

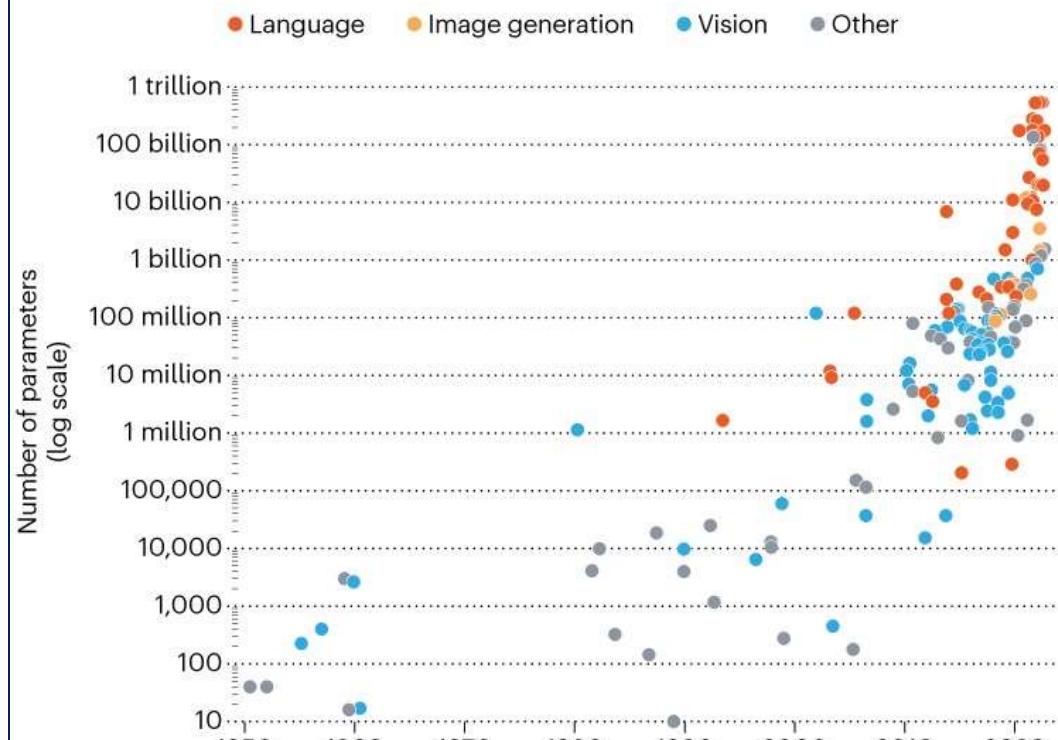
# Enabling massive adoption of TinyML models in edge devices thanks to In-Memory Computing

Michele Rossi

# Why Edge AI?

## THE DRIVE TO BIGGER AI MODELS

The scale of artificial-intelligence neural networks is growing exponentially, as measured by the models' parameters (roughly, the number of connections between their neurons)\*.

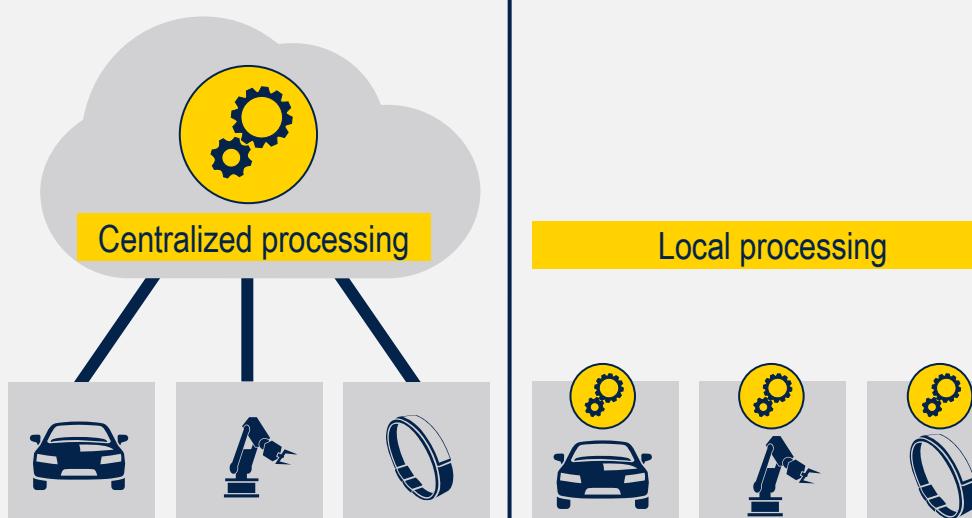


\*Sparse' models, which have more than one trillion parameters but use only a fraction of them in each computation, are not shown.

