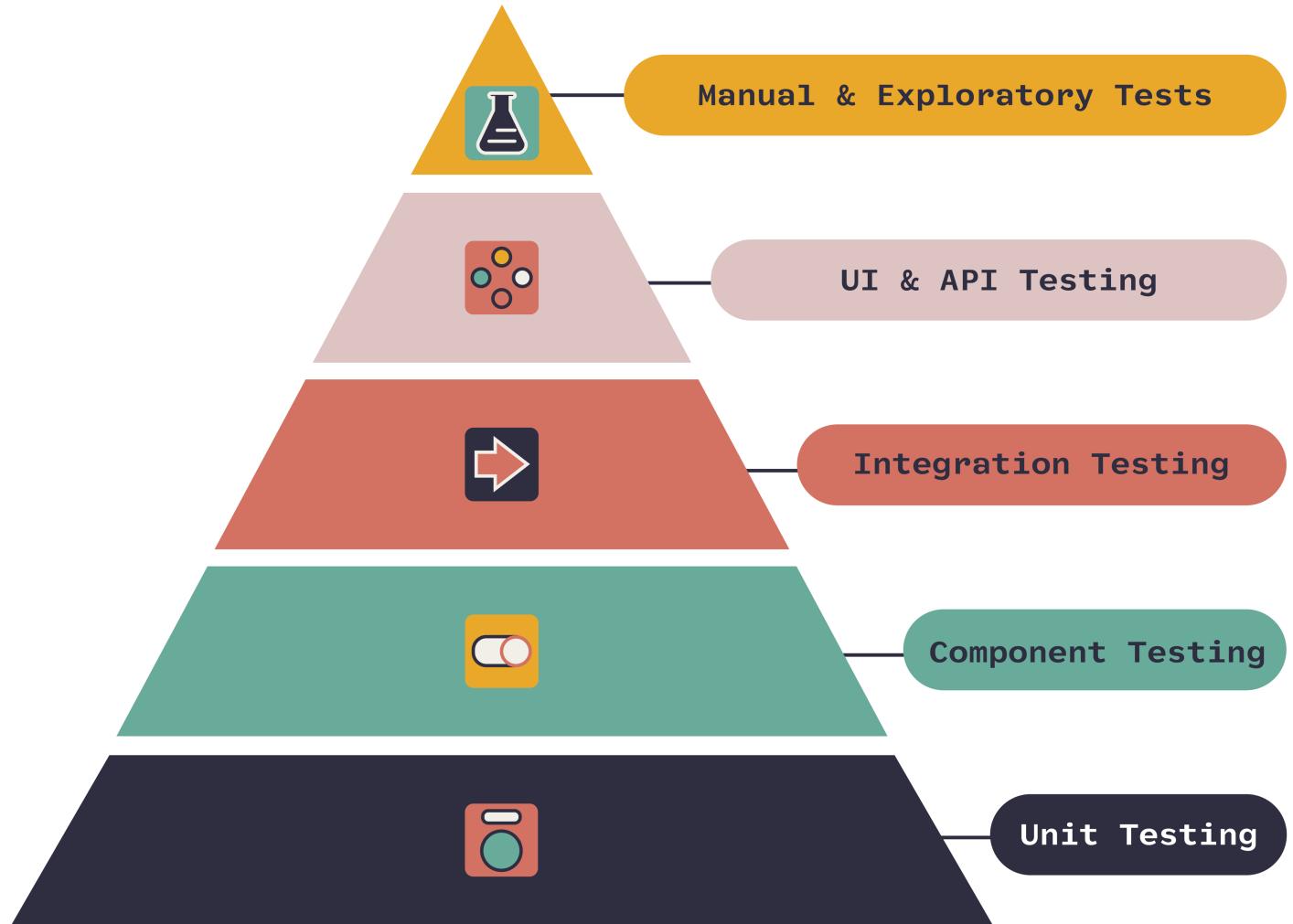
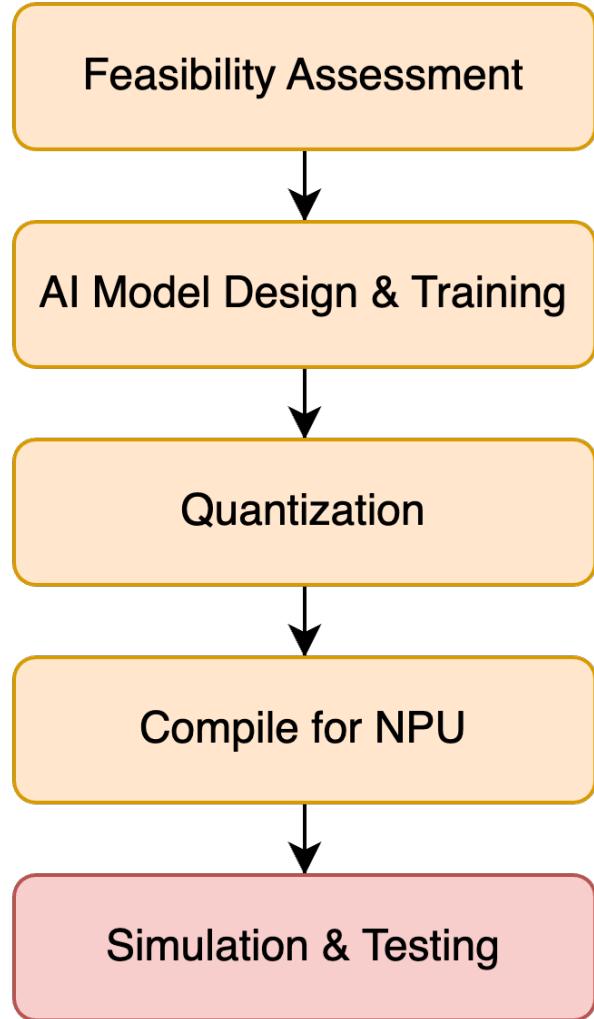
- Applications programmed against spec and verified against test cases.



Sounds familiar?



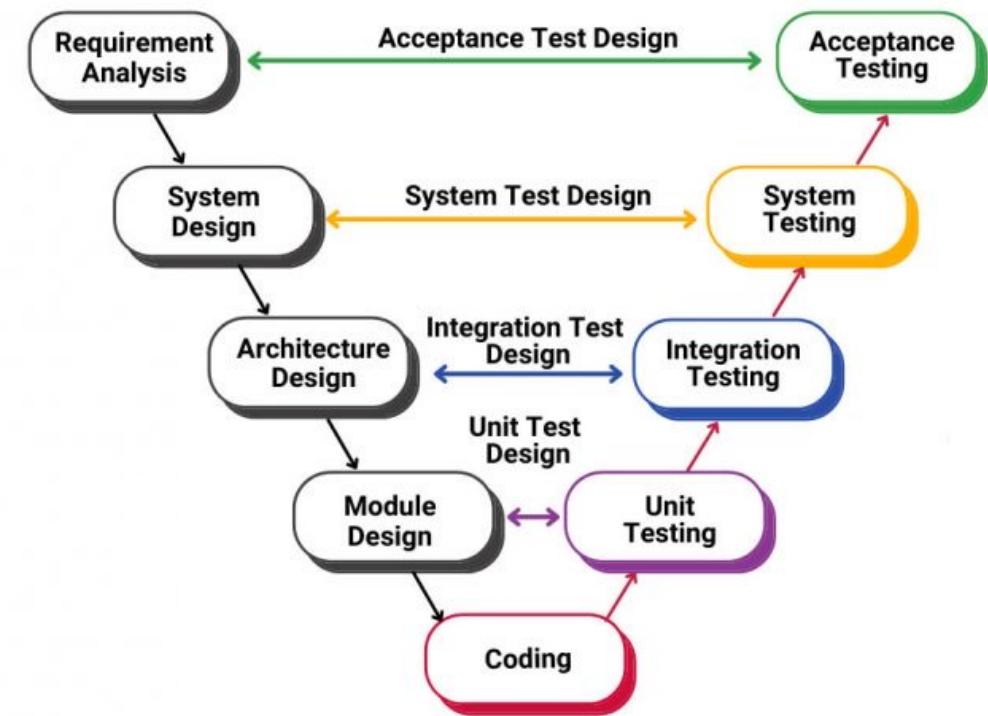
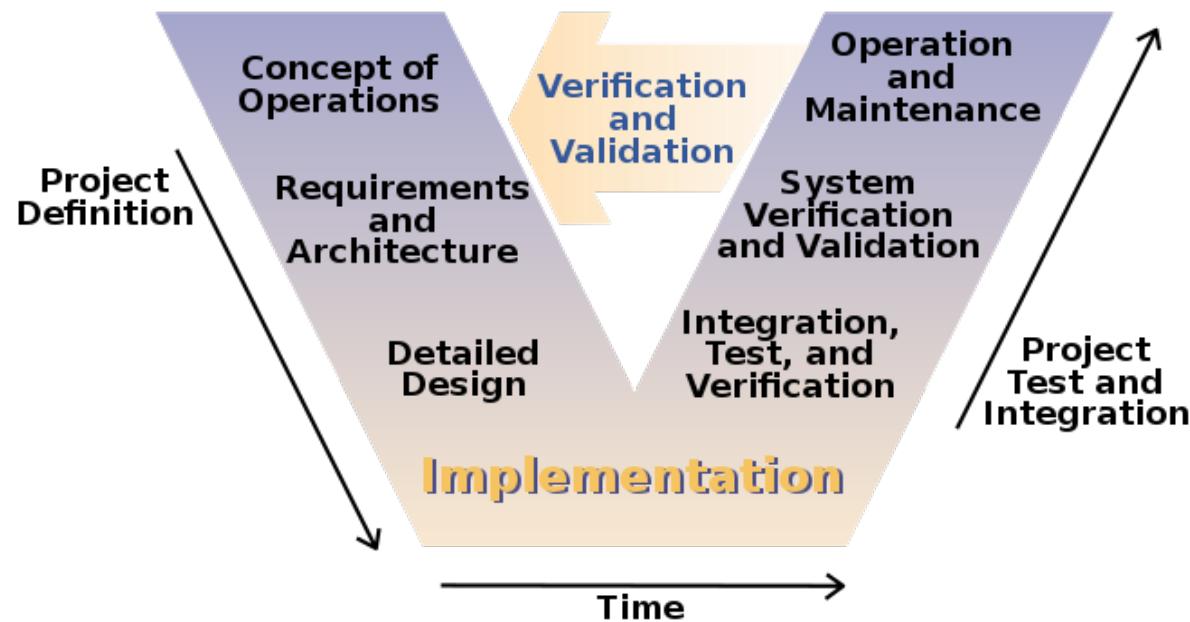
Importance of Testing



