**PULP PLATFORM**
Open Source Hardware, the way it should be!

# *Open-Source Heterogeneous Computing with Cluster-Coupled Accelerators: A Neural Engine Case Study*

Arpan Prasad[*], Gianna Paulin[*], **Yvan Tortorella[+]**, Luca Bertaccini[*], Luca Benini[*+], Francesco Conti[+]

[*]Integrated Systems Laboratory, ETH Zürich
[+]Energy-Efficient Embedded Systems Laboratory, University of Bologna

{prasadar, pauling, lbertaccini, lbenini}@iis.ee.ethz.ch
{yvan.tortorella,f.conti}@unibo.it

ETH zürich

http://pulp-platform.org     @pulp_platform     https://www.youtube.com/pulp_platform

# PULP includes Cores+Interco+IO+__HWPE__ → Open Platform
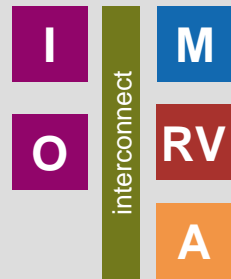
## RISC-V Cores

| RI5CY 32b | Ibex 32b | Snitch 32b | Ariane + Ara 64b |
|---|---|---|---|

## Peripherals

| JTAG | SPI |
|---|---|
| UART | I2S |
| DMA | GPIO |

## Interconnect

- Logarithmic interconnect
- APB – Peripheral Bus
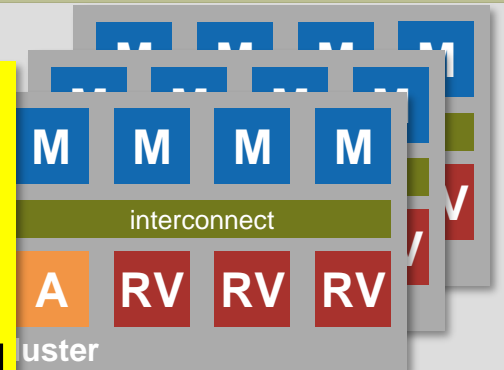- AXI4 – Interconnect

## Platforms

**Open-Source Hardware > Open-Source ISA:**
- core implementations
- system-level integration
- peripherals, __accelerators__

__Everything Designed in System Verilog__

**Single Core**
- PULPissimo

**Multi-core**
- Mr. Wolf
- Vega

**Multi-cluster**
- HERO
- Open Piton

IOT → HPC

## Hardware Processing Elements

| NE (DNN) | NTX (ML) | FFT | RedMulE (MatMul) |
|---|---|---|---|

PULP is silicon proven

*37 SoCs & counting*

*(25% with HWPEs)*

http://asic.ethz.ch

# PULP Cluster = Parallel RISC-V + Shared-L1 Scratchpad

**Tightly Coupled Data Memory**

| Mem | Mem | Mem | Mem | Mem |
|-----|-----|-----|-----|-----|
| Mem | Mem | Mem | Mem | Mem |

**bank interleaving** to maximize available bandwidth in typical parallel computing scenarios

**Cluster Interconnect**

| RISC-V core | RISC-V core | RISC-V core | RISC-V core |

**Cluster**

# DMA for data transfer + EU to coordination cores



**Tightly Coupled Data Memory**

Mem  Mem  Mem  Mem  Mem

Mem  Mem  Mem  Mem  Mem

DMA

L2 Mem

interconnect

Cluster Interconnect

Event Unit

RISC-V core  RISC-V core  RISC-V core  RISC-V core
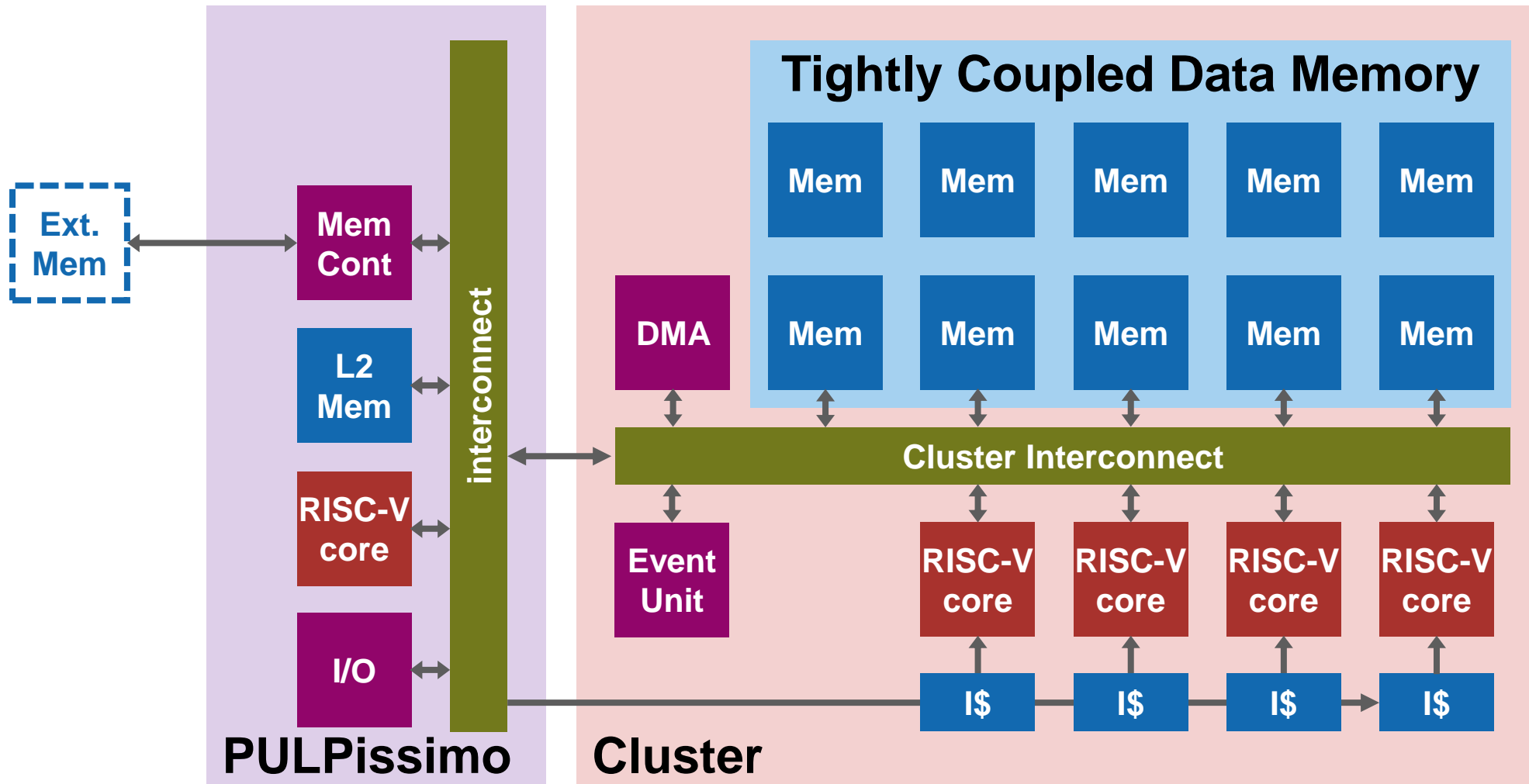
**Cluster**
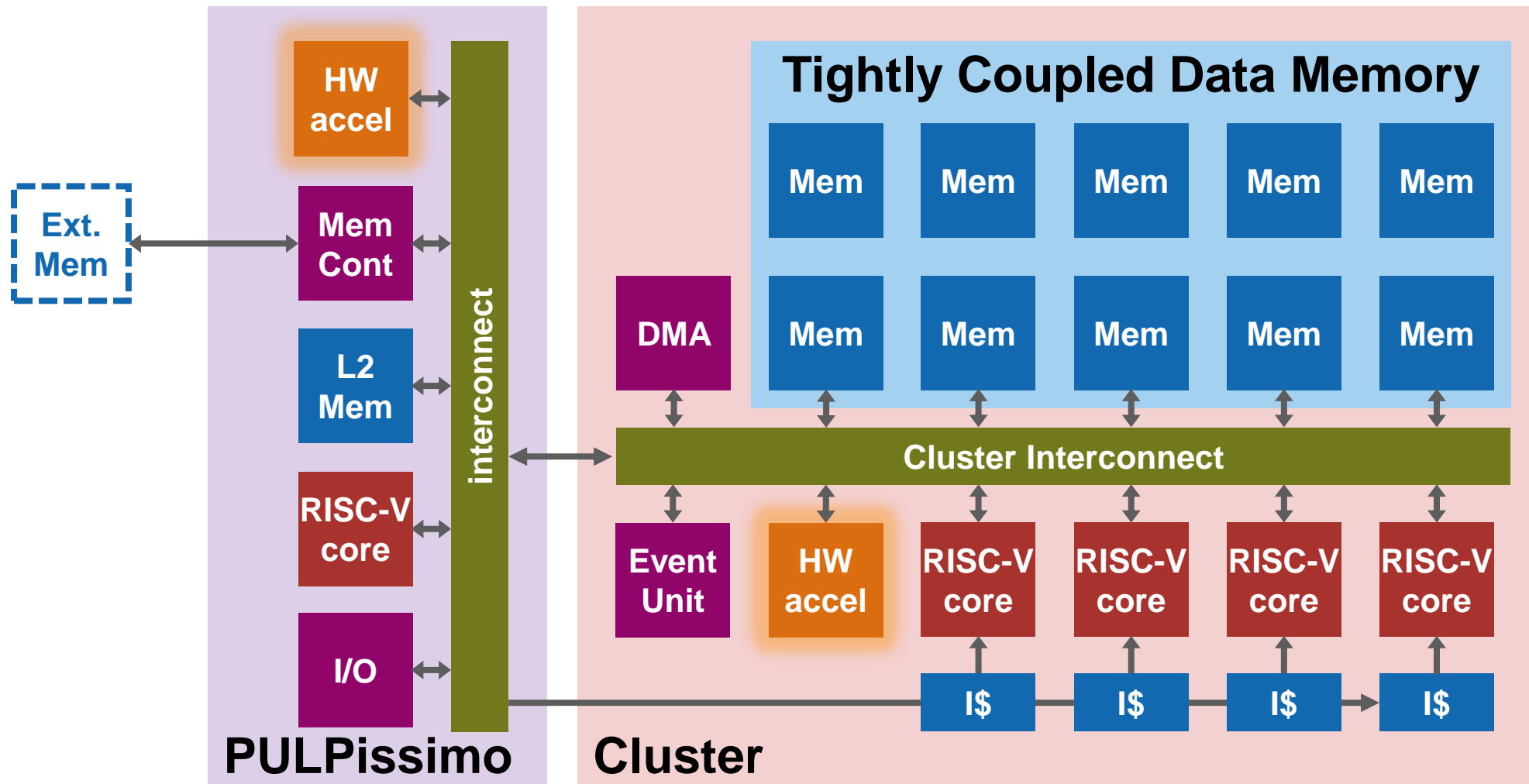
# There is a (shared) instruction cache that fetches from L2
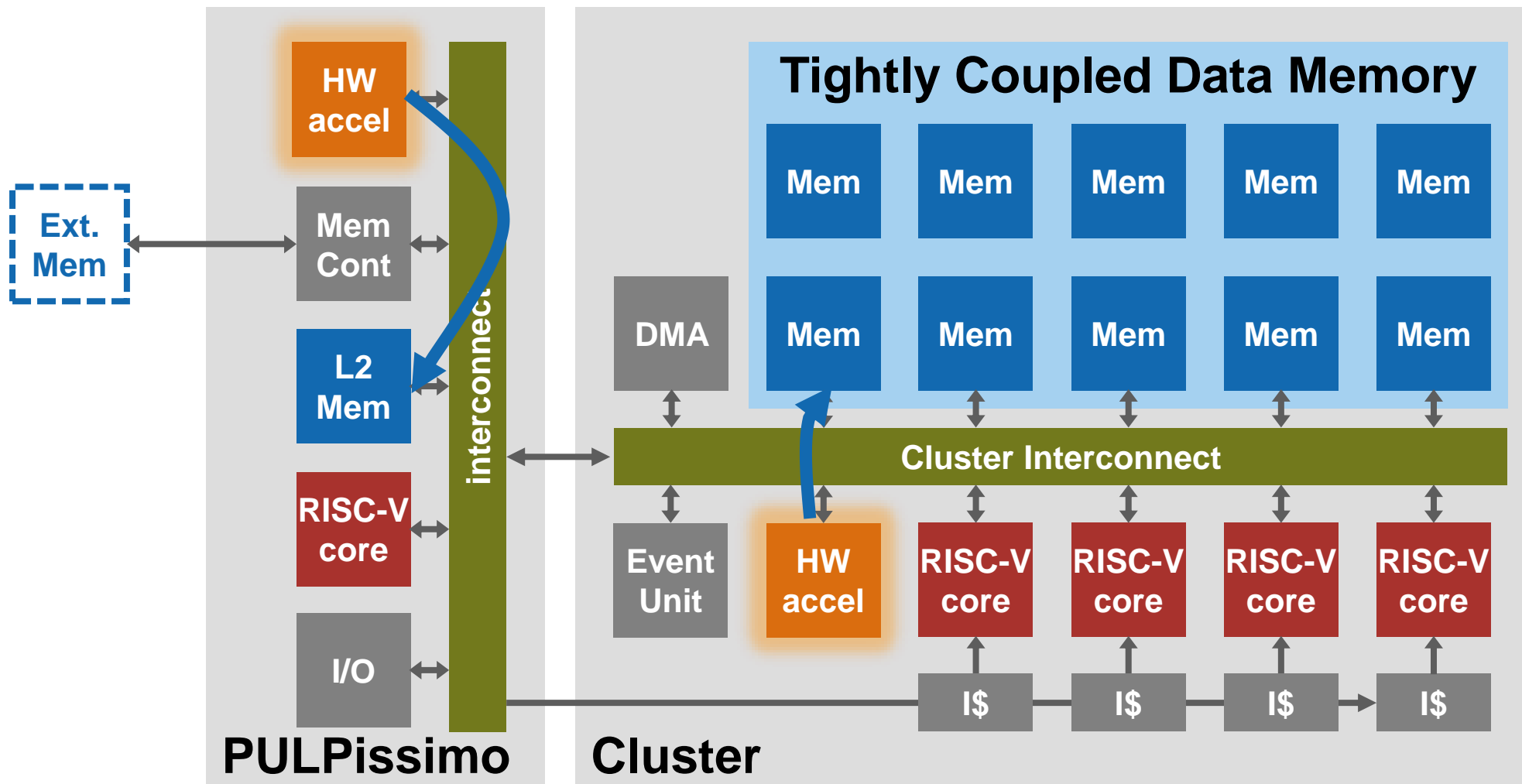
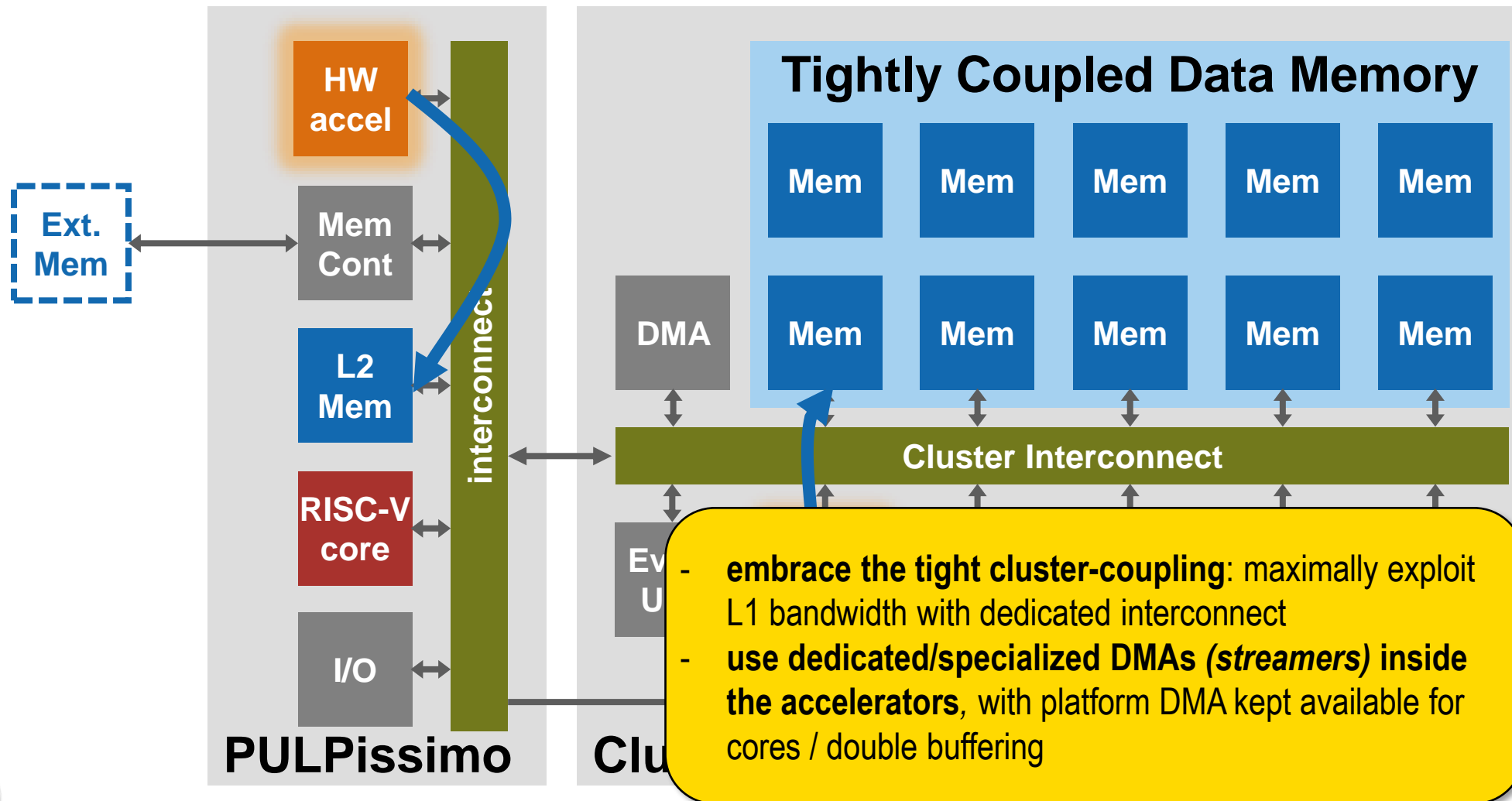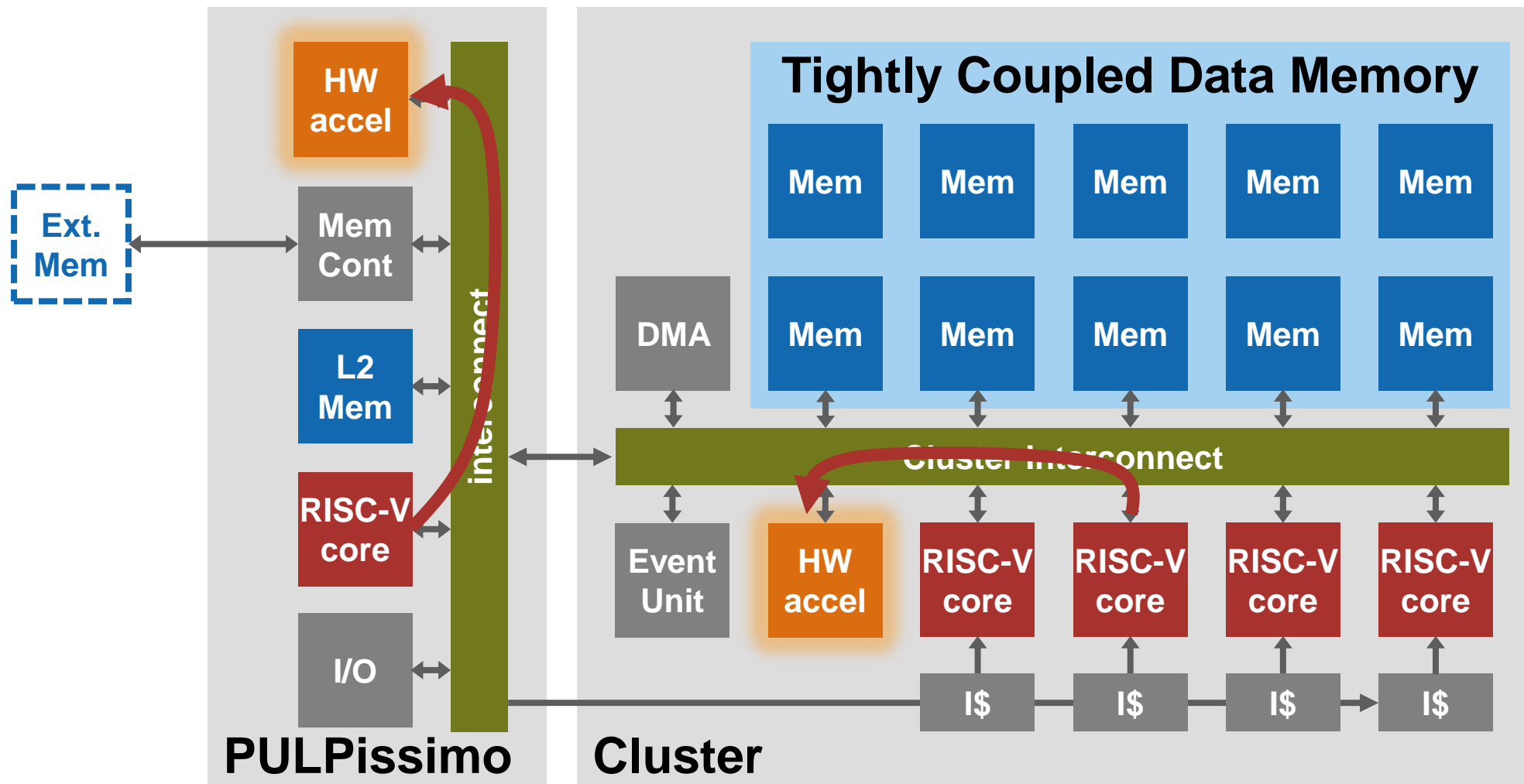# Microcontroller System (PULPissimo) for I/O