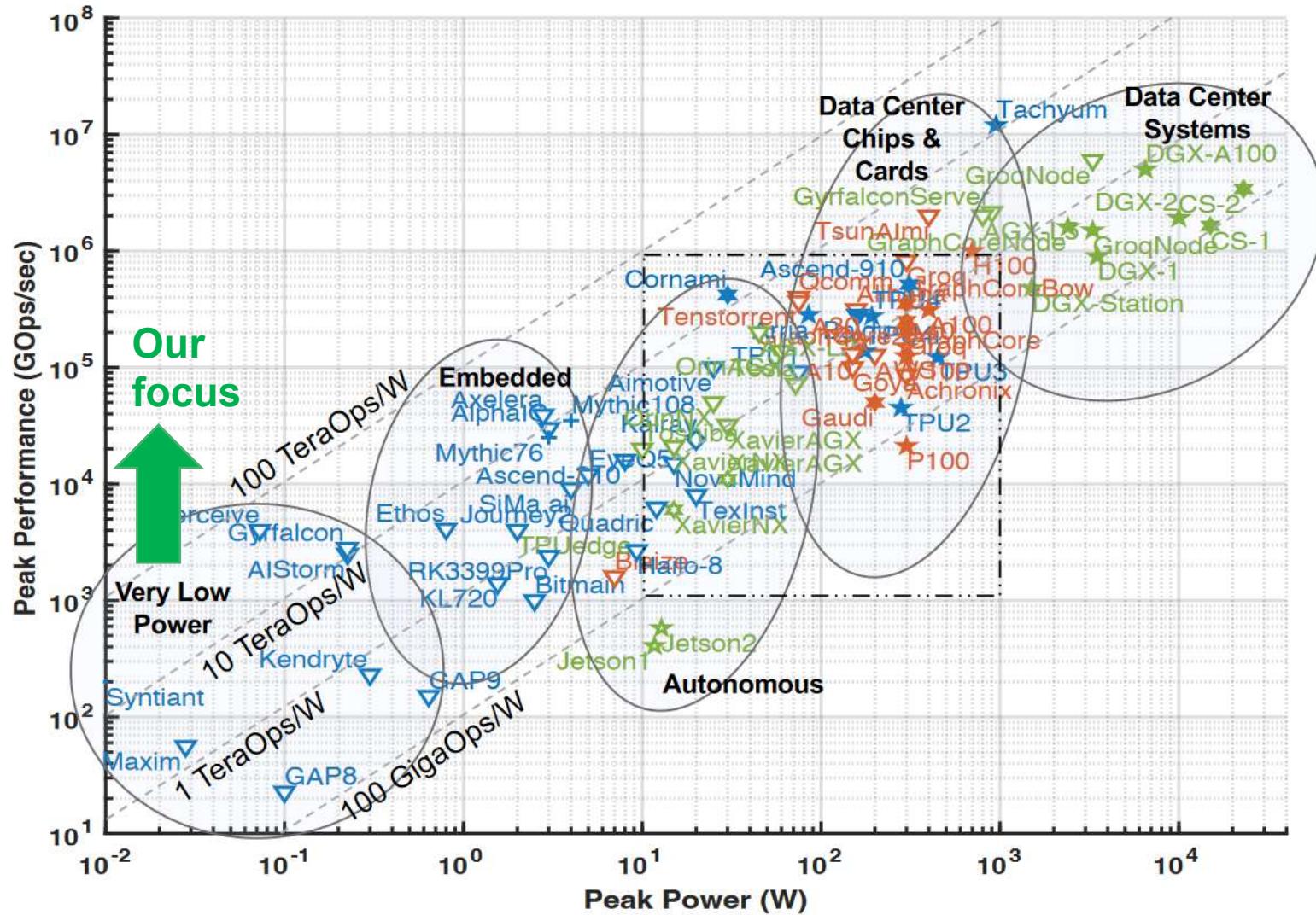
©nature



# Cloud vs. Edge

| Icon  | Cloud computing            | Edge computing  | Icon |
|---|----------------------------|-----------------|------|
|   | Expensive                  | Affordable      |      |
|   | Inefficient                | Efficient       |      |
|   | Remote                     | On-prem         |      |
|   | Privacy at risk            | Privacy granted |      |
|   | Multiple points of failure | Reliable        |      |
|  |                            |                 |      |

# From the Cloud to the Edge



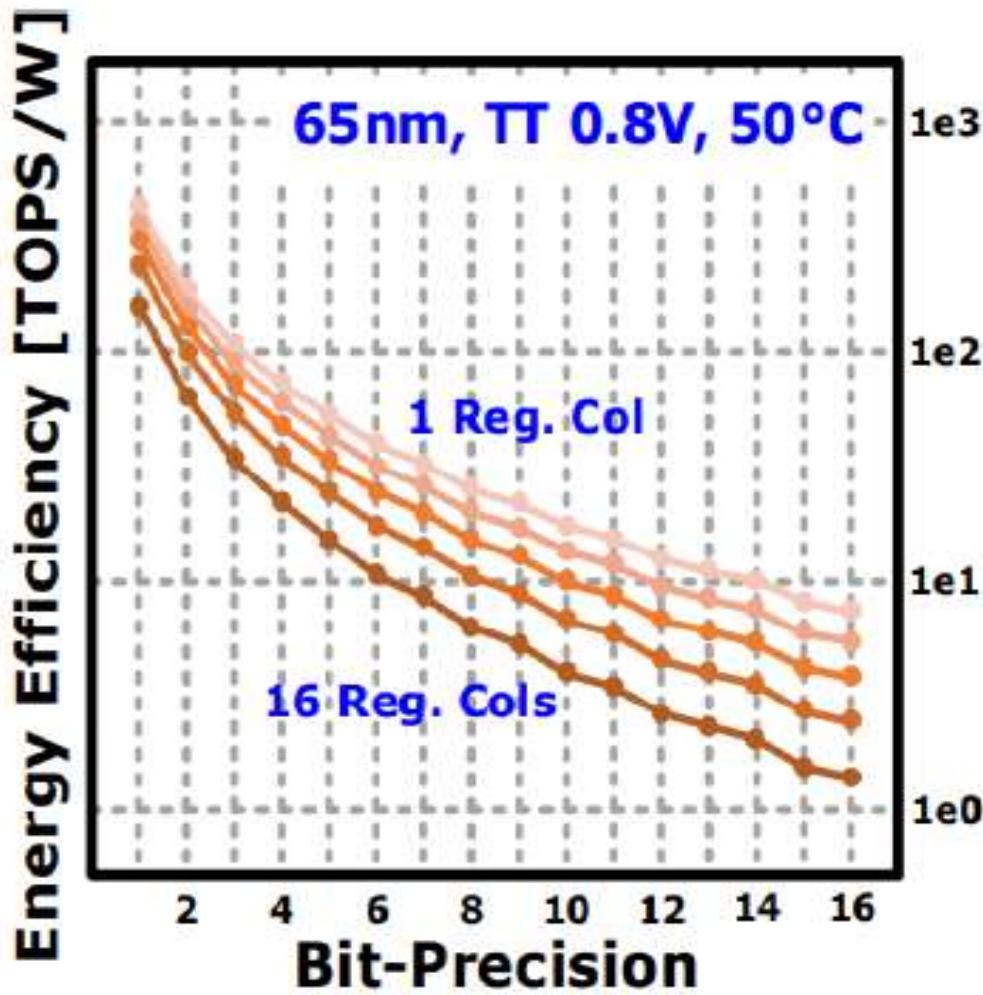
## Limits of edge devices:

- memory
  - area
  - performances
  - battery-powered

# Solutions:

- TinyML
  - NPUs
  - New technologies

# Quantization



google/qkeras

QKeras: a quantization deep learning library for  
Tensorflow Keras



20  
Contributors

9  
Used by

442  
Stars

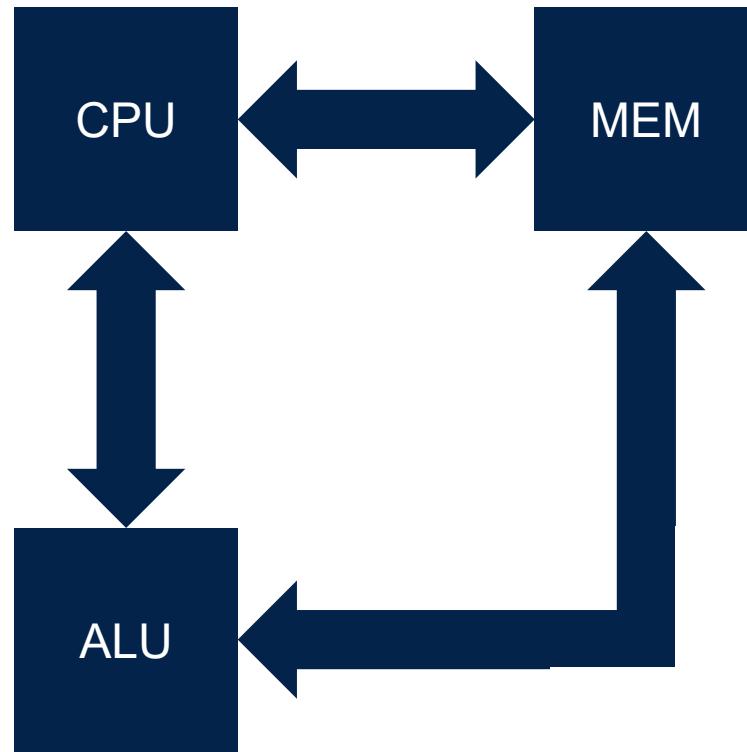
89  
Forks



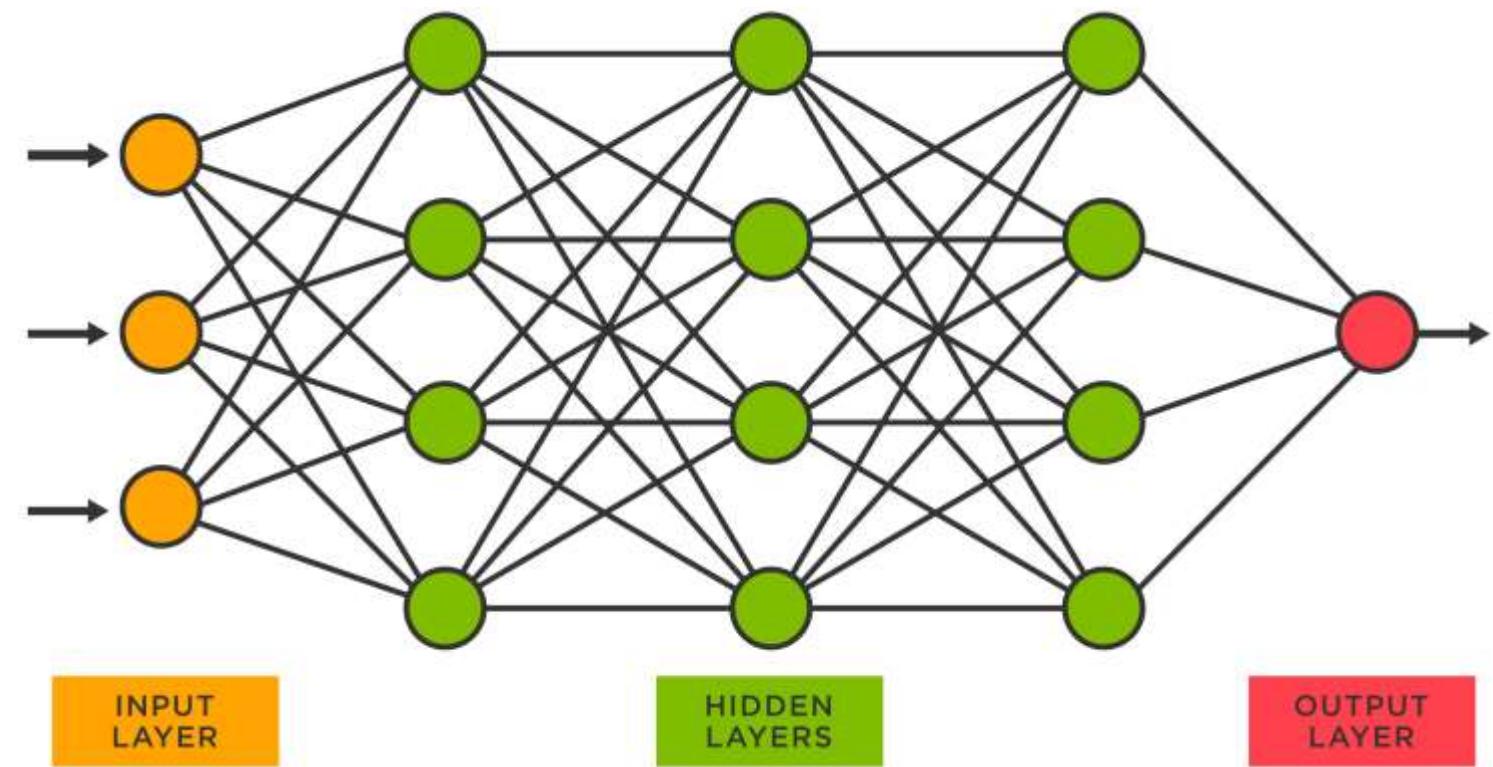
<https://github.com/google/qkeras>

# Von Neumann Bottleneck

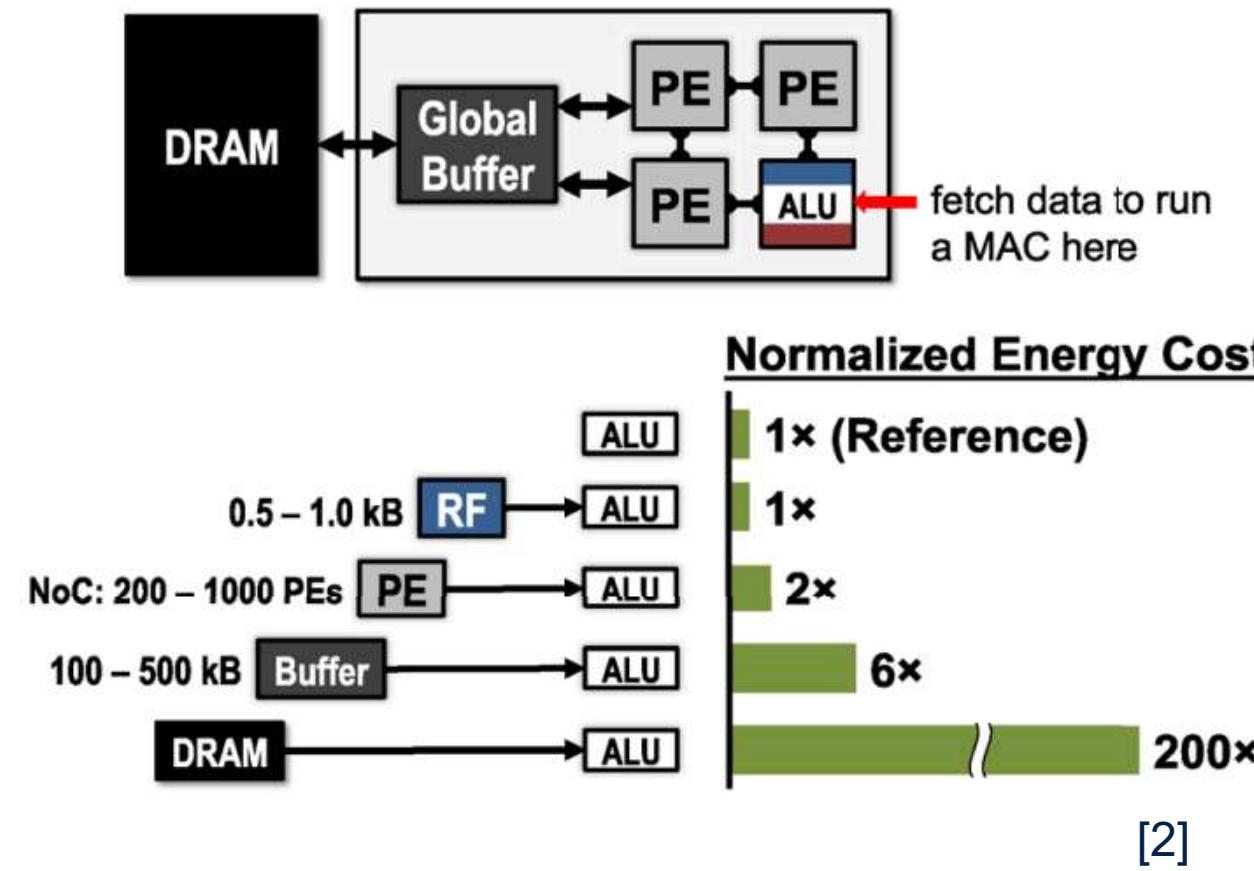
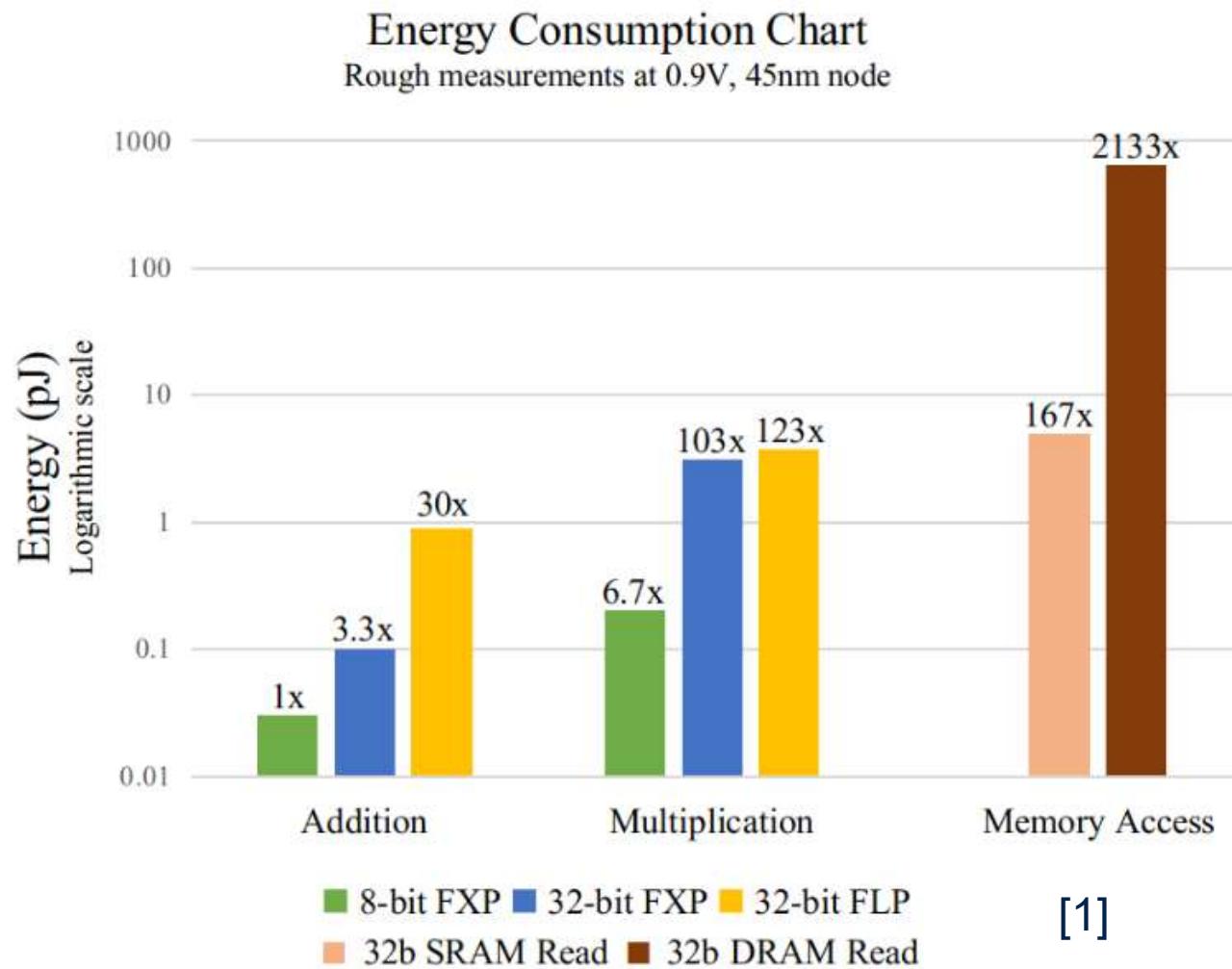
Von Neumann Architecture