V-Model

Standard software development process

- Verification: Are we building the product right (according to spec)?
- Validation: Are we building the right product (meet user requirements/expectations)?



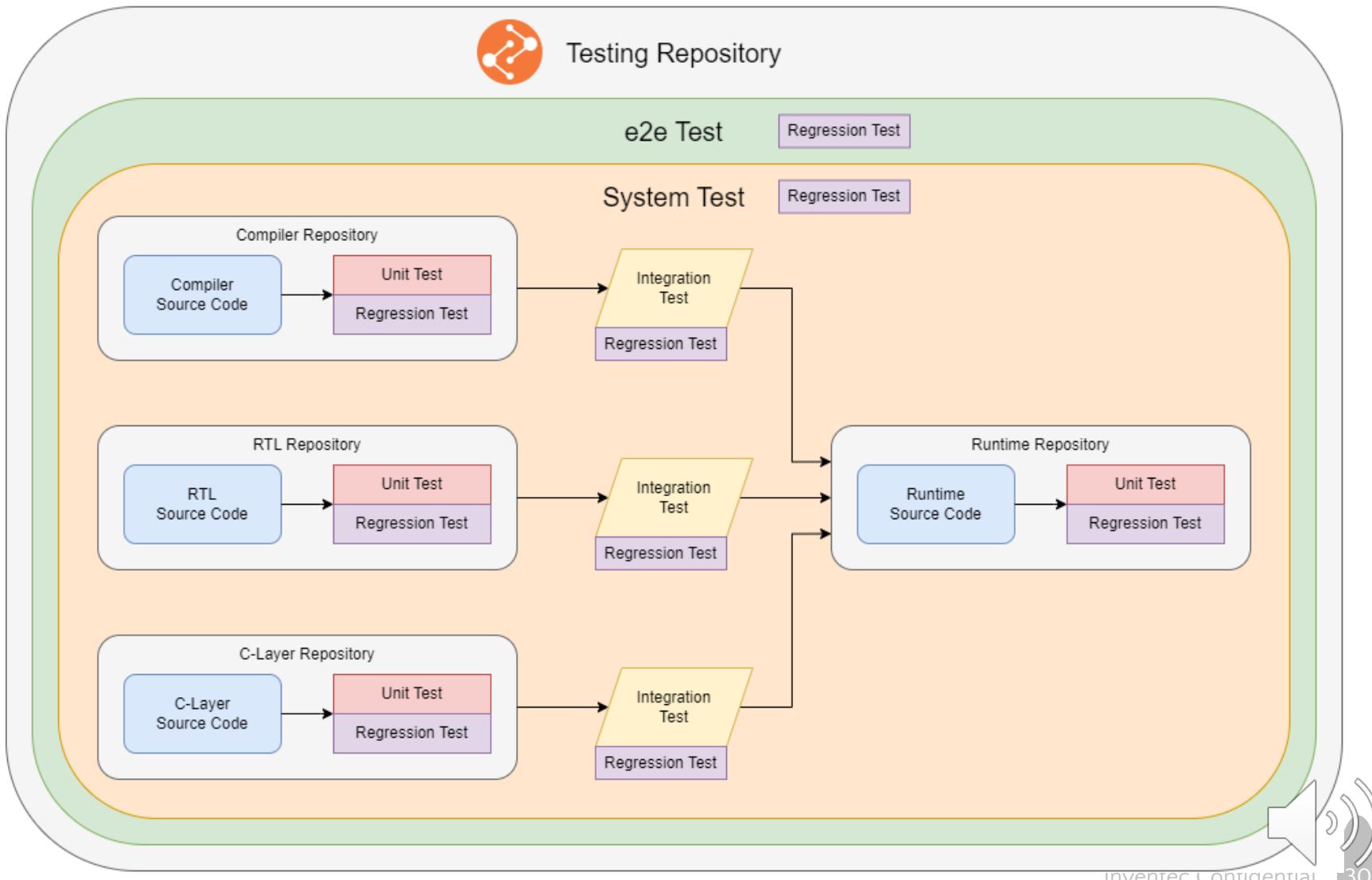
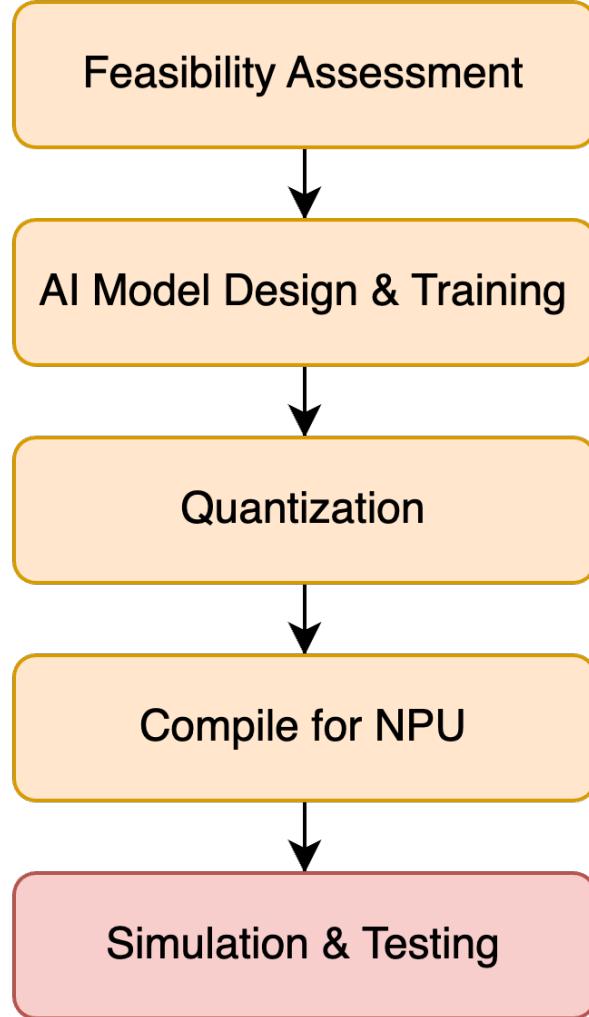
[https://en.wikipedia.org/wiki/V-model_\(software_development\)](https://en.wikipedia.org/wiki/V-model_(software_development))

<https://www.coleyconsulting.co.uk/v-model-verification-and-validation.htm>

Slide credit: Jacob



Testing System



Outline

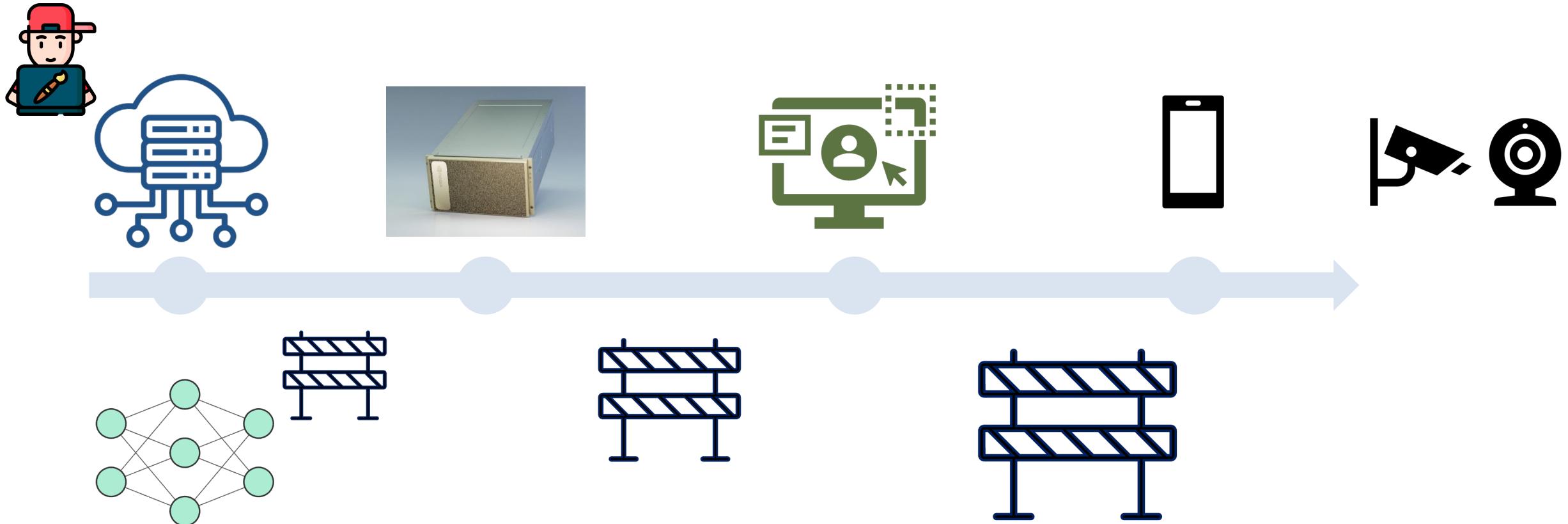
- Overview
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- **Hardware design**
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design