# How to couple HW Accerators?

# 1. High-Bandwidth, Latency-Tolerant access to Memory

# 1. High-Bandwidth, Latency-Tolerant access to Memory



**Tightly Coupled Data Memory**

- **embrace the tight cluster-coupling**: maximally exploit L1 bandwidth with dedicated interconnect
- **use dedicated/specialized DMAs** *(streamers)* **inside the accelerators**, with platform DMA kept available for cores / double buffering

# 2. Low-Overhead SW Control + Communication

# 2. Low-Overhead SW Control + Communication



- **memory-mapped control** with regular LD/ST like any other peripheral in the system (DMAs, I/O)
- standardized **ctrl interface** between different accelerators with common *(trigger, status, ...)* and specialized commands

# 3. Standardized Interface → Redesign only Datapath+Ctrl



**PULPissimo** block: Data Path + Intf, interconnect, Ext. Mem, Mem Cont, L2 Mem, RISC-V core, I/O

**Cluster** block: Tightly Coupled Data Memory (Mem × 11), DMA, Cluster Interconnect, Event Unit, Intf / Data Path, RISC-V core × 4, I$ × 4

# 1+2+3 = PULP Hardware Processing Engines (HWPEs)



**Tightly Coupled Data Memory**

PULPissimo: HWPE, Mem Cont, L2 Mem, RISC-V core, I/O, interconnect, Ext. Mem

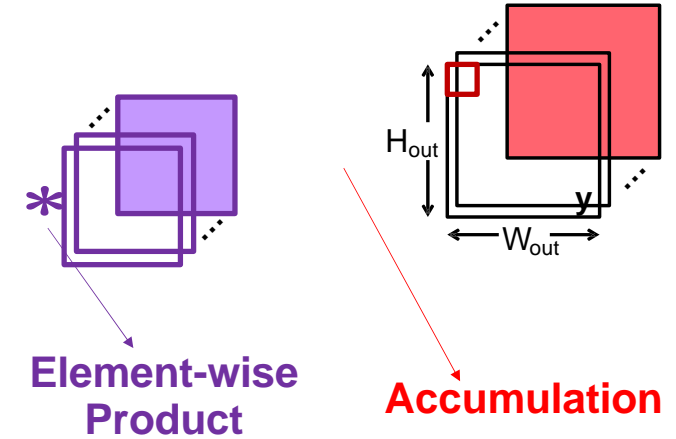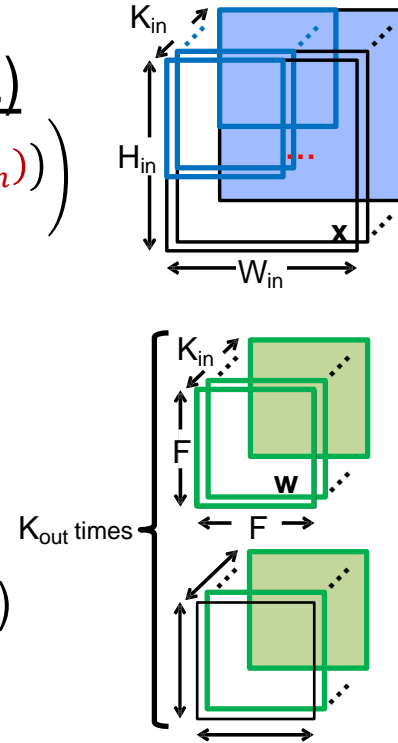Cluster: DMA, Event Unit, HWPE, RISC-V core, Cluster Interconnect, I$