Neural Network Architecture

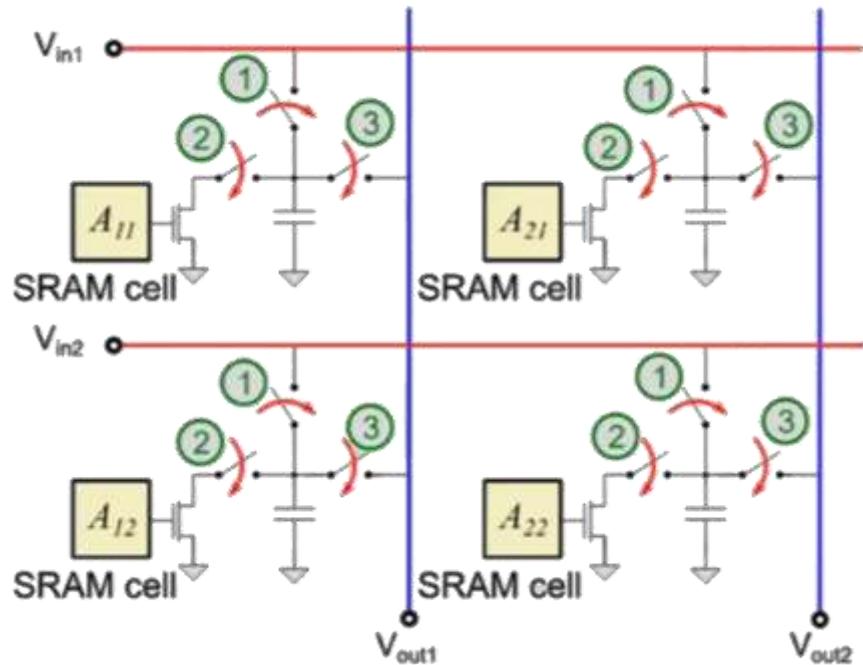


# Memory Wall



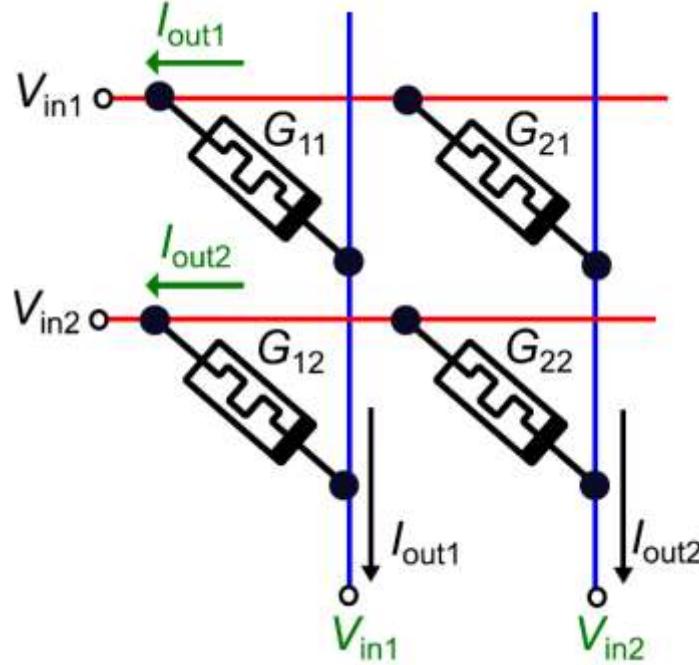
# In-Memory Computing

## Analog



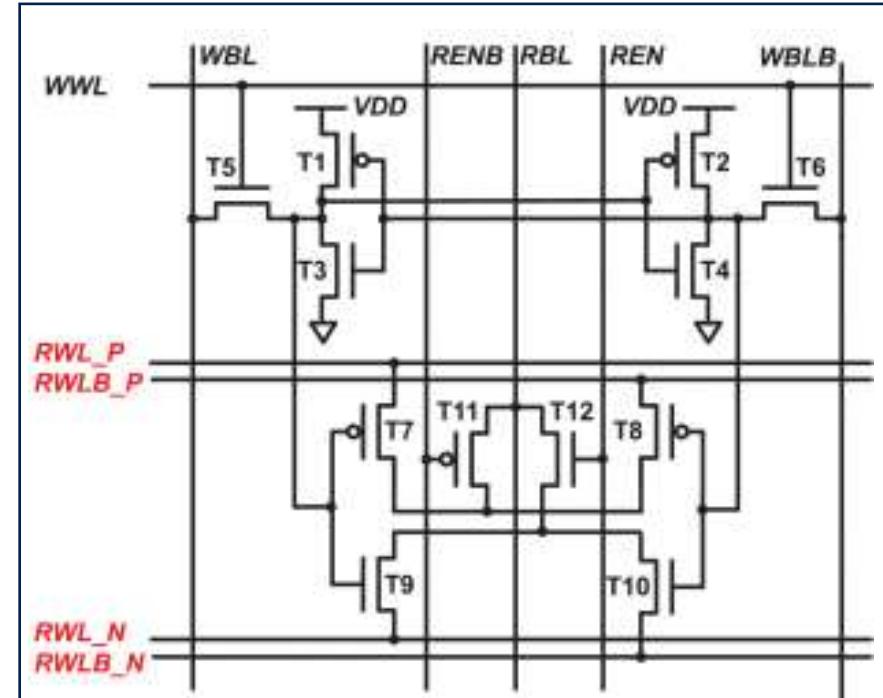
Capacitor-based

[4]