Software Engineer

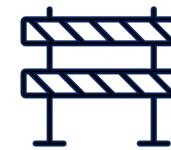
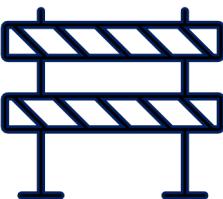
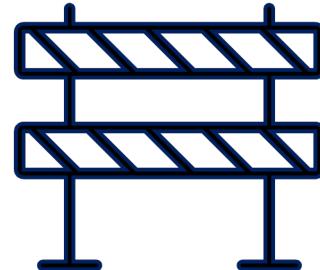
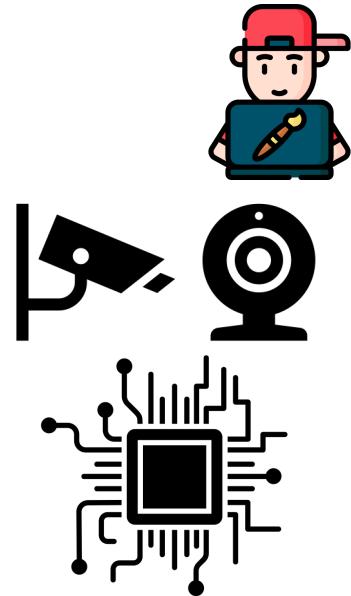
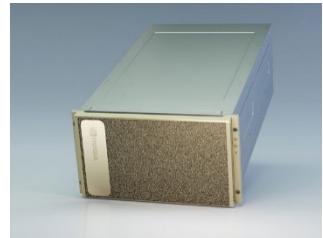
New challenges

- Edge-of-the-edge



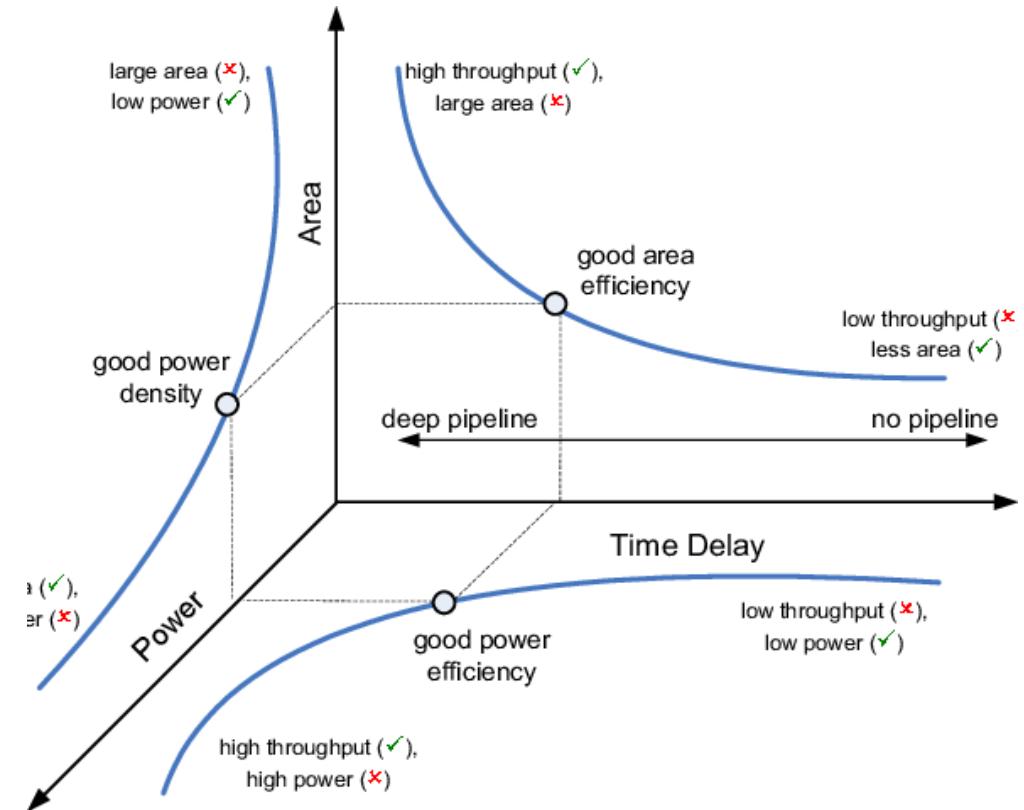
Hardware Engineer

New challenges



Measurement of Efficiency

- Hardware-specific
 - Area (Gate count)
 - Number of computing unit
 - Memory footprint
 - Time Delay
 - Power (W/s)
- Test metric – might be very different
 - Accuracy: The quality of the result for a given task
 - Throughput: Inferences per second.
 - Latency: The time between the input data arrives and the result is generated.
 - Energy: Total energy of executing a task.

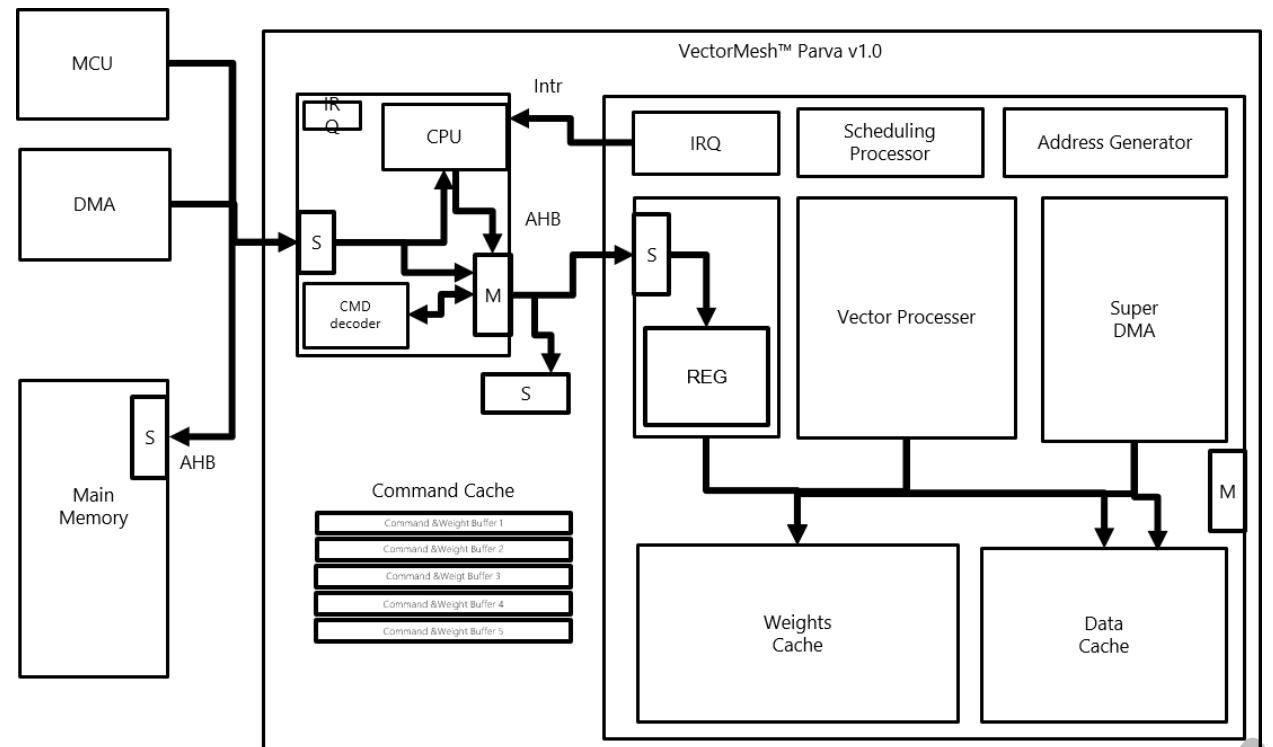
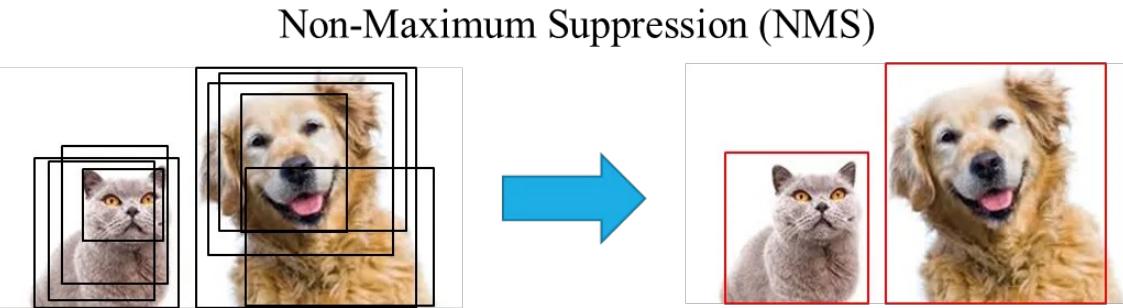