# Case Study: the Neural Engine

## Flexible weight precision Neural Engine (NE)

$$\mathbf{y}(k_{out}) = \boldsymbol{quant}\left( \sum_{i=0..M} \sum_{filter} \sum_{k_{in}} 2^i \big( \mathbf{W_{bin}}(k_{out}, k_{in}, i) \otimes \mathbf{x}(k_{in}) \big) \right)$$
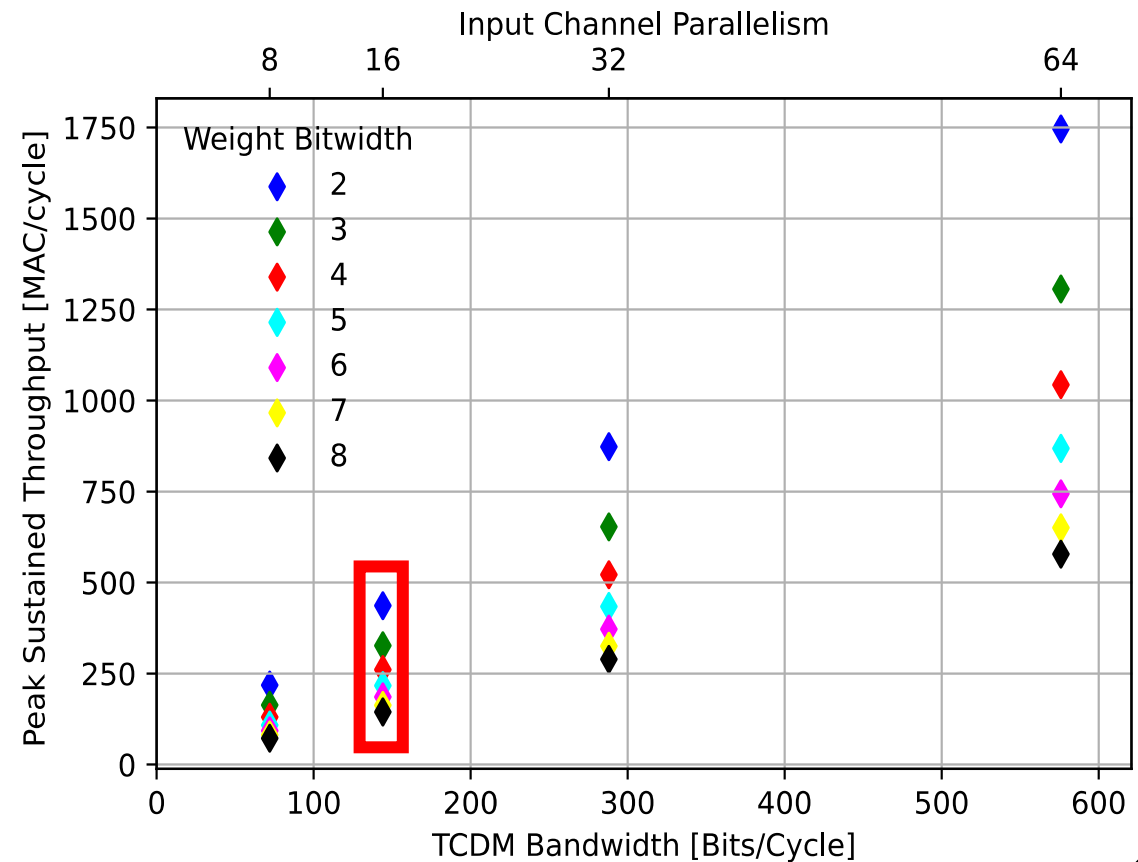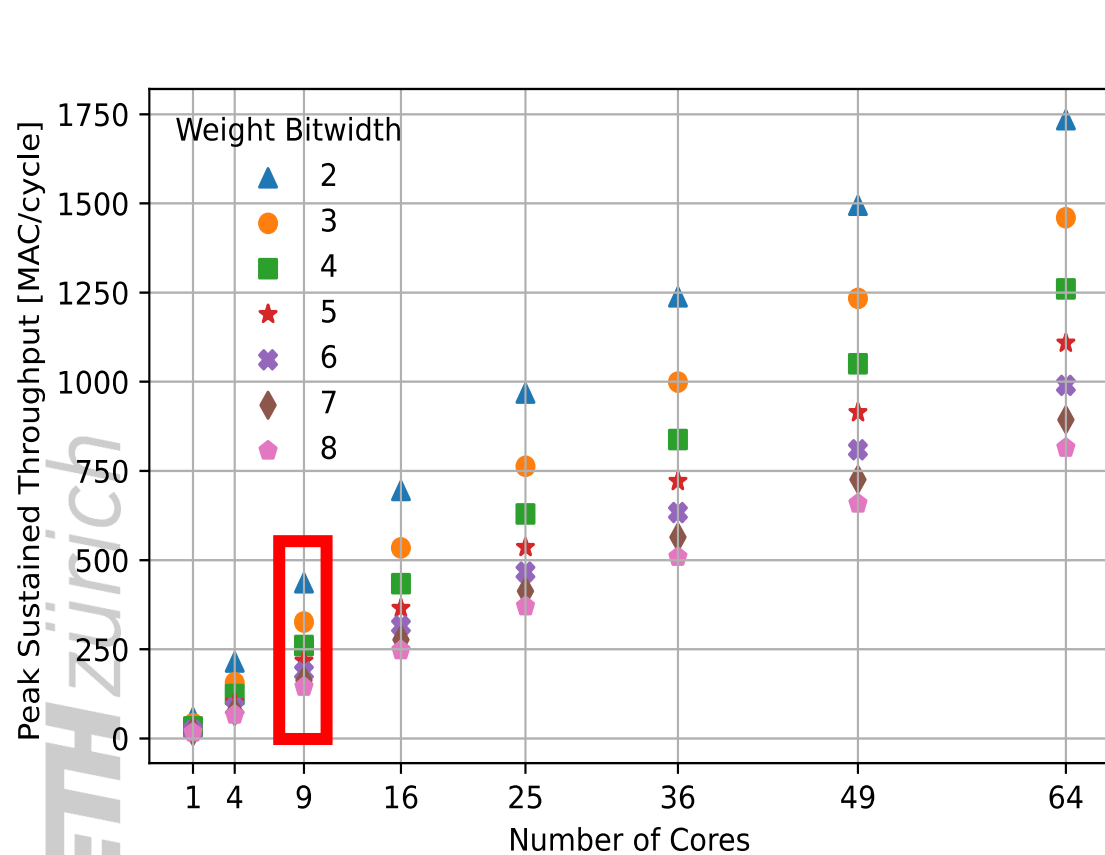
Key NE idea:

- many **cores**, each dedicated to the receptive field of 1 pixel across **OCP** out channels and **ICP** input channels
- splitting weights in **single bitwise** contributions: flexible weight prec (2-8b), fixed activation prec (8b)
- multiple iterations guided by SW tiling
- realized with **1x8b multipliers + adder trees**



**Element-wise Product**

**Accumulation**

# NE Performance Scalability (CONV 3x3)



Performance scalable by changing **number of cores** and **ICP** (design time)
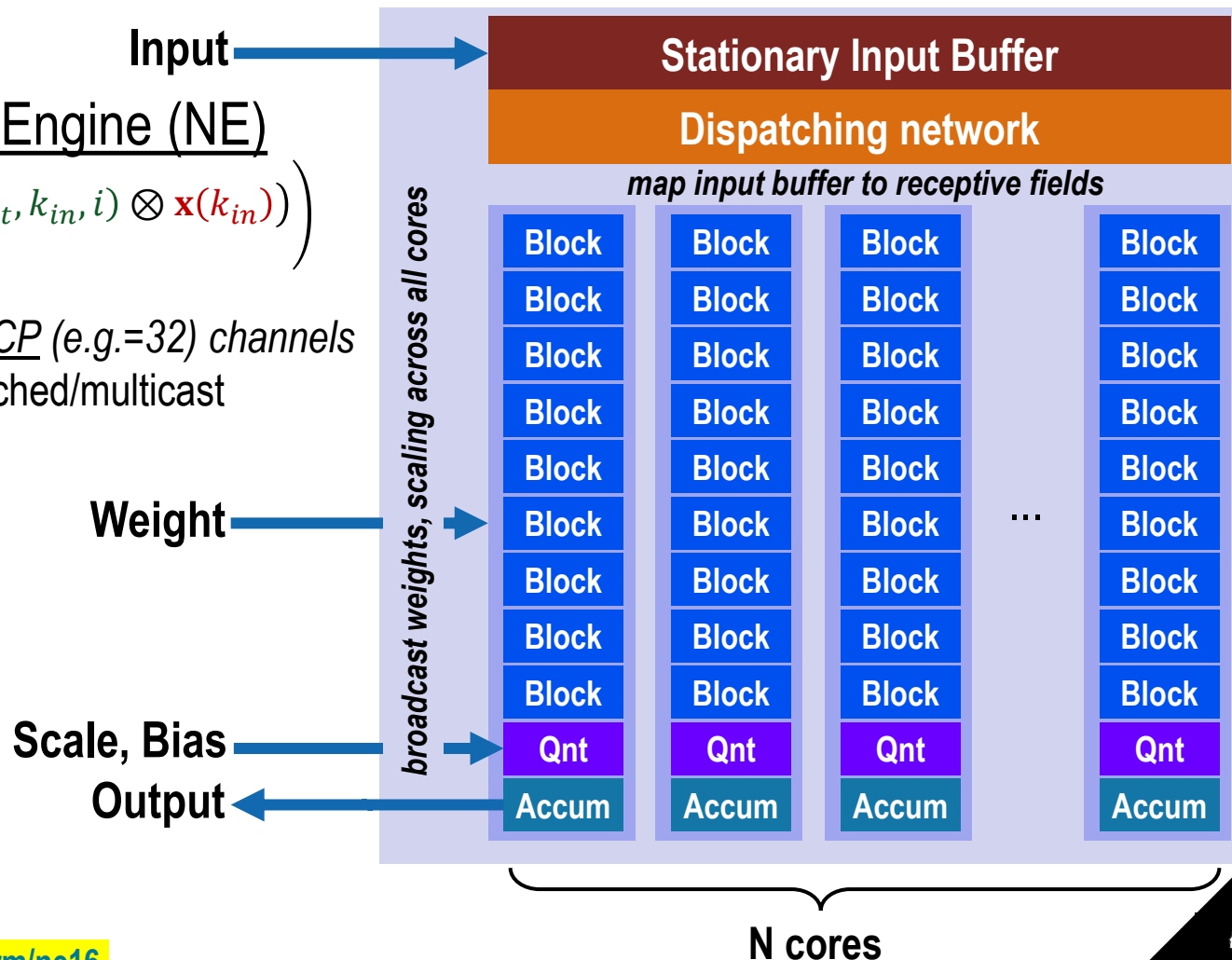or **weight bitwidth** (at run time)

# Case Study: the Neural Engine

## Flexible weight precision Neural Engine (NE)

$$\mathbf{y}(k_{out}) = \boldsymbol{quant}\left(\sum_{i=0..M}\sum_{filter}\sum_{k_{in}} 2^i\big(\mathbf{W_{bin}}(k_{out},k_{in},i) \otimes \mathbf{x}(k_{in})\big)\right)$$

**1 core = receptive field of 1 output px** across *OCP (e.g.=32) channels*
weights broadcast between all cores inputs dispatched/multicast
according to mode (3x3, 1x1, DepthWise)

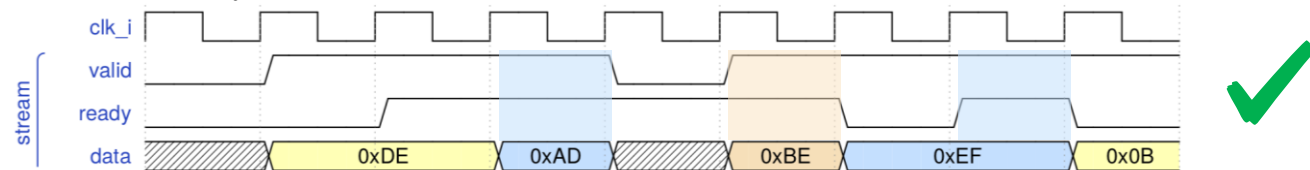**N cores, 9 blocks per core**
each block contributes
for *ICP (e.g.=32)*

Input

Weight

Scale, Bias

Output



**Stationary Input Buffer**

**Dispatching network**

*map input buffer to receptive fields*

*broadcast weights, scaling across all cores*

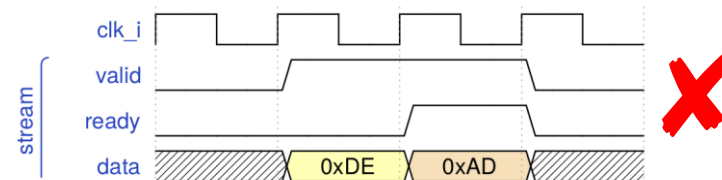| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | … | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Block | Block | Block | | Block |
| Qnt | Qnt | Qnt | | Qnt |
| Accum | Accum | Accum | | Accum |