Resistive

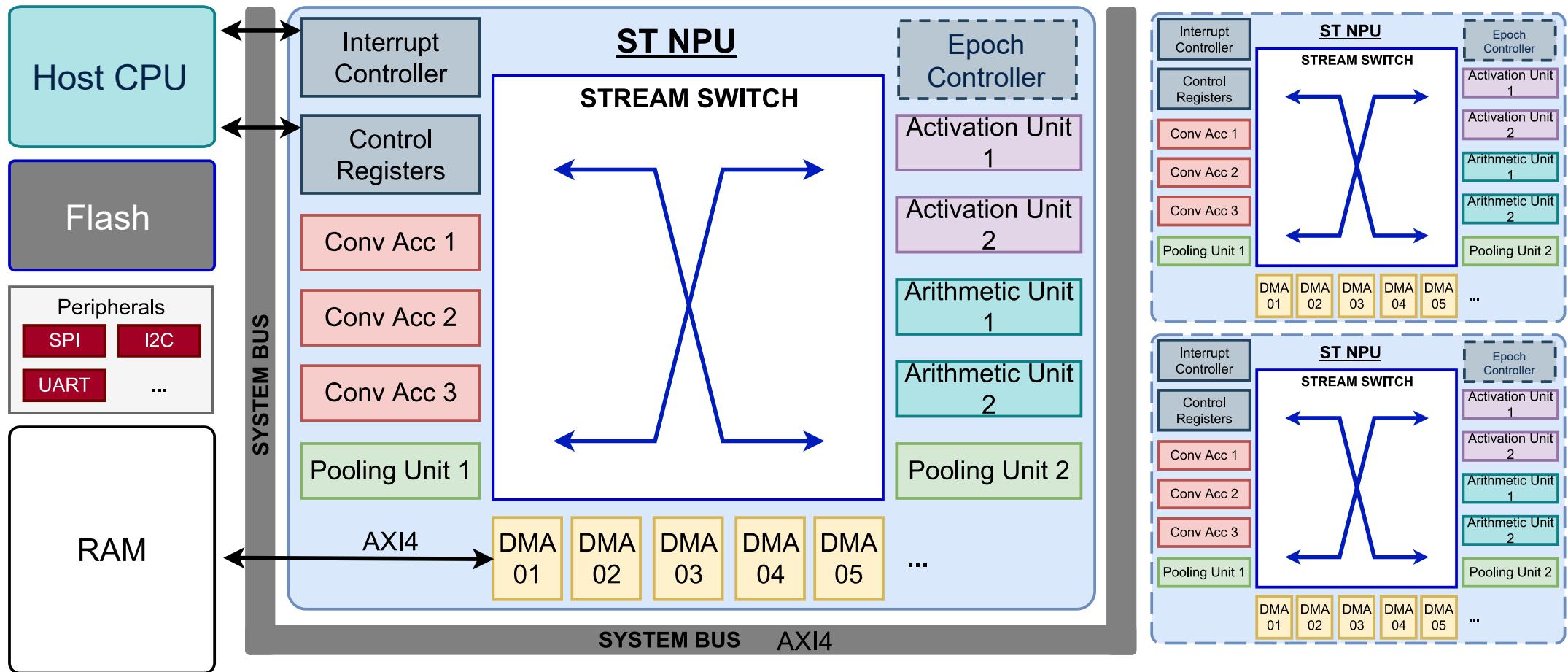
[4][5]



SRAM

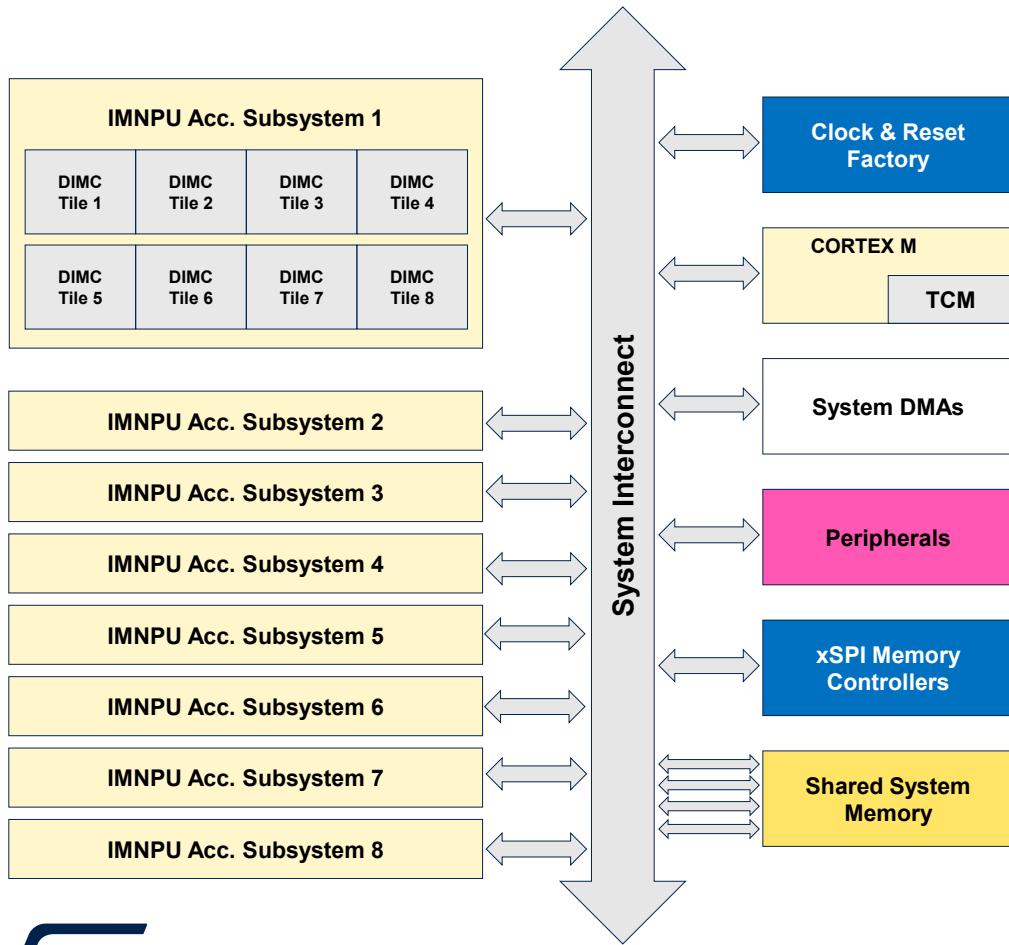
[6]