Outline

- Overview
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- Hardware design
 - **NPU**
 - ASIC (Application Specific Integrated Circuit)
- System-hardware co-design

NPU - Design Philosophy of RISC-V

- What computations do we need?
 - Image pre-processing (denoising, gamma, ...)
 - Neural network inference (CNN kernel computation, activation, ...)
 - Data post-processing (NMS, ID matching, ...)
- What kinds of hardware do we have?
 - NPU: CNN operation
 - DMA: resize, crop
 - MCU: runtime, NMS
 - ISP: image enhancement



Slide credit: Jacob

Outline

- Overview
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- Hardware design
 - NPU
 - **ASIC (Application Specific Integrated Circuit)**
- System-hardware co-design

Systolic Array

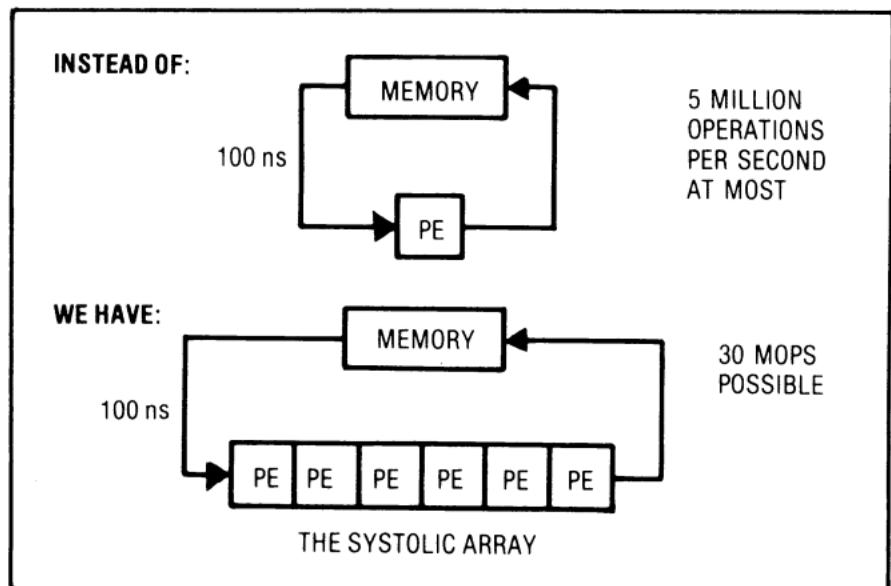
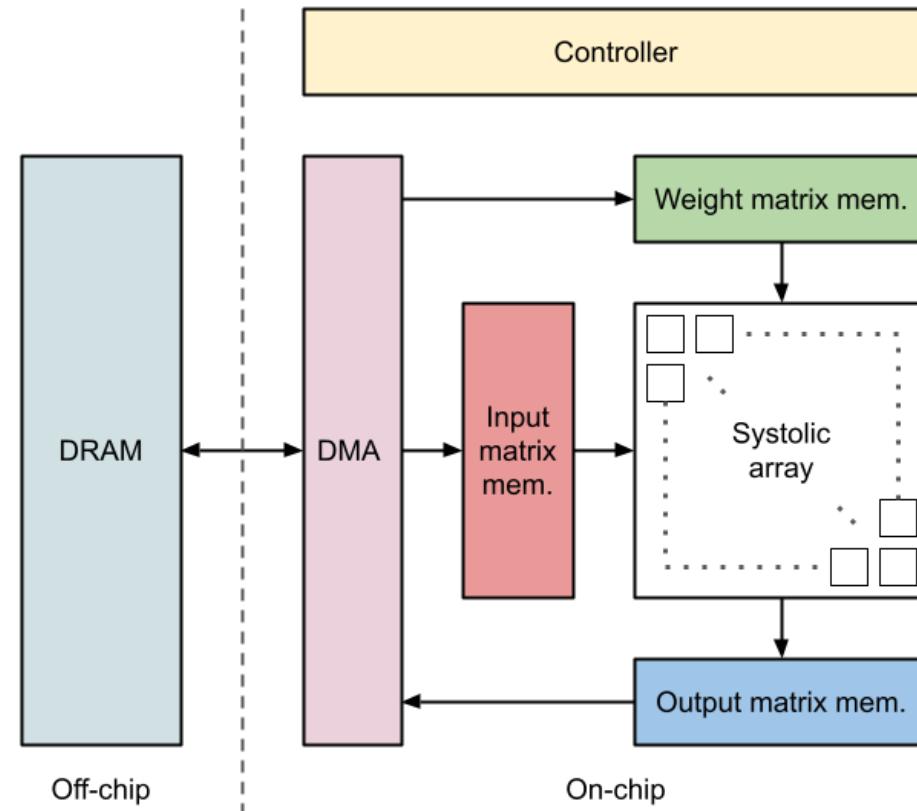
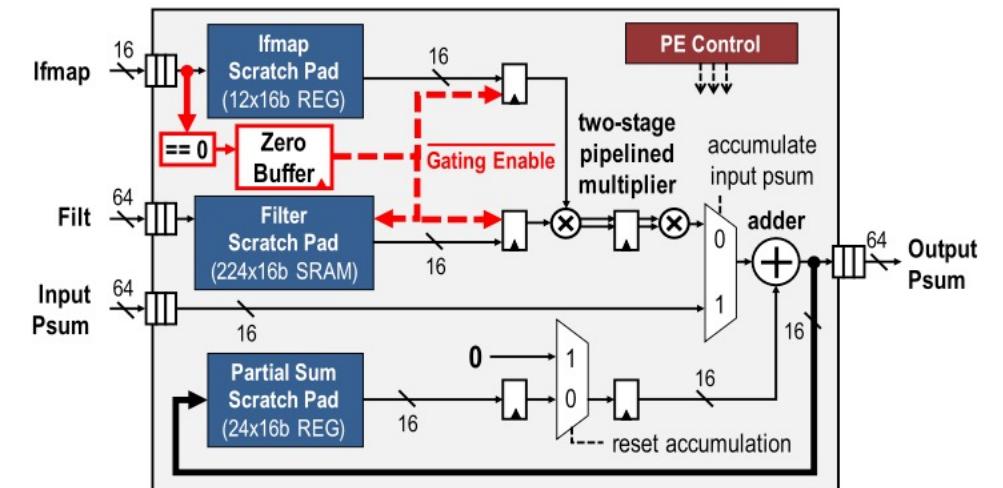
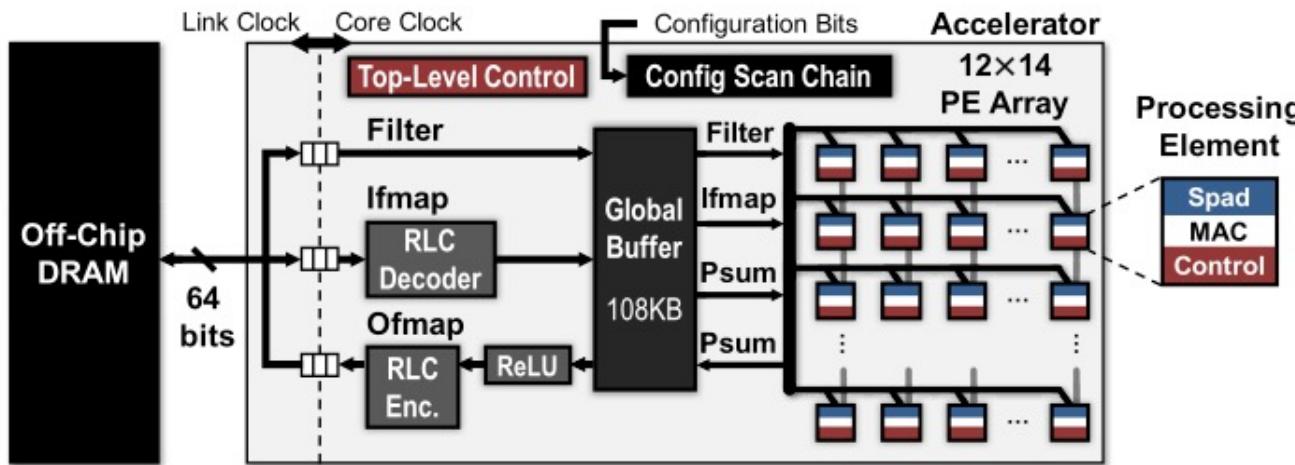