**N cores**

# Lightweight, Latency-Tolerant *HWPE-Streams*

**Monodirectional, minimal** (only handshake + data + optional strobe), no assumptions on content. Positional information carried by order of data packets instead of address → **meant to represent data inside accelerator**
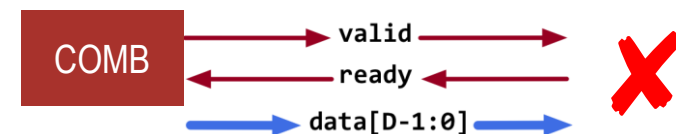
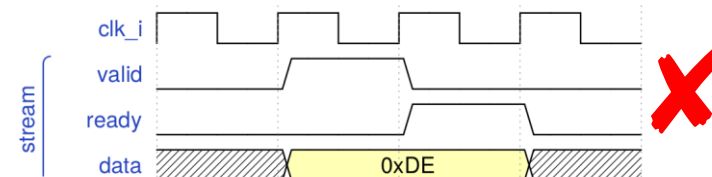**1.** handshake happens when valid & ready = 1

**2.** data/strobe can only change 1) following a handshake, 2) when valid is deasserted

**3.** assertion of valid (0 to 1) cannot depend combinationally on ready

**4.** deassertion of valid (1 to 0) can only happen in the cycle after a handshake
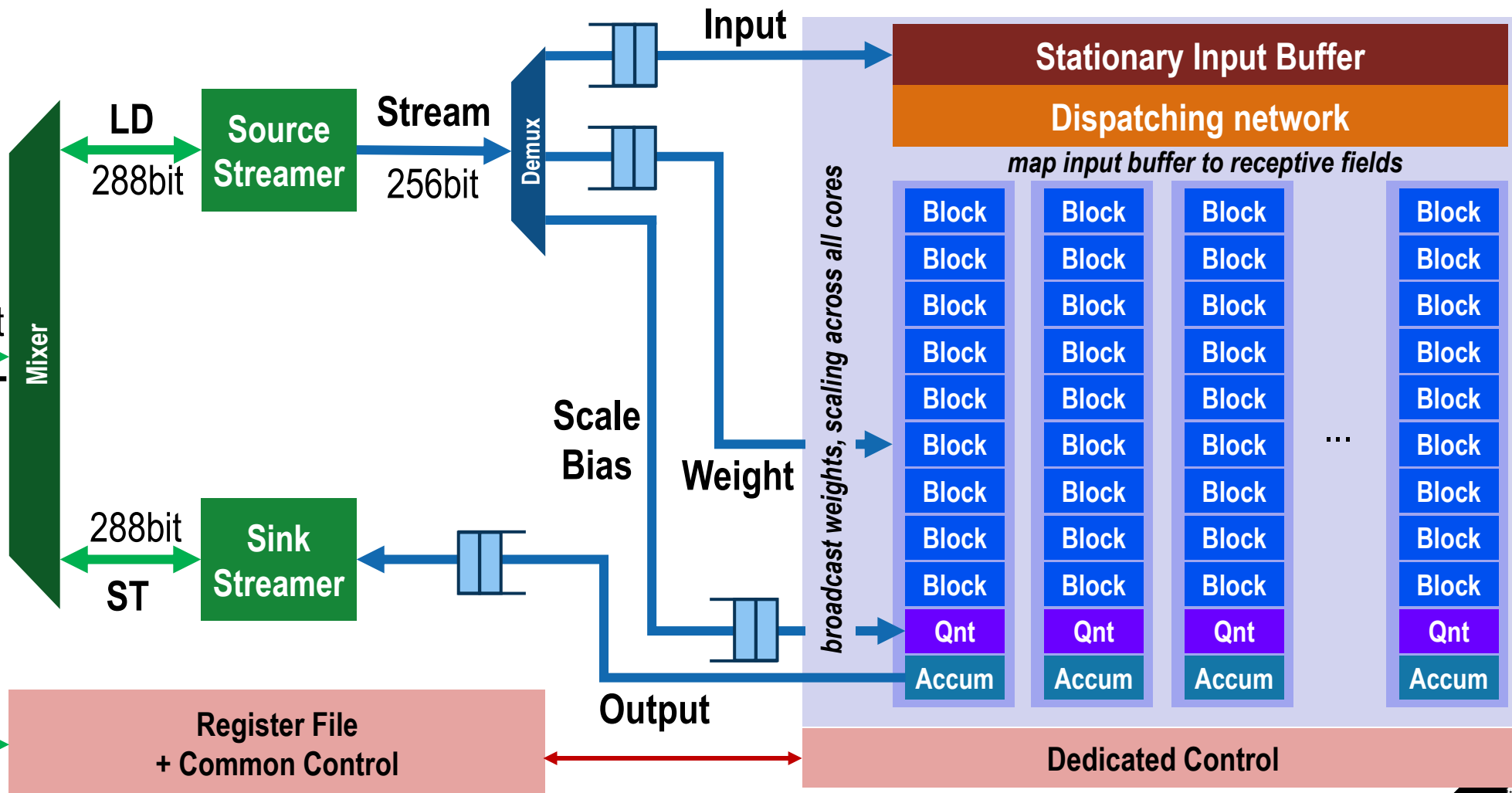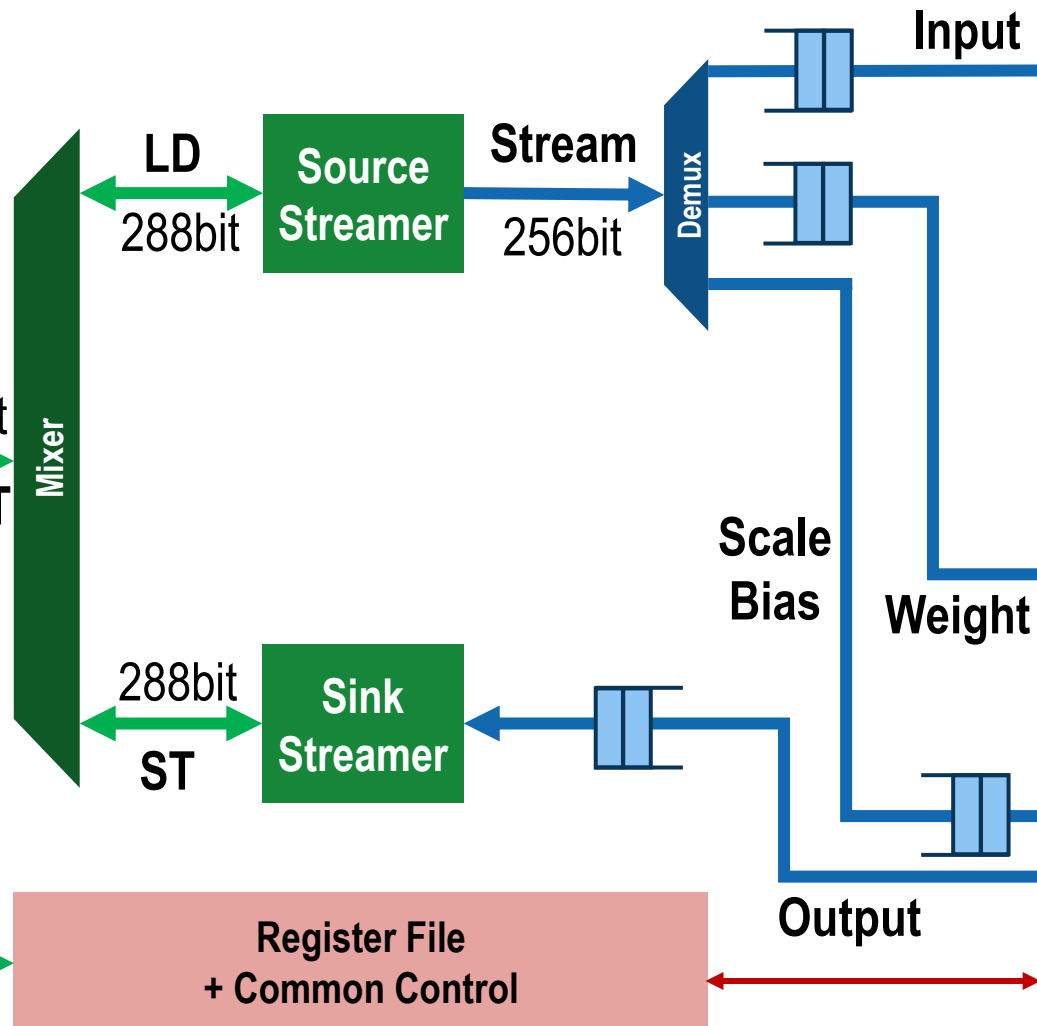
# Case Study: the Neural Engine

# Case Study: the Neural Engine



**Fully built with open-source HWPE IPs**
***https://hwpe-doc.readthedocs.io***

Interface designed by <u>composition</u> of open-source IPs for control + memory interfaces through ***HWPE-streams*** → specialized **streamers** convert from streams to mem (TCDM)
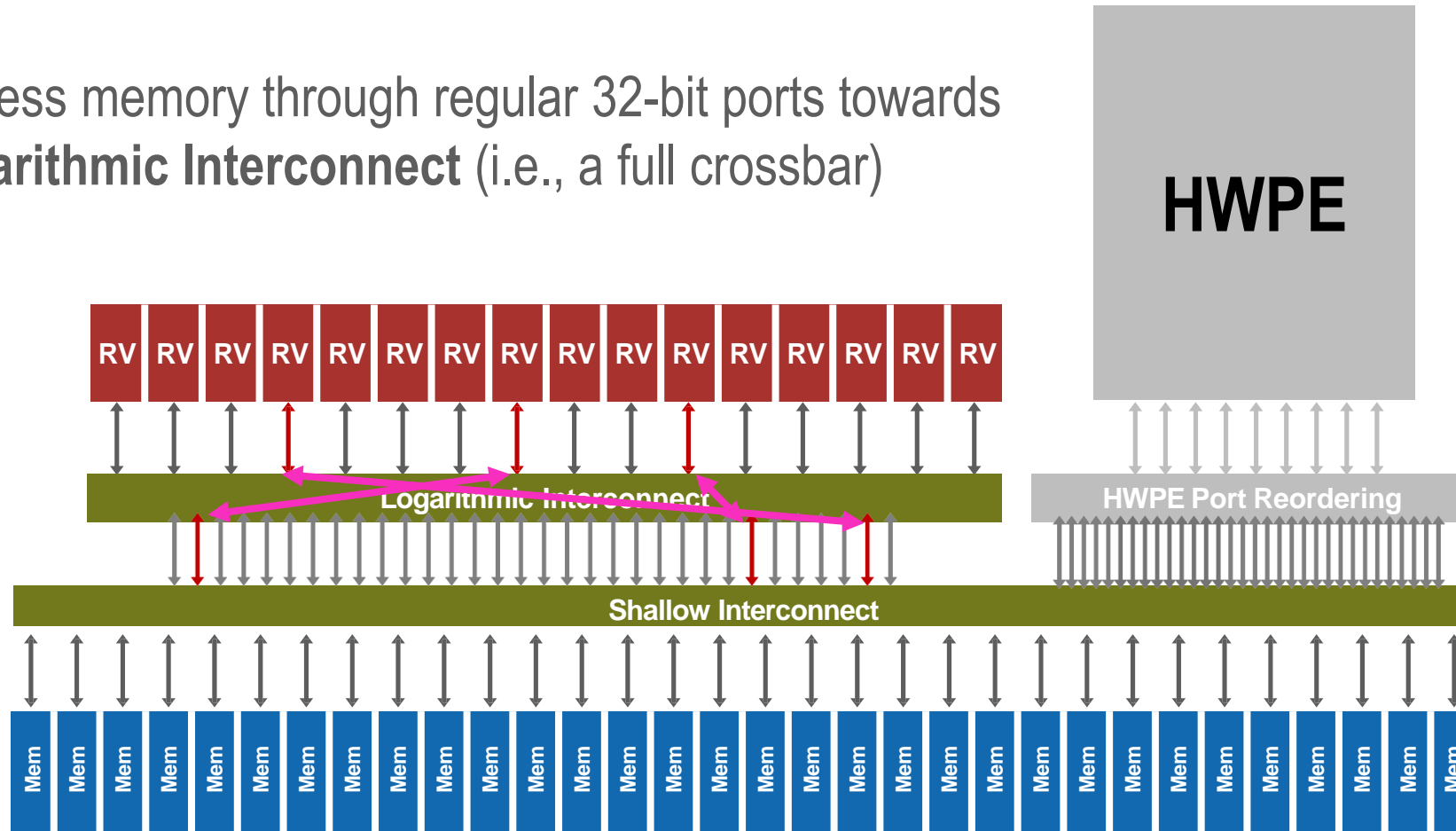
Memory access capabilities <u>equivalent or superior</u> to SW (misaligned access, delayed stores, high bandwidth) on part or all of the memory map