# Architecture Template

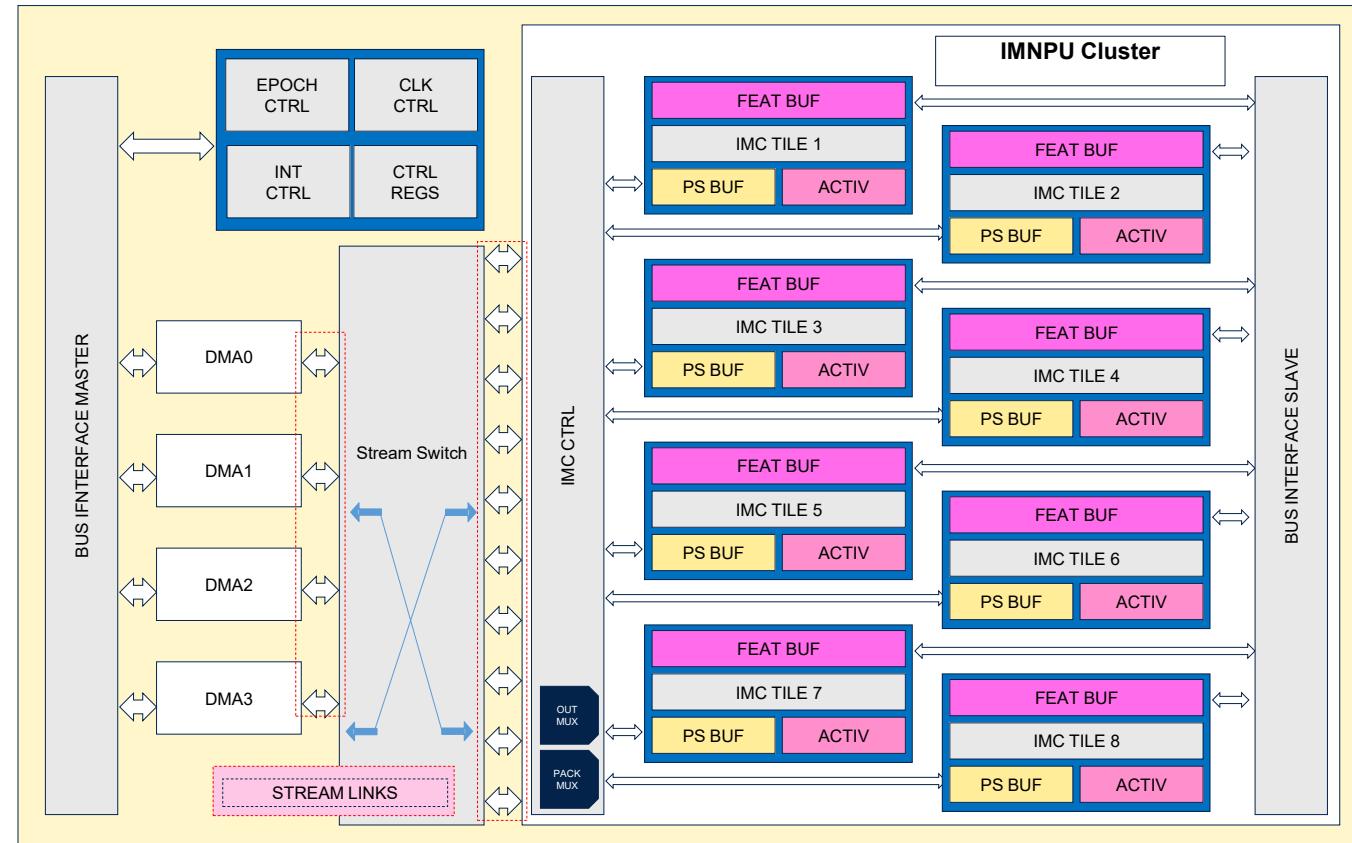


# IMNPU System-on-Chip

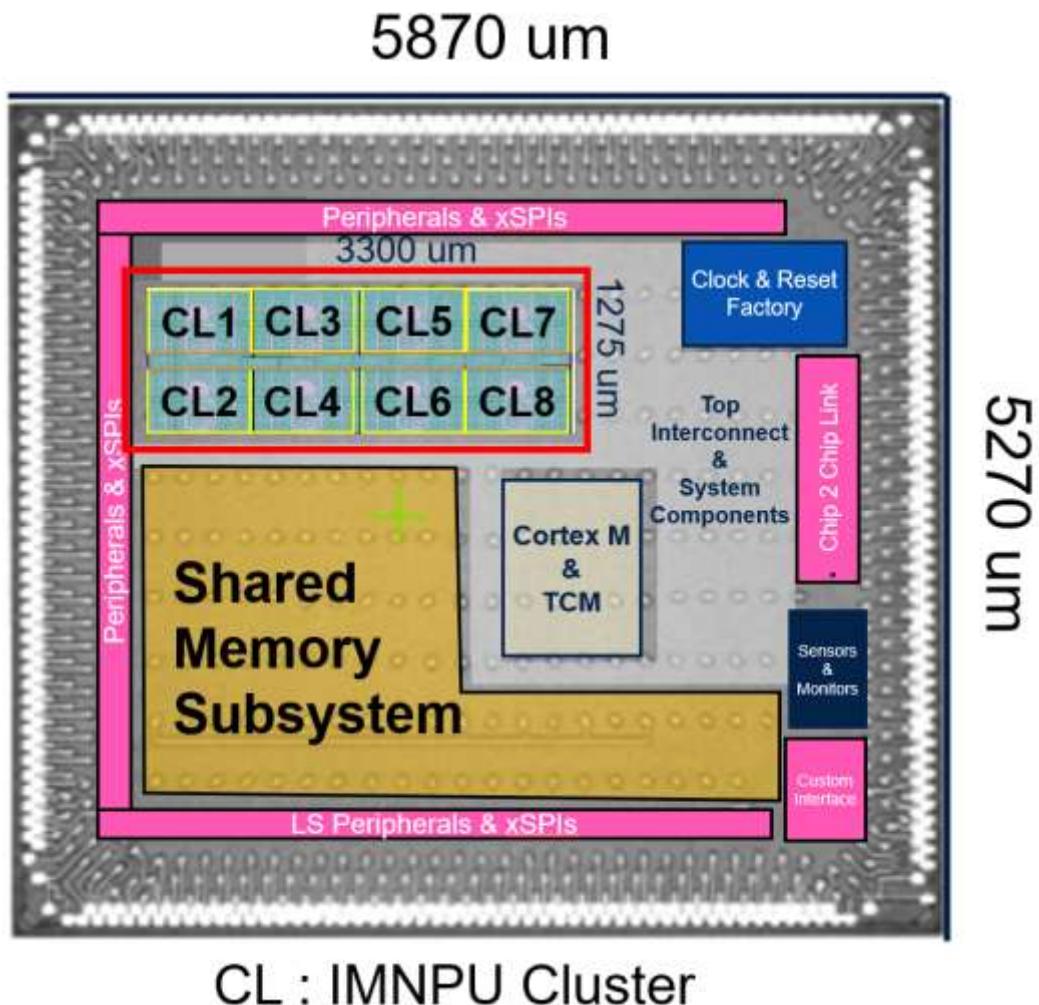
## IMNPU System-on-Chip



## IMNPU Accelerator Subsystem



# Prototype Chip



|   |
|---|
| Technology 18nm FDSOI   |
| Multi-Cluster IMNPU along with system interconnect: 4.2 mm <sup>2</sup> |
| Voltage range: 0.525-1.0V,<br>FBB 0-1.5V                                |
| IMC Capacity 2 Mb   |
| Computation: Deterministic  |
| Precision Mode: 1-4 bits  |
| 229 TOPS (Peak Performance) 1 bit Weight - 1bit Feature                 |
| 57 TOPS (Peak Performance) 4bit Weight - 4bit Feature                   |
| 310 TOPS/W (1 bit)  |
| 77 TOPS/W (4 bit)   |
| 54 TOPS/mm <sup>2</sup> (1 bit)   |
| 13.6 TOPS/mm <sup>2</sup> (4 bit)                                       |
| CNN, LSTM, RNN  |