Figure 1. Basic principle of a systolic system. https://blog.csdn.net/wangwangmoon_light/article/details/121890472 @wangwangmoon_light



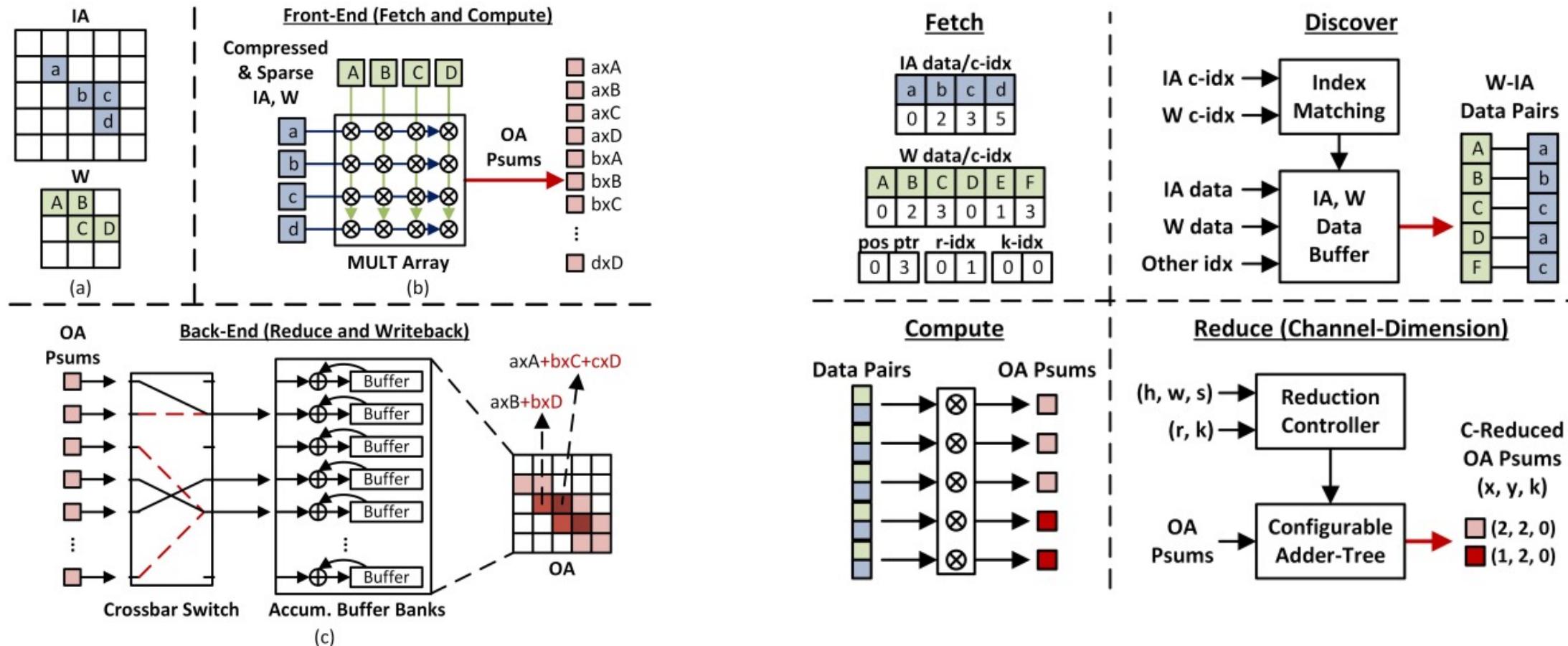
<https://towardsdatascience.com/accelerating-neural-networks-on-hardware-baa3c14cd5ba>

Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks

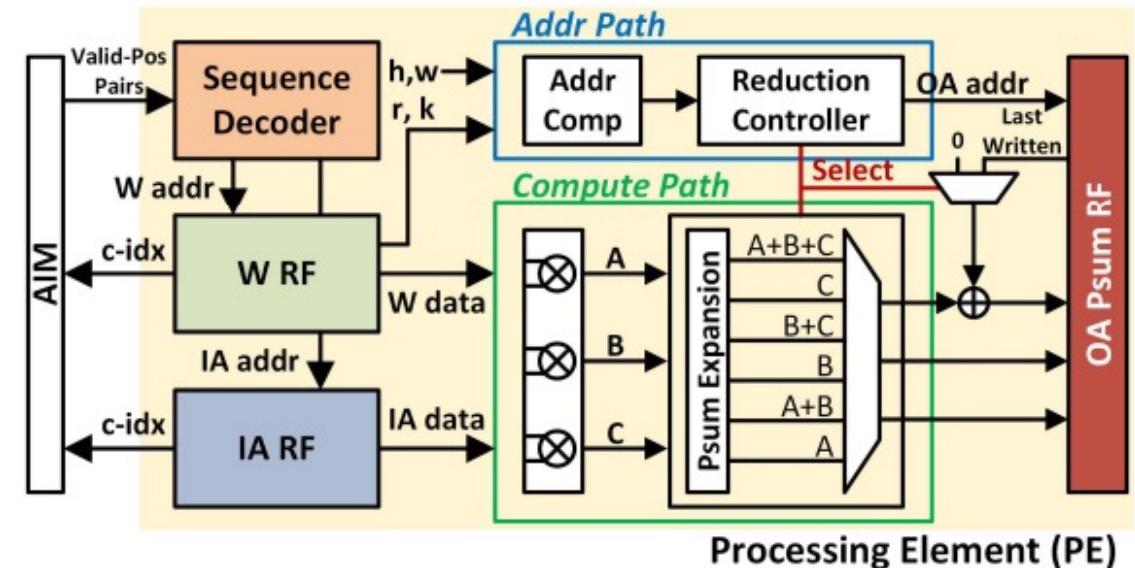
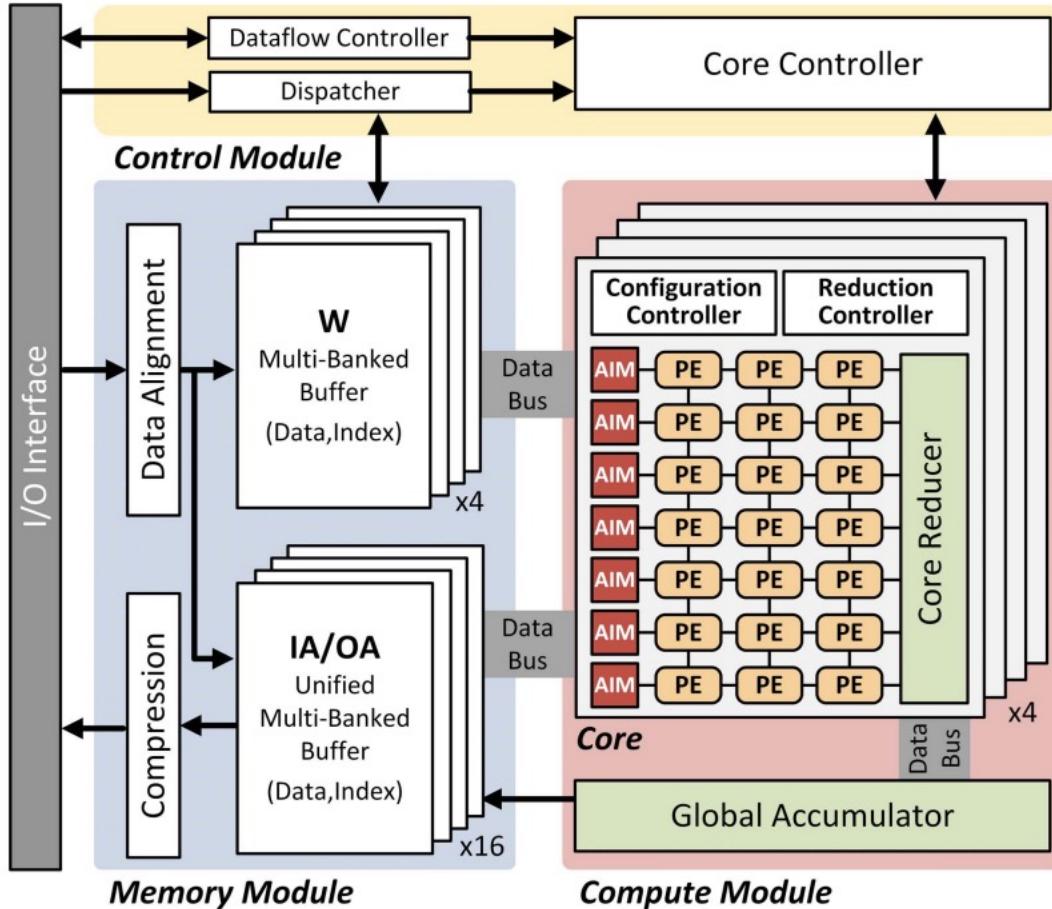


SNAP: An Efficient Sparse Neural Acceleration Processor for Unstructured Sparse Deep Neural Network Inference

- Channel last vs channel first



SNAP: An Efficient Sparse Neural Acceleration Processor for Unstructured Sparse Deep Neural Network Inference