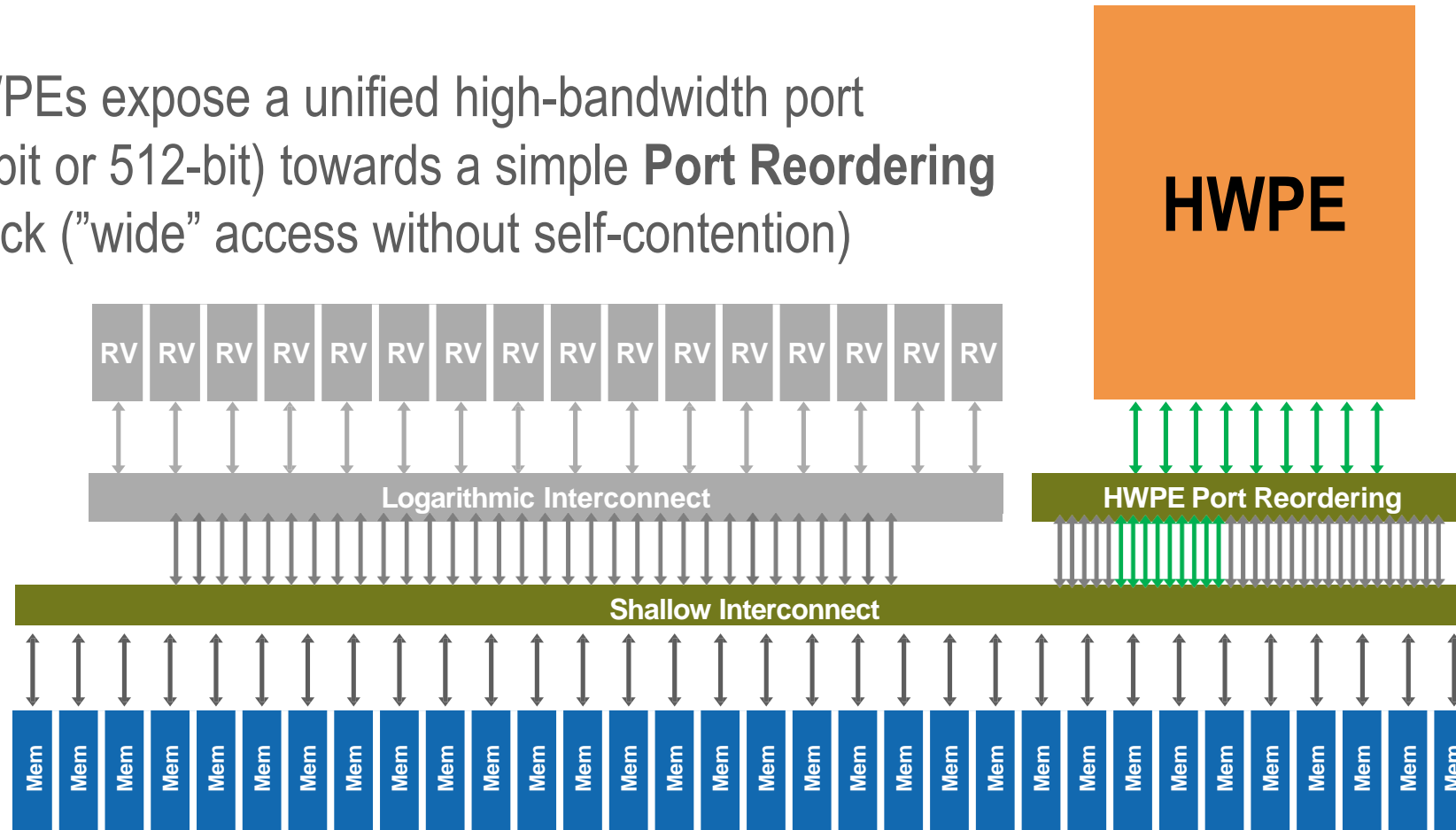# High-Bandwidth access: Heterogeneous Cluster Interconnect

Cores access memory through regular 32-bit ports towards **Logarithmic Interconnect** (i.e., a full crossbar)

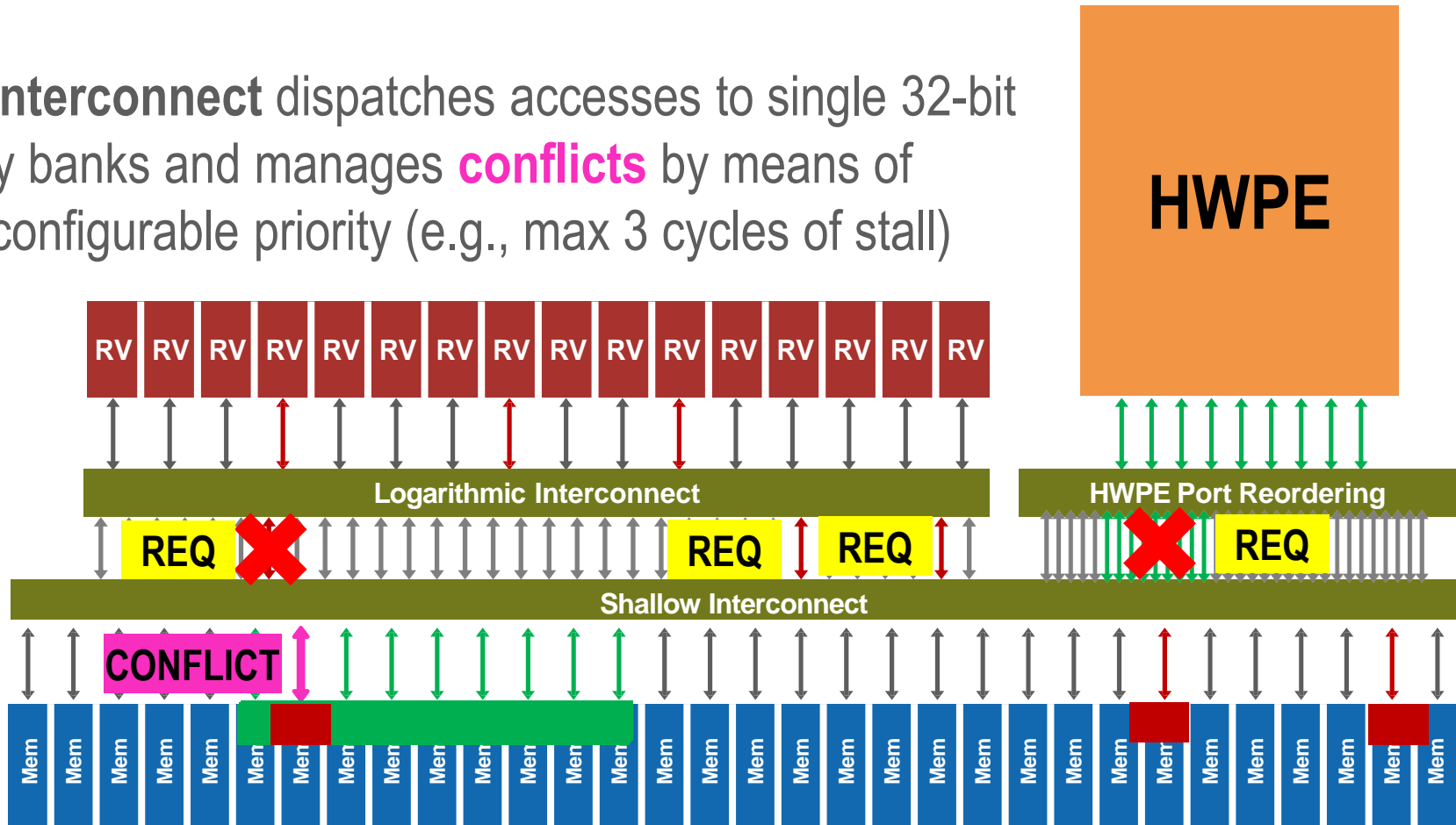# High-Bandwidth access: Heterogeneous Cluster Interconnect

HWPEs expose a unified high-bandwidth port
(e.g., 256-bit or 512-bit) towards a simple **Port Reordering**
block ("wide" access without self-contention)

# High-Bandwidth access: Heterogeneous Cluster Interconnect

A **Shallow Interconnect** dispatches accesses to single 32-bit memory banks and manages **conflicts** by means of rotating configurable priority (e.g., max 3 cycles of stall)
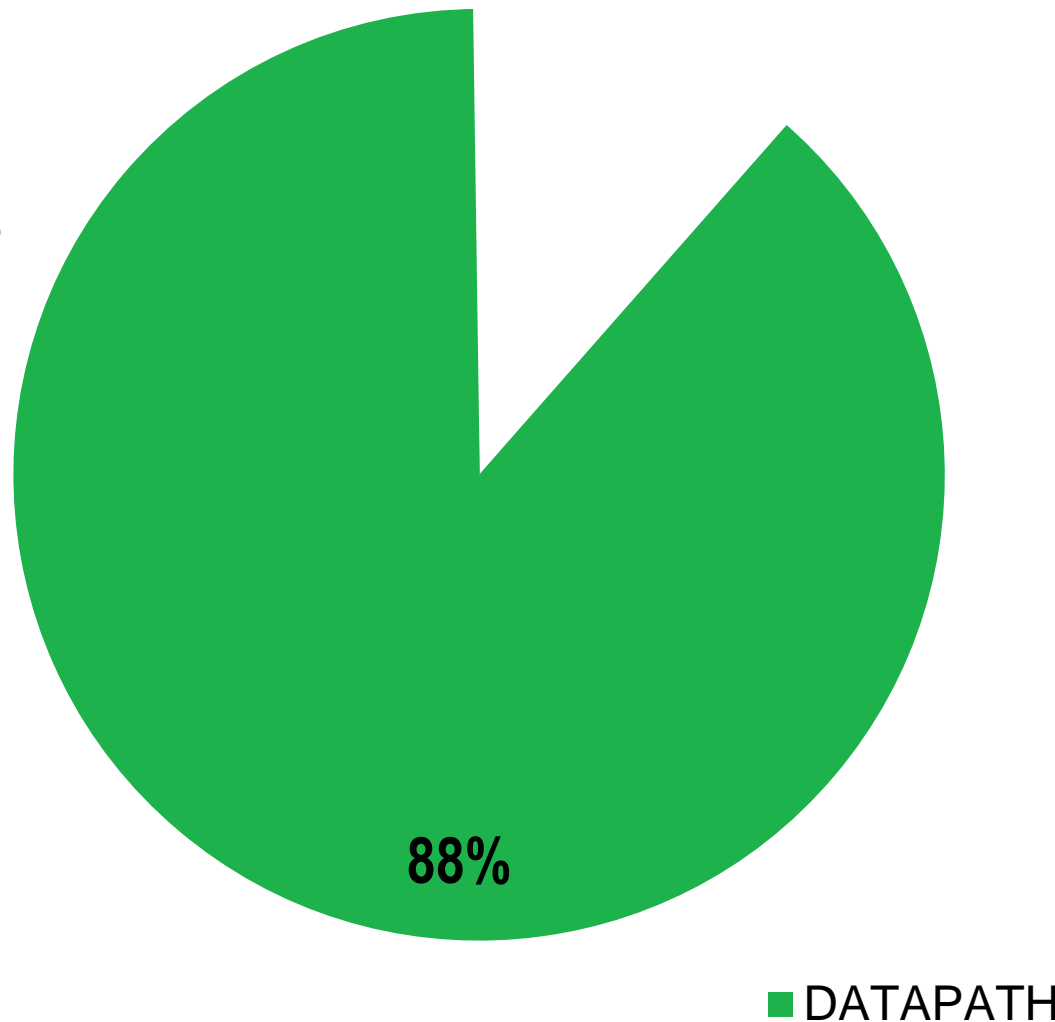
# Case study: NE area