# Measurements

## Ultra-low-power always-on staged inference application example

Endurance for a battery-operated device for video surveillance, a single cluster is always on clocked @ 10 MHz running a VGG16-like network (as in Fig. 16.7.4) while the rest of the SoC is powered down and selectively activated @ 400 MHz to process a 10x complexity network (e.g., a ResNet152) with a duty cycle of 1 to 100 at 0.525 VDD

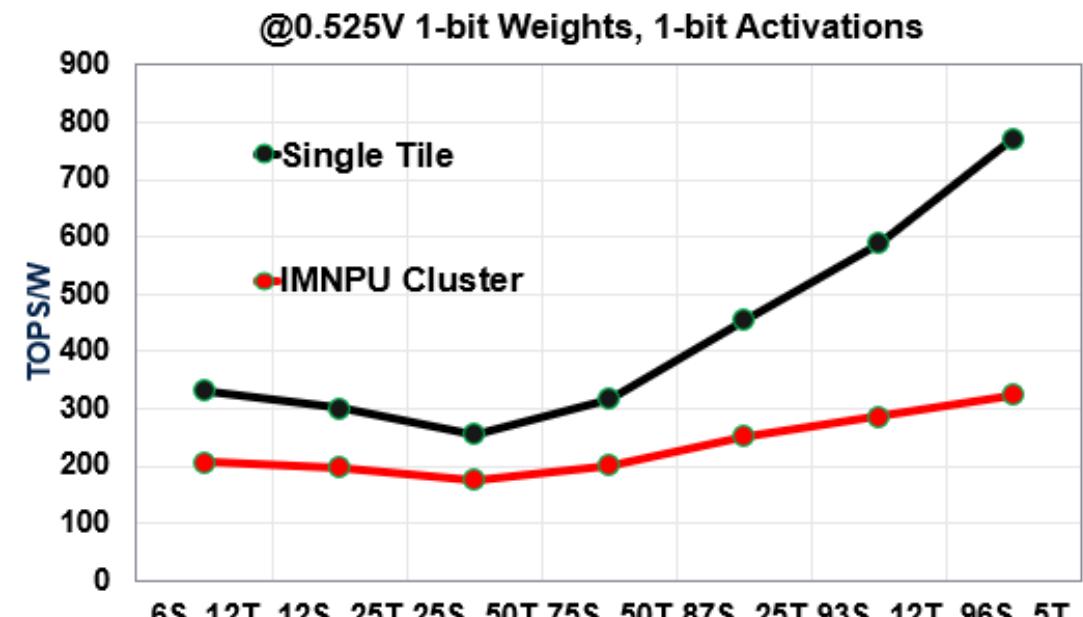
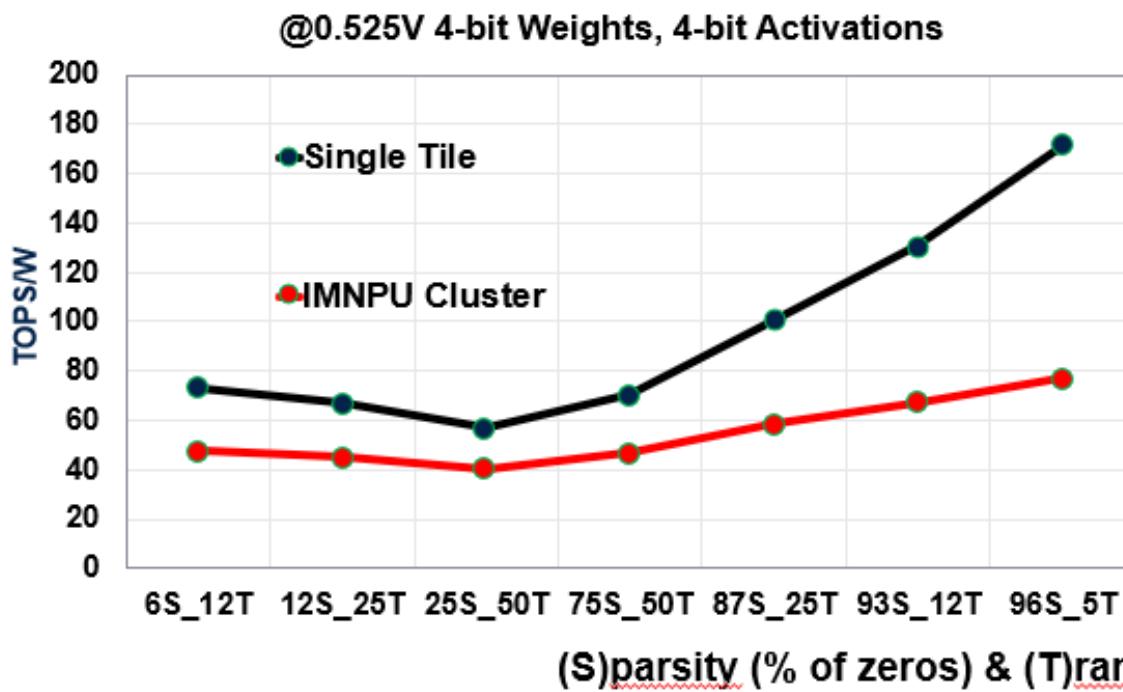
| configuration                 | MACSS/inf | Inf/sec | IMNPU Power   | Total Power <sup>1</sup> | Battery endurance <sup>2</sup><br>(1/100 duty cycle) |
|-------------------------------|-----------|---------|---------------|--------------------------|--|
| 1 cluster @ 10 MHz, 0.0V FBB  | 1.25 GOPS | 10      | 267 <u>uW</u> | 567 <u>uW</u>            | 363 days   |
| 8 clusters @ 400MHz, 0.3V FBB | 12.5 GOPS | 30      | 8.0 <u>mW</u> | 12.0 <u>mW</u>           |  |

(1) Estimated power includes a portion of shared memory, IOs, clock, and external sensor interface

(2) 6000 mA/h battery capacity assumed (e.g., 2 AA 1.5v batteries)

# Measurements

## Energy efficiency for tile and IMNPU



Charts show the trend of diminishing return of TOPS/W gains for a single DIMC instance not translating to proportional TOPS/W improvement at the IMNPU cluster level due to the scalar processing and data movement overheads.

# Our technology starts with You



Find out more at [www.st.com](http://www.st.com)

© STMicroelectronics - All rights reserved.

ST logo is a trademark or a registered trademark of STMicroelectronics International NV or its affiliates in the EU and/or other countries.

For additional information about ST trademarks, please refer to [www.st.com/trademarks](http://www.st.com/trademarks).

All other product or service names are the property of their respective owners.