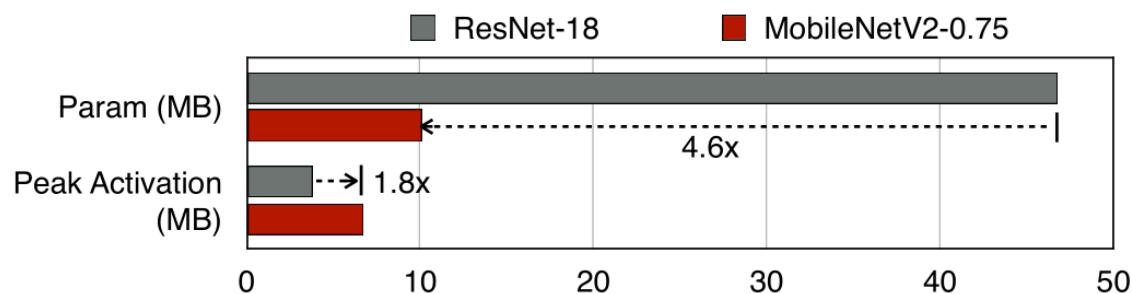
Outline

- Overview
- AI on Chip (AloC) team work flow
- Quantization
- Testing System
- Hardware design
 - NPU
 - ASIC (Application Specific Integrated Circuit)
- **System-hardware co-design**

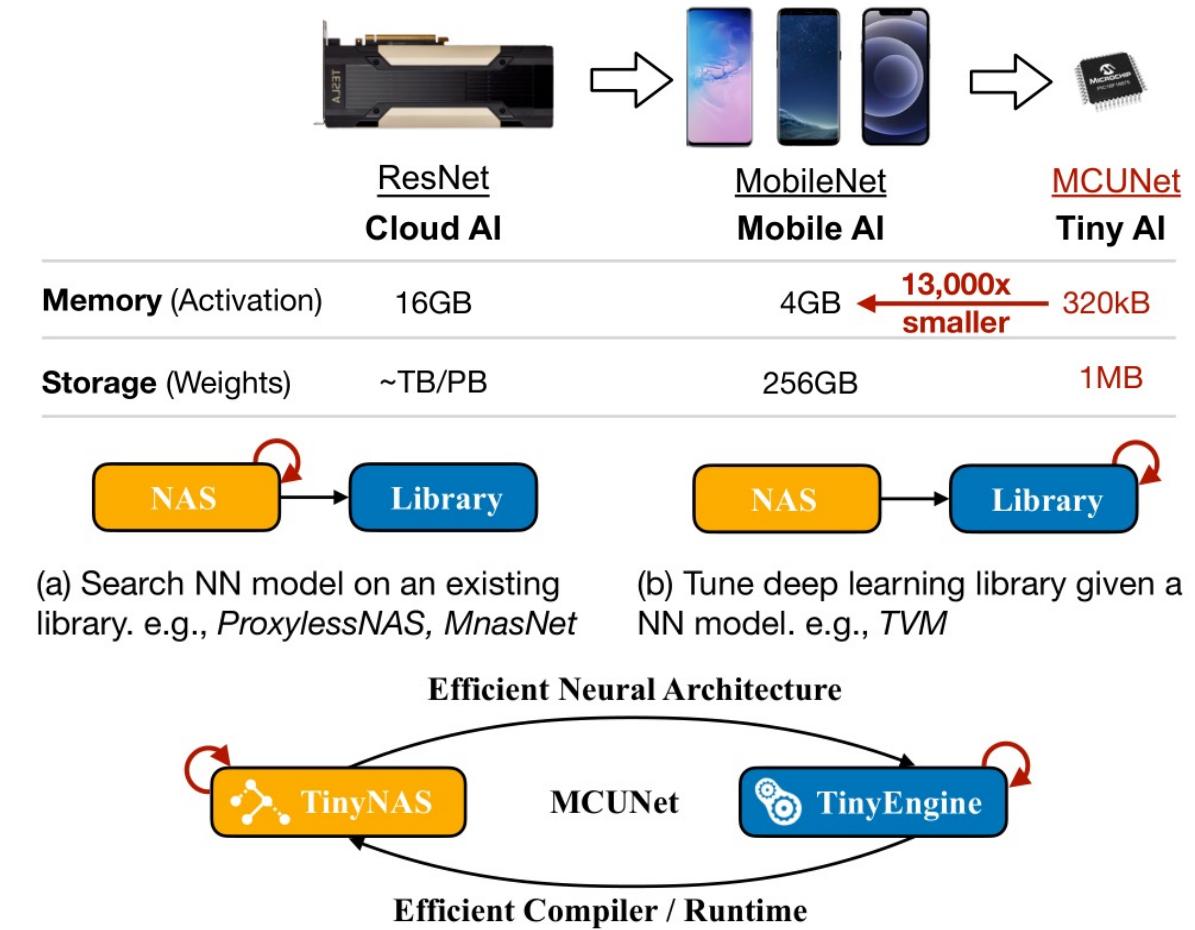
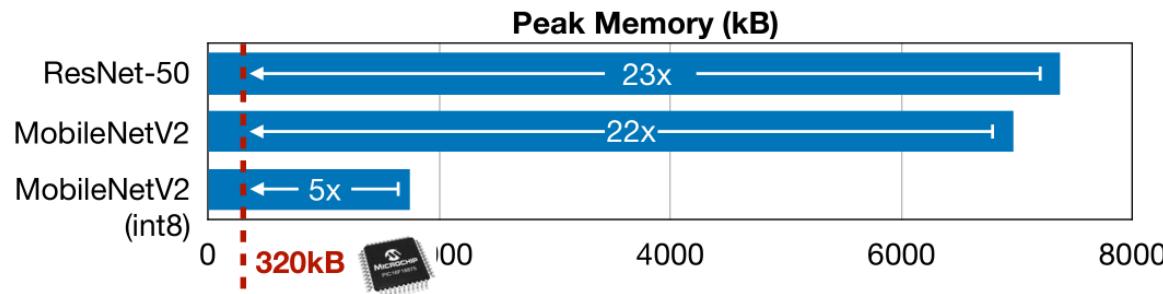
MCUNet: Tiny Deep Learning on IoT Devices