Considering **NE**
with **ICP=16**, **OCP=32**, **9 cores**

Datapath accounts for **88%**
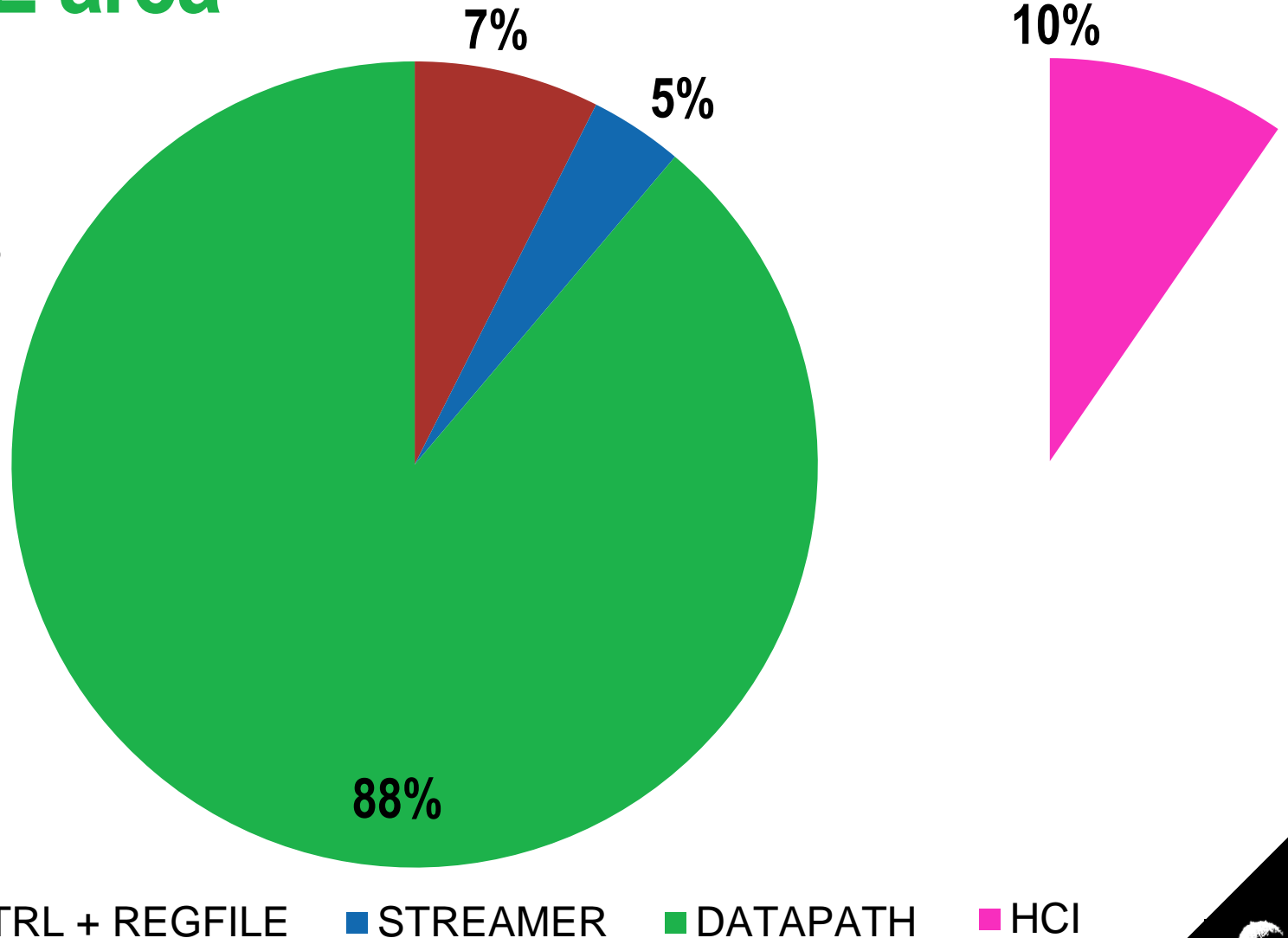of area of full accelerator

88%

■ DATAPATH

# Case study: NE area
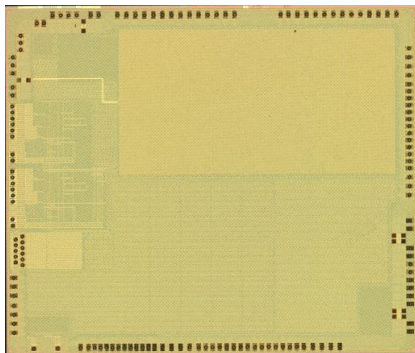
Considering **NE** with **ICP=16**, **OCP=32**, **9 cores**

Datapath accounts for **88%** of area of full accelerator

What about full HCI overhead? Log intc + HWPE intc → **~10%** of accelerator area!

7%

5%

10%

88%

■ CTRL + REGFILE  ■ STREAMER  ■ DATAPATH  ■ HCI

# HW-acceleration playground: success stories



*ETH zürich*

**GREENWAVES** TECHNOLOGIES

**Vega 22nm**, *ISSCC'21 with HW Conv Engine (UNIBO + ETHZ + GreenWaves)*



*ETH zürich*

Angelo Garofalo
Matteo Perotti
Luca Valente
Yvan Tortorella
Alessandro Nadalini
Alexandre Levisse
Riccardo Gandolfi
Carla Ohanesian
Francesco Conti
Andrew Simon
Pietro Maltoni
Marco Rios

**DARKSIDE**

**ECHOES** a PULP chip by Mattia Sinigaglia, Luca Valente, Luca Bertaccini, Angelo Garofalo, Francesco Conti, Davide Rossi

**Darkside 65nm**, *accepted to ESSCIRC'22* + **Echoes 65nm** *Master student tape-outs with many HWPEs: Transposer, FFT, Systolic Array, Depthwise Conv... (UNIBO+ETHZ)*

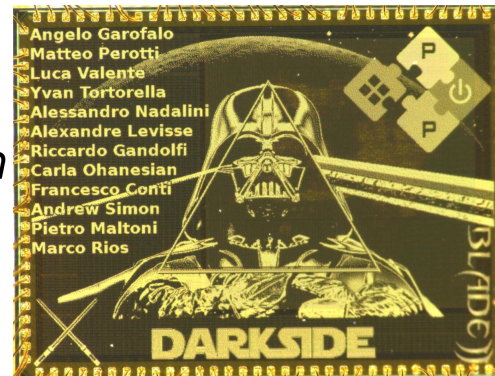# HW-acceleration playground: success stories

**ETH** zürich

GREENWAVES TECHNOLOGIES

**ETH** zürich

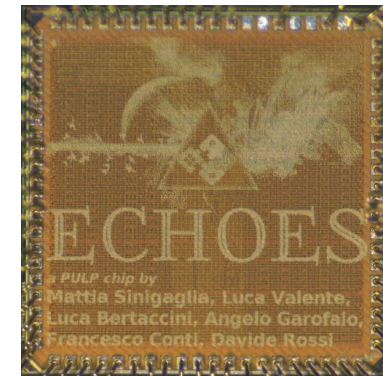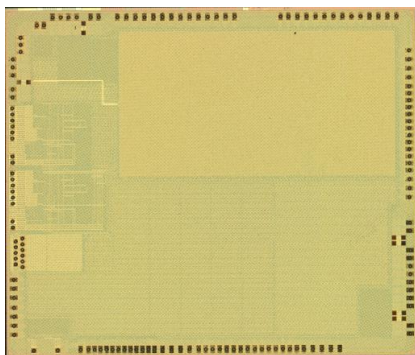**Vega 22nm**, *ISSCC'21 with HW Conv Engine (UNIBO + ETHZ + GreenWaves)*

**Darkside 65nm**, *accepted to ESSCIRC'22* + **Echoes 65nm**
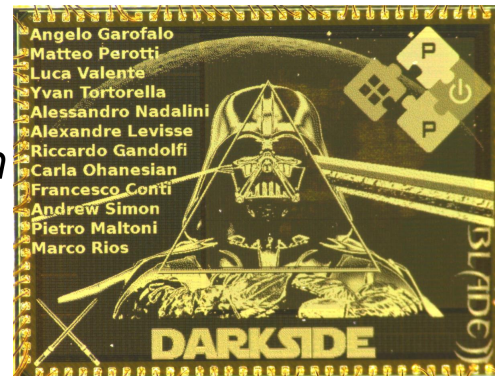*Master student tape-outs with many HWPEs: Transposer, FFT, Systolic Array, Depthwise Conv... (UNIBO+ETHZ)*

## Not only ETHZ/UNIBO:

GAP 9

**GAP9 SoC (commercial)**
*with NE16 (GreenWaves)*

KU LEUVEN imec
**DIANA 22nm**, *ISSCC'22 PULPissimo + DNN AiMC (KU Leuven + IMEC) [Ueyoshi et al., ISSCC'22]*

KU LEUVEN MAGICS
**TinyVers 22nm**, *VLSI'22 (to appear) PULPissimo + DNN HWPE (KU Leuven + Magics)*

# Take-home message

- **A set of templates, SystemVerilog IPs, and helpful tools that can be used to simplify accelerator development & usage**
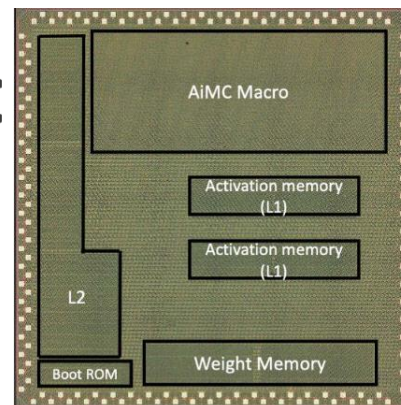
- **Main research activities now targeted at easing <u>usage</u> of HWPEs, not only developing new ones**

- **HWPE home page [https://hwpe-doc.readthedocs.io](https://hwpe-doc.readthedocs.io)**

- **Main PULP home page [https://pulp-platform.org](https://pulp-platform.org)**

- **Not only RTL: GVSOC supports HWPEs [https://github.com/pulp-platform/pulp-sdk](https://github.com/pulp-platform/pulp-sdk)**

# PULP
## Parallel Ultra Low Power

**HWPE Team:** Francesco Conti, Yvan Tortorella, Arpan Prasad, Luca Bertaccini, Gianna Paulin, Alessio Burrello, Luka Macan, Luca Benini
**And all the PULP team:** Luca Benini, Alessandro Capotondi, Alessandro Ottaviano, Alessio Burrello, Alfio Di Mauro, Andrea Borghesi, Andrea Cossettini, Andreas Kurth, Angelo Garofalo, Antonio Pullini, Arpan Prasad, Bjoern Forsberg, Corrado Bonfanti, Cristian Cioflan, Daniele Palossi, Davide Rossi, Fabio Montagna, Florian Glaser, Florian Zaruba, Francesco Conti, Georg Rutishauser, Germain Haugou, Gianna Paulin, Giuseppe Tagliavini, Hanna Müller, Luca Bertaccini, Luca Valente, Luca Colagrande, Manuel Eggimann, Manuele Rusci, Marco Guermandi, Matheus Cavalcante, Matteo Perotti, Matteo Spallanzani, Michael Rogenmoser, Moritz Scherer, Moritz Schneider, Nazareno Bruschi, Nils Wistoff, Pasquale Davide Schiavone, Paul Scheffler, Philipp Mayer, Robert Balas, Samuel Riedel, Sergio Mazzola, Sergei Vostrikov, Simone Benatti, Stefan Mach, Thomas Benz, Thorir Ingolfsson, Tim Fischer, Victor Javier Kartsch Morinigo, Vlad Niculescu, Xiaying Wang, Yichao Zhang, Frank K. Gürkaynak, all our past collaborators and many more that we forgot to mention
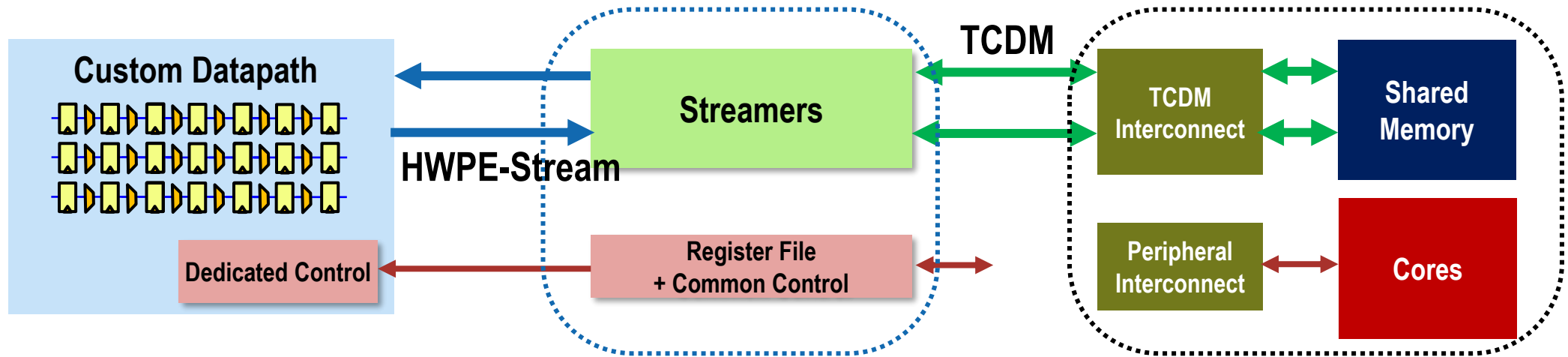
🌐 **http://pulp-platform.org** 🐦 **@pulp_platform**



Angelo Garofalo
Matteo Perotti
Luca Valente
Yvan Tortorella
Alessandro Nadalini
Alexandre Levisse
Riccardo Gandolfi
Carla Ohanesian
Francesco Conti
Andrew Simon
Pietro Maltoni
Marco Rios

BL(ADE)

DARKSIDE

# BACKUP SLIDES

# HWPEs are "first-class citizens" in the cluster



**Custom Datapath**

**HWPE-Stream**

**Dedicated Control**

**Streamers**

**Register File + Common Control**

**TCDM**

**TCDM Interconnect**

**Peripheral Interconnect**

**Shared Memory**

**Cores**

**Built with open-source HWPE IPs**
*https://hwpe-doc.readthedocs.io*

**Open-source PULP cluster**
*https://github.com/pulp-platform/pulp*

*Divide et Impera* design approach for HWPEs

Interface designed by <u>composition</u> of open-source IPs for control + memory interfaces through **HWPE-streams** → specialized **streamers** convert from streams to mem (TCDM)

Memory access capabilities <u>equivalent or superior</u> to SW (misaligned access, delayed stores, high) on part or all of the memory map
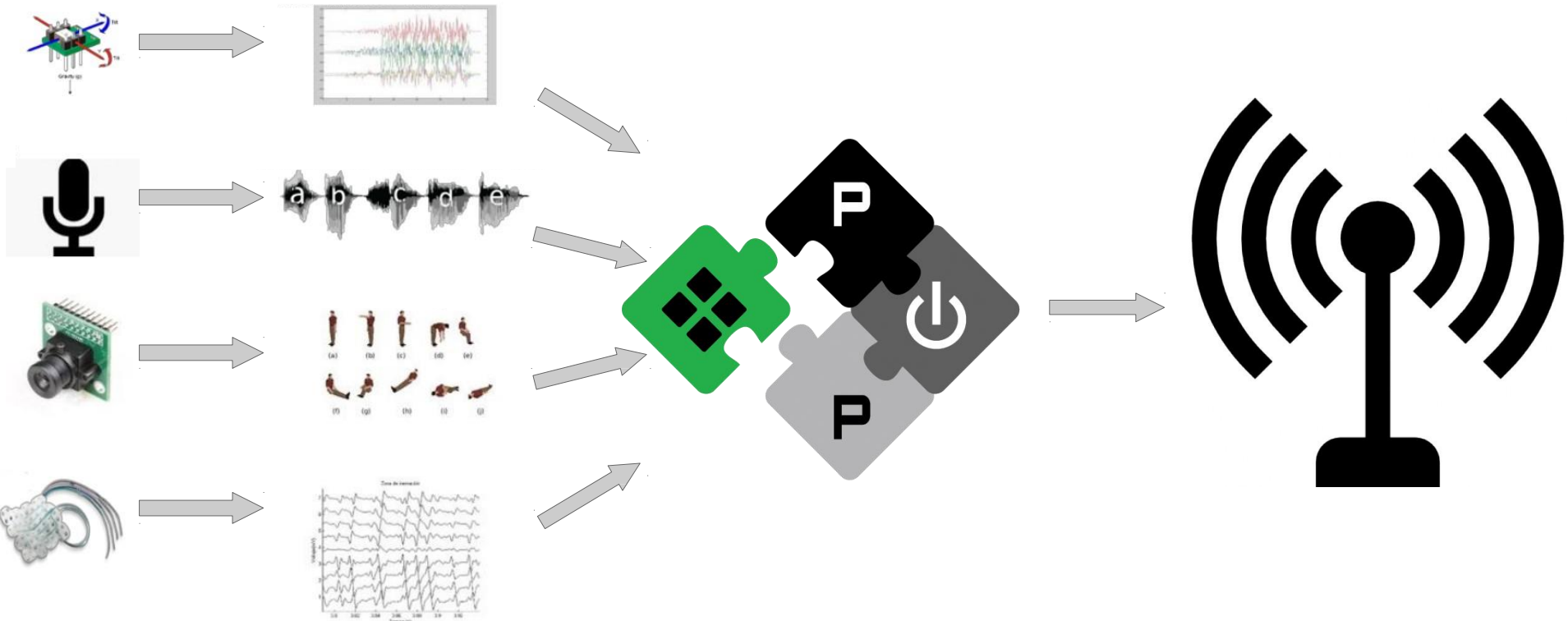
# Edge Processing

**Sensors**  **Signals**            **Processing**         **Transmission**


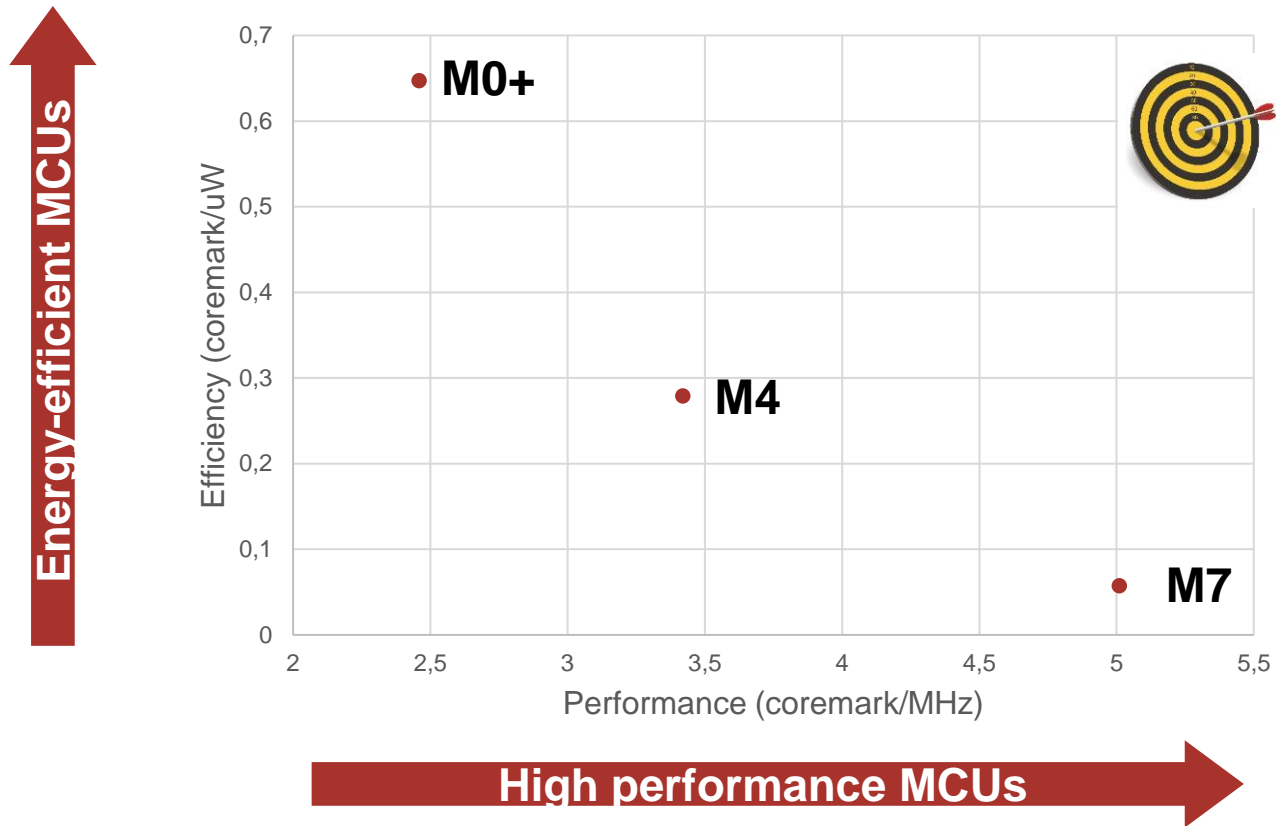
100μW - 1mW          1mW - 10mW          1mW (idle) - 50mW (active)

# Energy efficiency @ GOPS is the Challenge

ARM Cortex-M MCUs: M0+, M4, M7 (40LP, typ, 1.1V)*

**Energy-efficient MCUs** ↑

Efficiency (coremark/uW) vs Performance (coremark/MHz)

- M0+ (≈2.5, ≈0.65)
- M4 (≈3.4, ≈0.28)
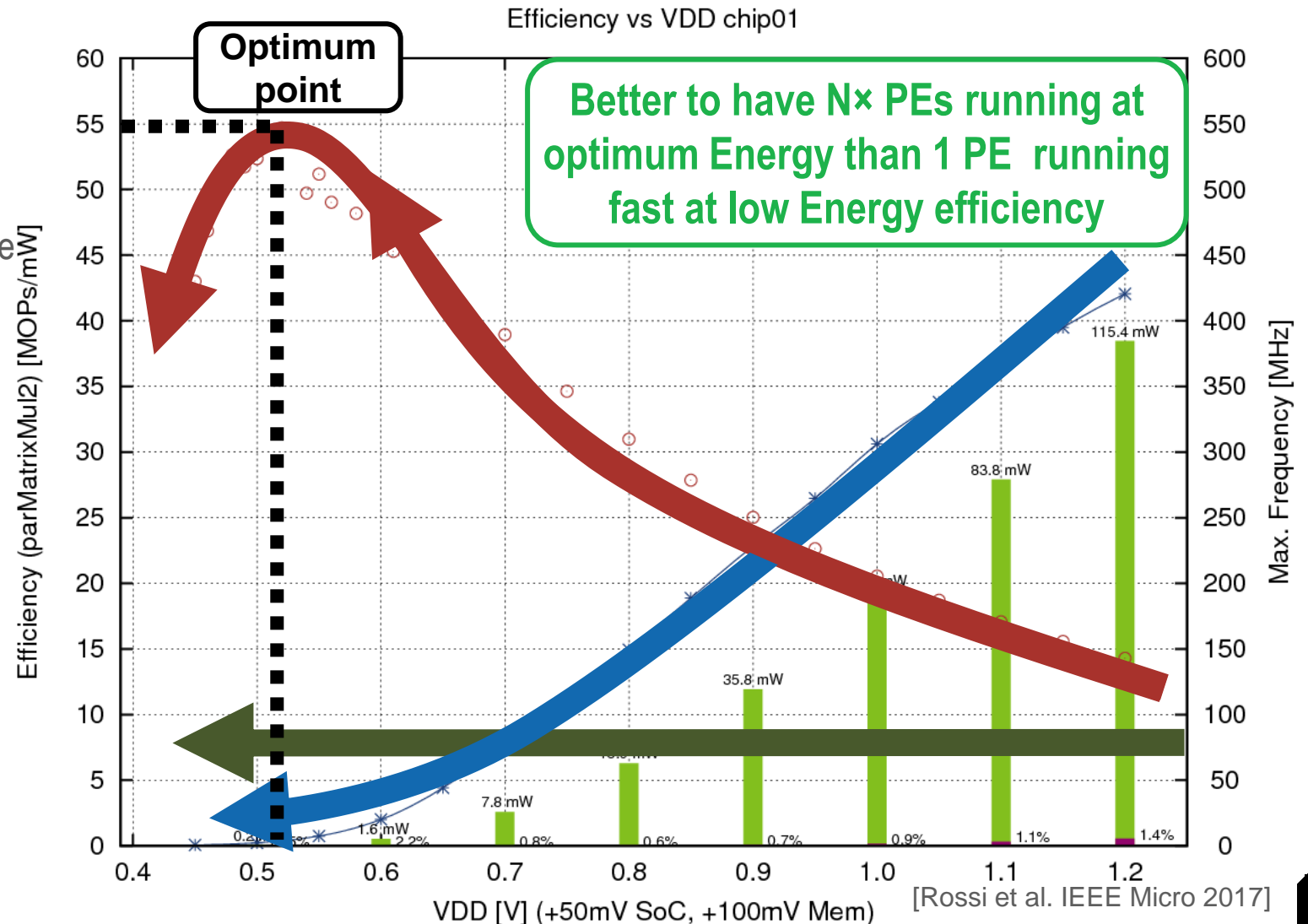- M7 (≈5, ≈0.06)