System-Algorithm Co-design

- Existing Methods Reduce Model Size, but not the Activation Size

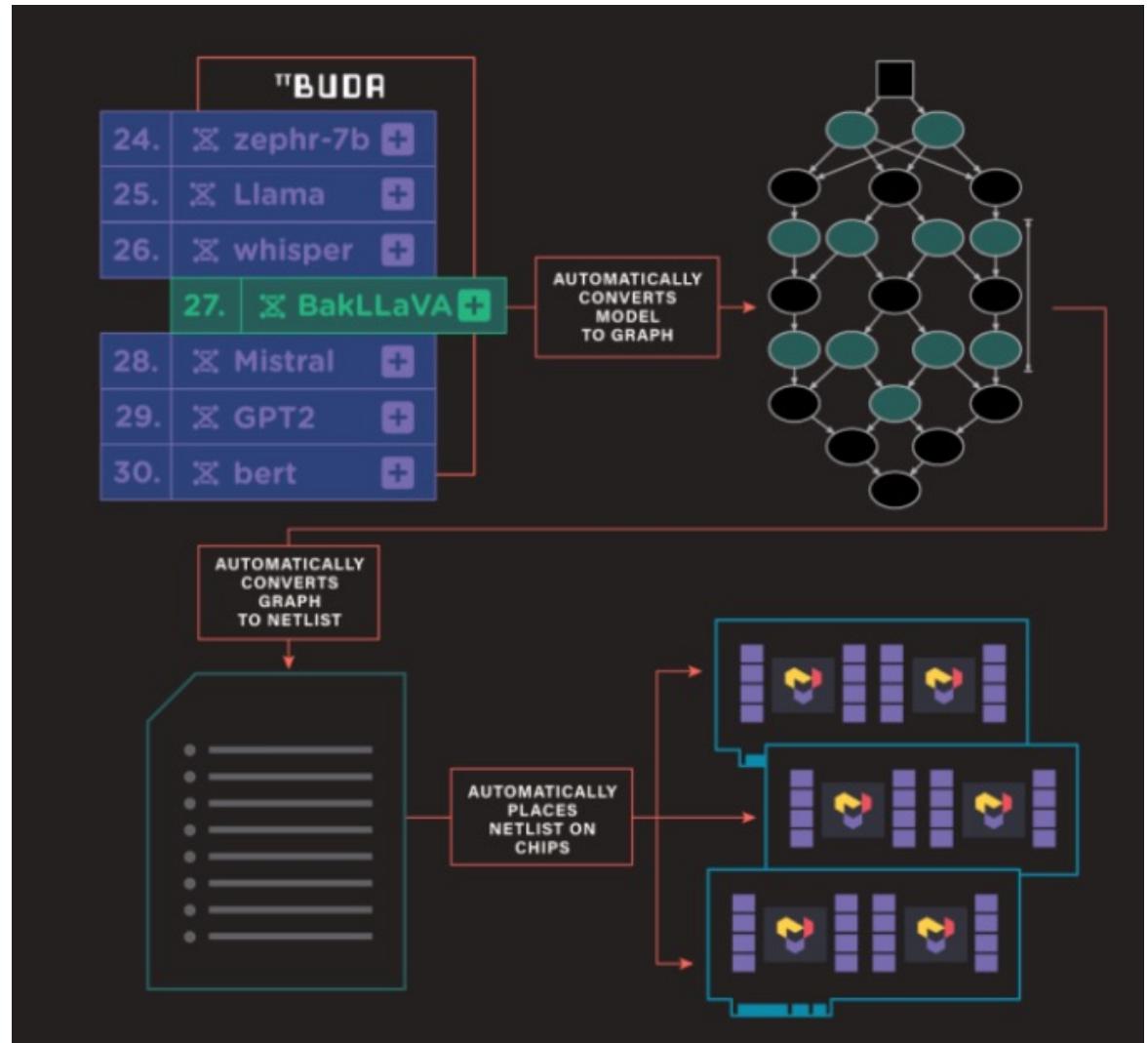
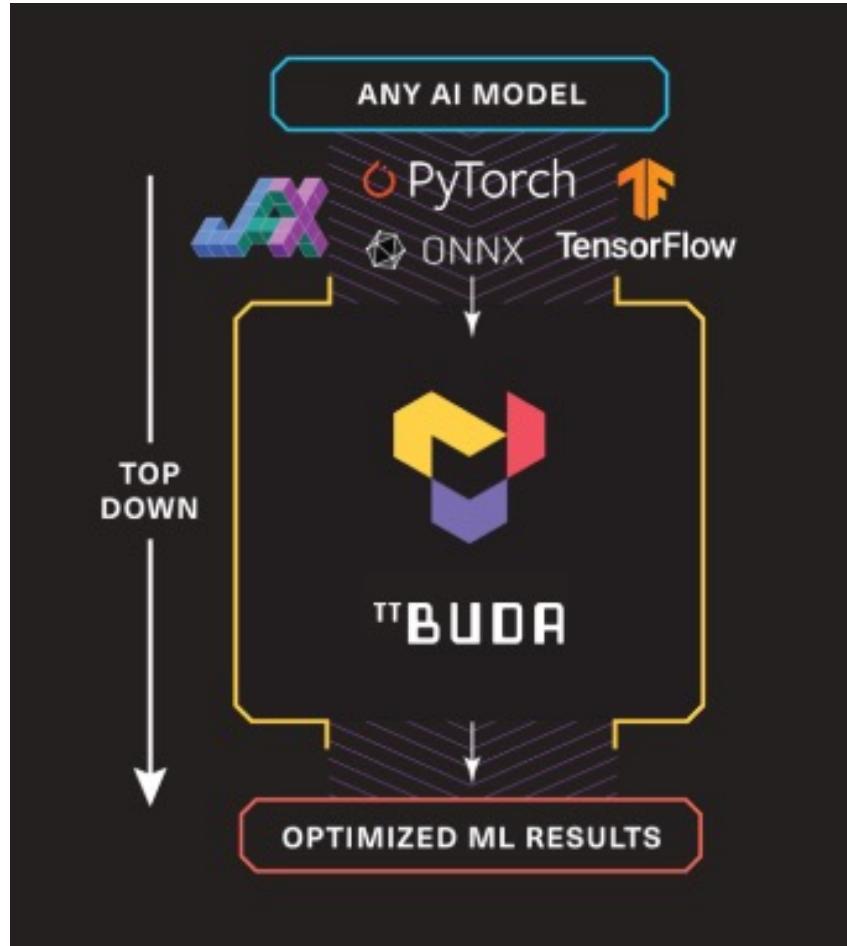


- Existing network **CANNOT** fit the tight memory constraints on MCU



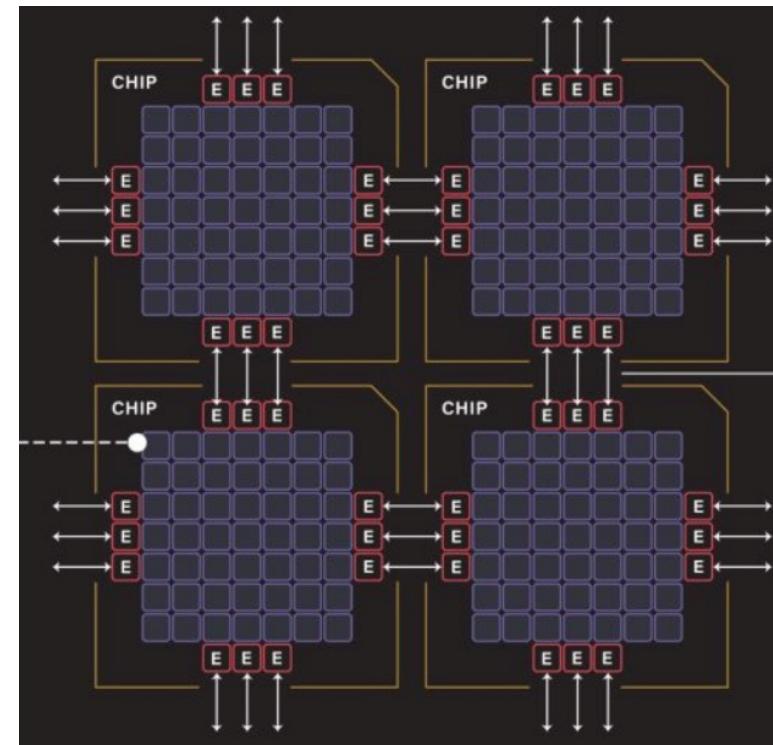
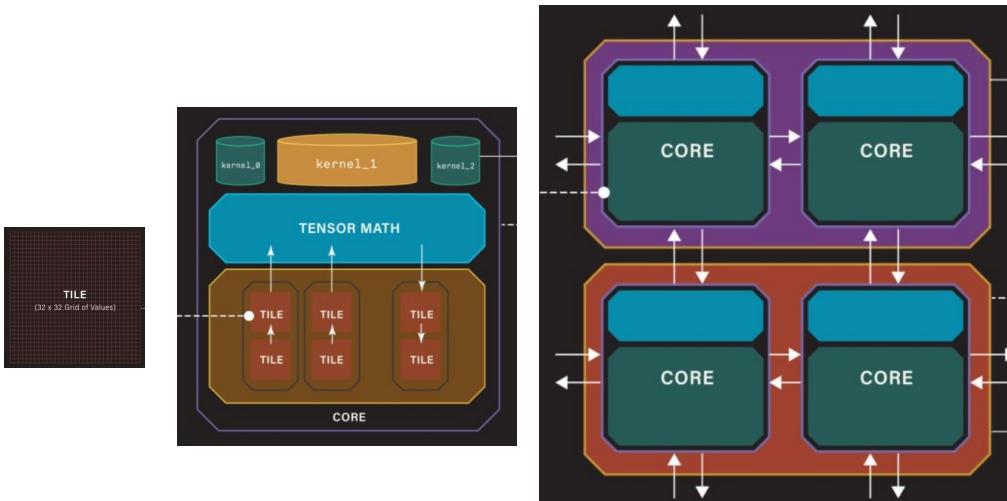
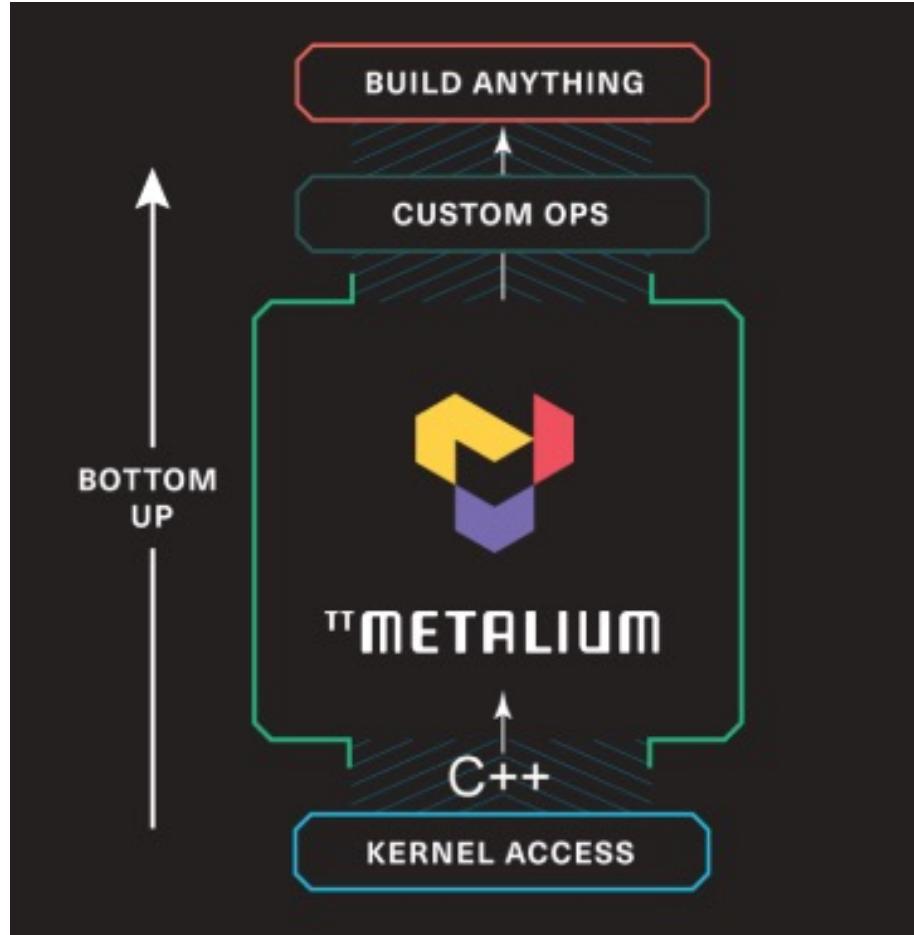
tenstorrent

Top down



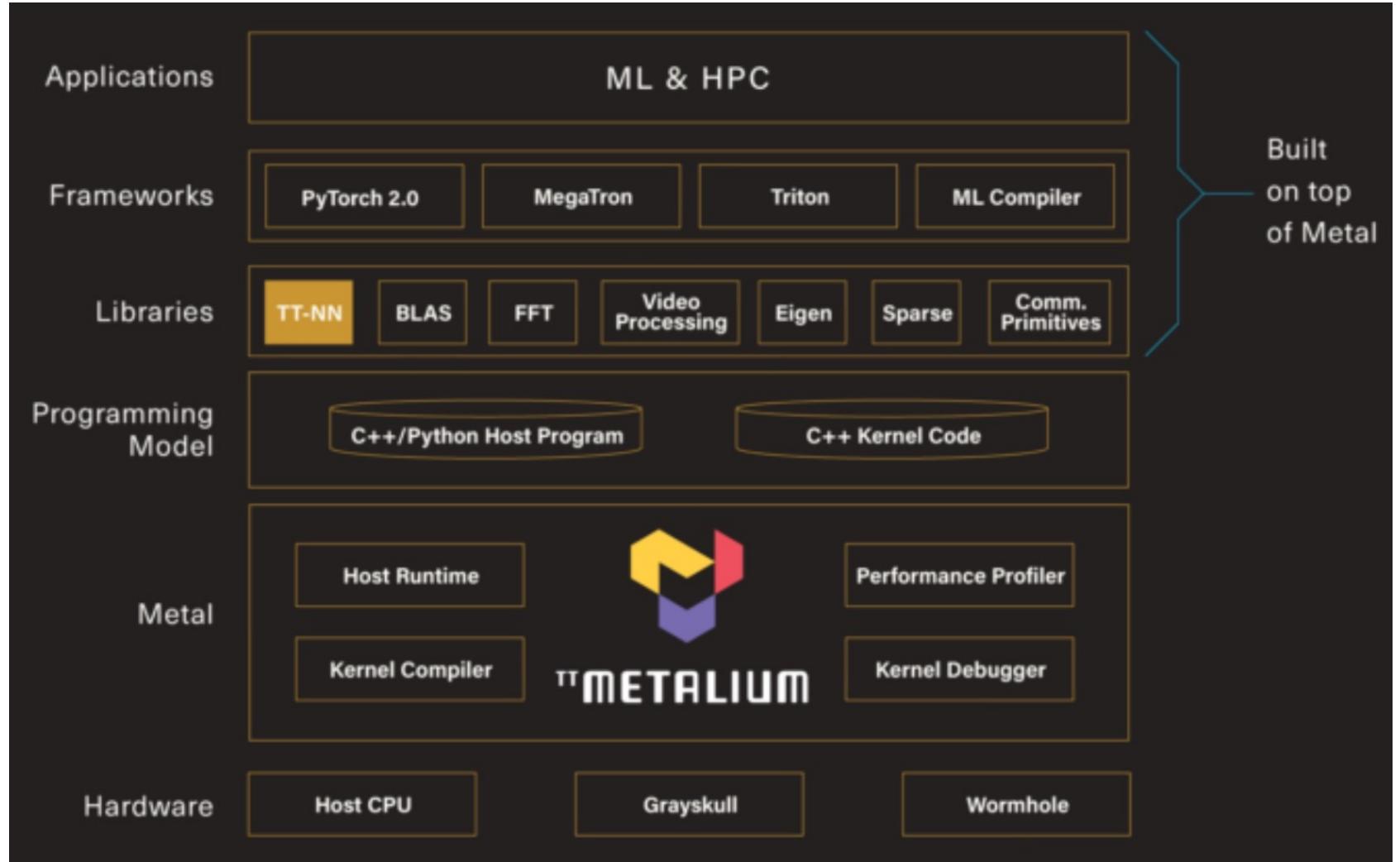
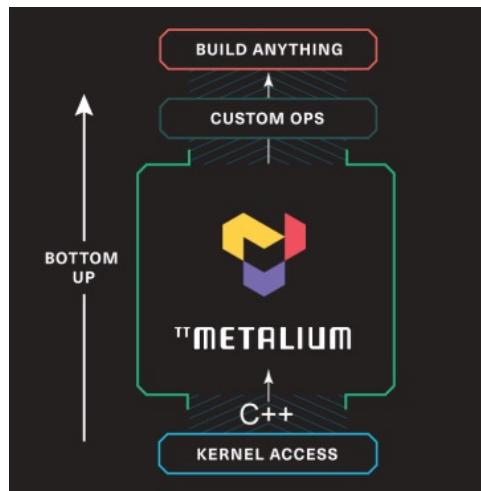
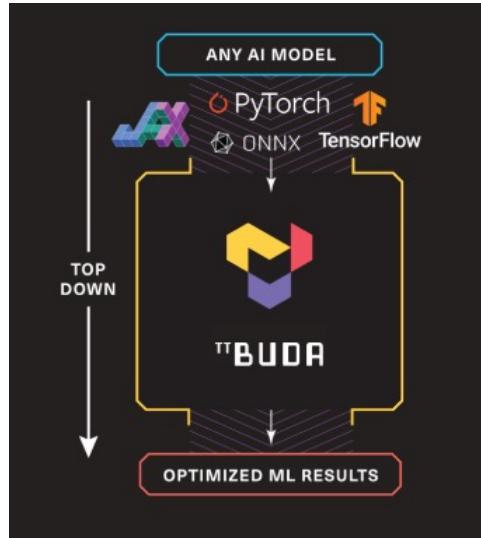
tenstorrent

Bottom Up



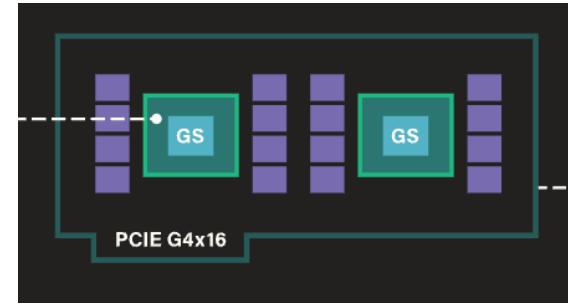
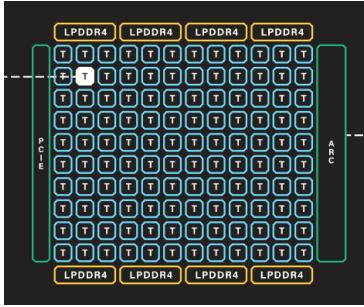
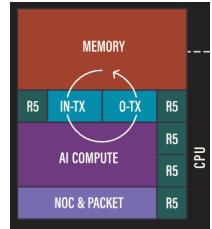
tenstorrent

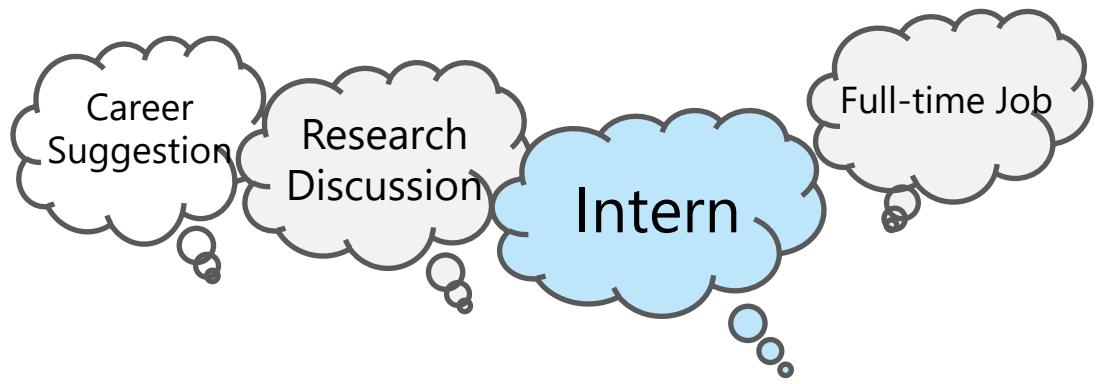
Integration



AI Cloud

Scale Up





Welcome to Contact Us!

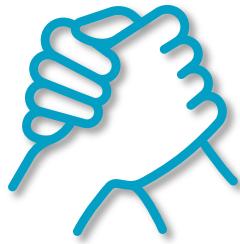
<https://ai.inventec.com>



Shang-Jui (Ray) Kuo (郭尚睿)
Inventec AI 研究員
kuo.raysj@inventec.com

Jeng-Lin (John) Li | 李政霖
Inventec 資深資料科學家
li.johncl@inventec.com

Chia-Ching (Jacob) Lin | 林家慶
Inventec 資深資料科學家
lin.jacob@inventec.com



Thank you!