**High performance MCUs** →
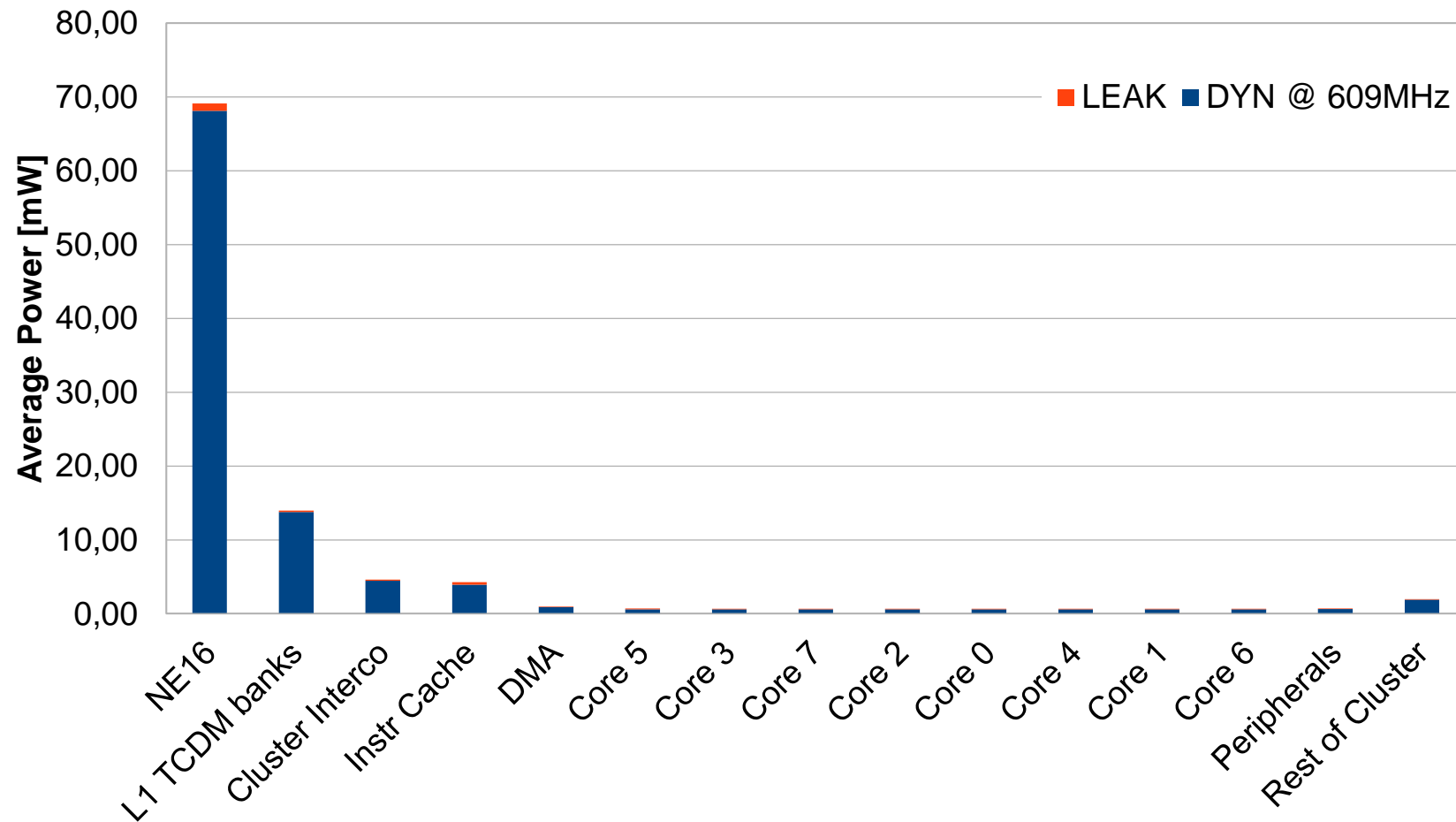
## How??

*data from ARMs web*

# Parallel & Accelerated + Near-threshold

- As VDD decreases, operating speed decreases

- However efficiency increases→ more work done per Joule

- Until leakage effects start to dominate

- Put more units in parallel to get performance up and keep them busy with a parallel workload

ML is massively parallel and scales well (P/S ↑ with NN size)



Efficiency vs VDD chip01

**Optimum point**

**Better to have N× PEs running at optimum Energy than 1 PE running fast at low Energy efficiency**

Efficiency (parMatrixMul2) [MOPs/mW]

Max. Frequency [MHz]

VDD [V] (+50mV SoC, +100mV Mem)

[Rossi et al. IEEE Micro 2017]

# NE Power (ICP = 16)

Memory-mapped, typically connected to Peripheral Interconnect or Peripheral Demux.
Executes a **queue** of **jobs** (typically 2 entries).

| Control Registers | | |
|---|---|---|
| | *trigger* | starts execution on HWPE, unlocks the ctrl |
| | *acquire* | starts a job offload on HWPE, returns a job ID handle, locks the ctrl |
| | *finished jobs* | returns no. of completed jobs |
| | *status* | each byte returns status of a job in the queue (1=working, 0=idle) |
| | *running job* | returns ID of the currently running job |
| | *soft clear* | clears the HWPE |
| | *offloader ID* | each byte returns ID+1 of a job in the queue |
| | *sw synch* | triggers an implementation-specific internal HWPE event (0 to 7 depending on written data) |
| **Generic Registers** | | |
| **Job-Dependent Registers** | | |

*Control registers used for standard HWPE control model, common to all HWPEs*

*Generic registers and Job-dependent registers used for runtime parameters specific of a HW accelerator. Generics are constant in all jobs*
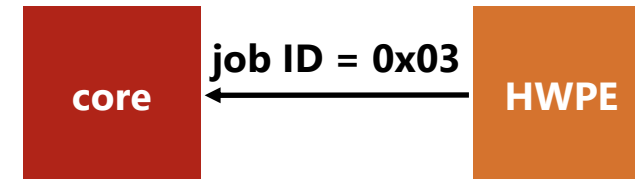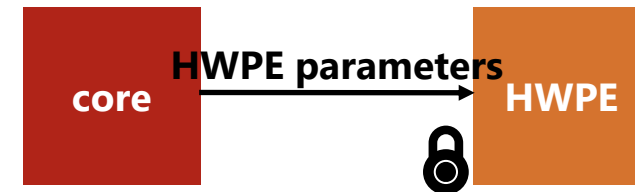
Memory-mapped, typically connected to Peripheral Interconnect or Peripheral Demux.
Executes a **queue** of **jobs** (typically 2 entries).

**Control Registers**

| | |
|---|---|
| *trigger* | starts execution on HWPE, unlocks the ctrl |
| *acquire* | starts a job offload on HWPE, returns a job ID handle, locks the ctrl |
| *finished jobs* | returns no. of completed jobs |
| *status* | each byte returns status of a job in the queue (1=working, 0=idle) |
| *running job* | returns ID of the currently running job |
| *soft clear* | clears the HWPE |
| *offloader ID* | each byte returns ID+1 of a job in the queue |
| *sw synch* | triggers an implementation-specific internal HWPE event (0 to 7 depending on written data) |

**Generic Registers**
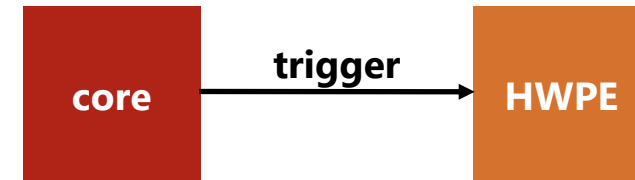
**Job-Dependent Registers**

**core** ← **job ID = 0x03** ← **HWPE**

Memory-mapped, typically connected to Peripheral Interconnect or Peripheral Demux.
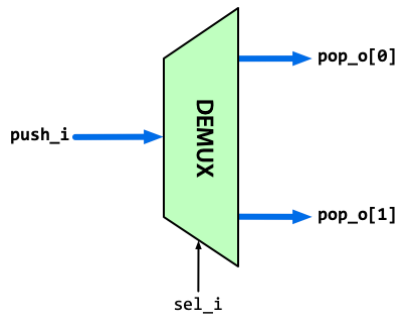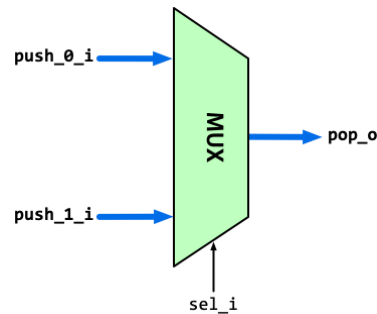Executes a **queue** of **jobs** (typically 2 entries).

**Control Registers**

| | | |
|---|---|---|
| | *trigger* | starts execution on HWPE, unlocks the ctrl |
| | *acquire* | starts a job offload on HWPE, returns a job ID handle, locks the ctrl |
| | *finished jobs* | returns no. of completed jobs |
| | *status* | each byte returns status of a job in the queue (1=working, 0=idle) |
| | *running job* | returns ID of the currently running job |
| | *soft clear* | clears the HWPE |
| | *offloader ID* | each byte returns ID+1 of a job in the queue |
| | *sw synch* | triggers an implementation-specific internal HWPE event (0 to 7 depending on written data) |
| **Generic Registers** | | |
| **Job-Dependent Registers** | | |

**core** — **HWPE parameters** → **HWPE**

# HWPE Control

Memory-mapped, typically connected to Peripheral Interconnect or Peripheral Demux.
Executes a **queue** of **jobs** (typically 2 entries).

| Control Registers | | |
|---|---|---|
| | *trigger* | starts execution on HWPE, unlocks the ctrl |
| | *acquire* | starts a job offload on HWPE, returns a job ID handle, locks the ctrl |
| | *finished jobs* | returns no. of completed jobs |
| | *status* | each byte returns status of a job in the queue (1=working, 0=idle) |
| | *running job* | returns ID of the currently running job |
| | *soft clear* | clears the HWPE |
| | *offloader ID* | each byte returns ID+1 of a job in the queue |
| | *sw synch* | triggers an implementation-specific internal HWPE event (0 to 7 depending on written data) |
| **Generic Registers** | | |
| **Job-Dependent Registers** | | |

core → **trigger** → HWPE

# *HWPE-Stream*: a lightweight protocol for streams

**Monodirectional**, **minimal** (only handshake + data + optional strobe), no assumptions on content. Positional information carried by *order* of data packets as opposed to *address*.
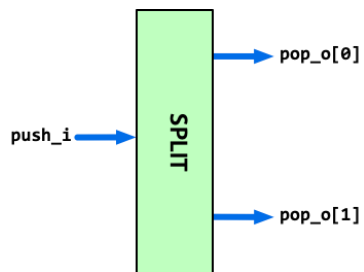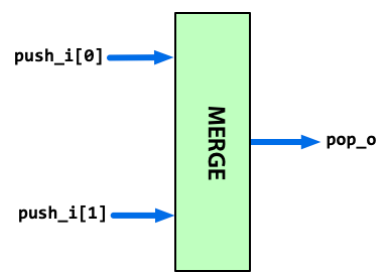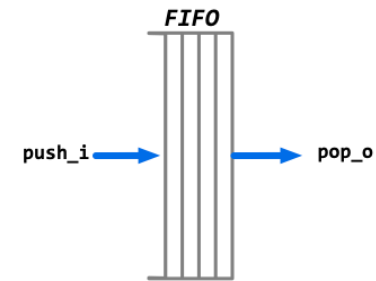


Demultiplexer
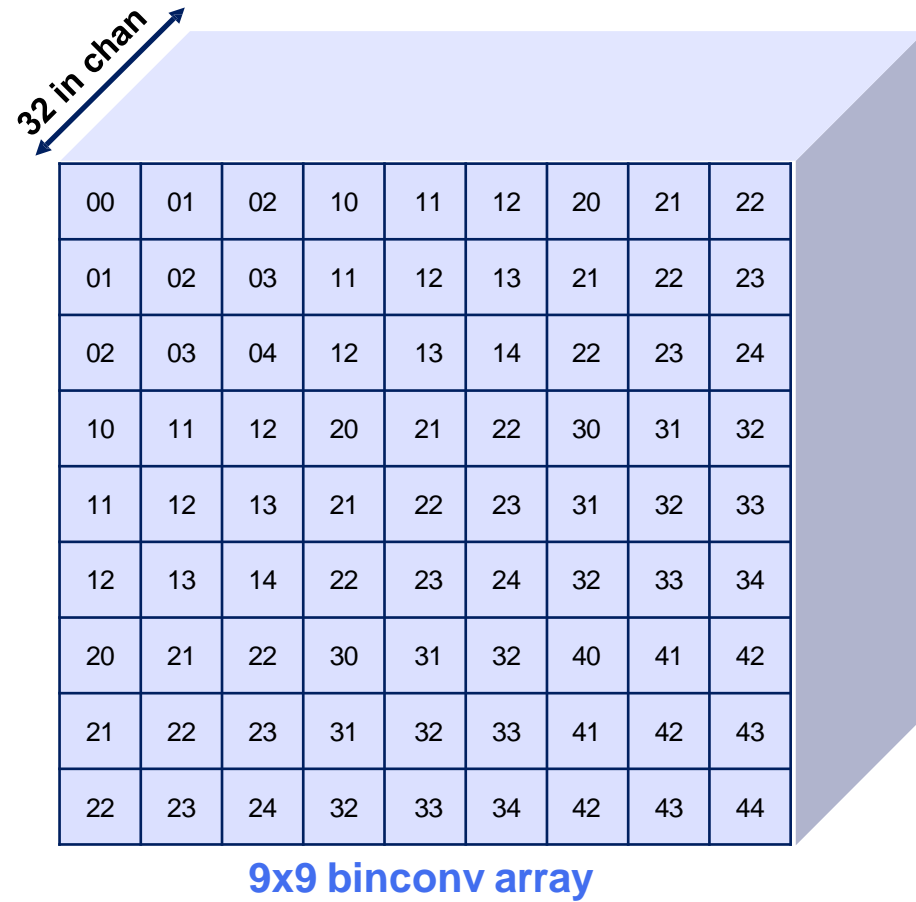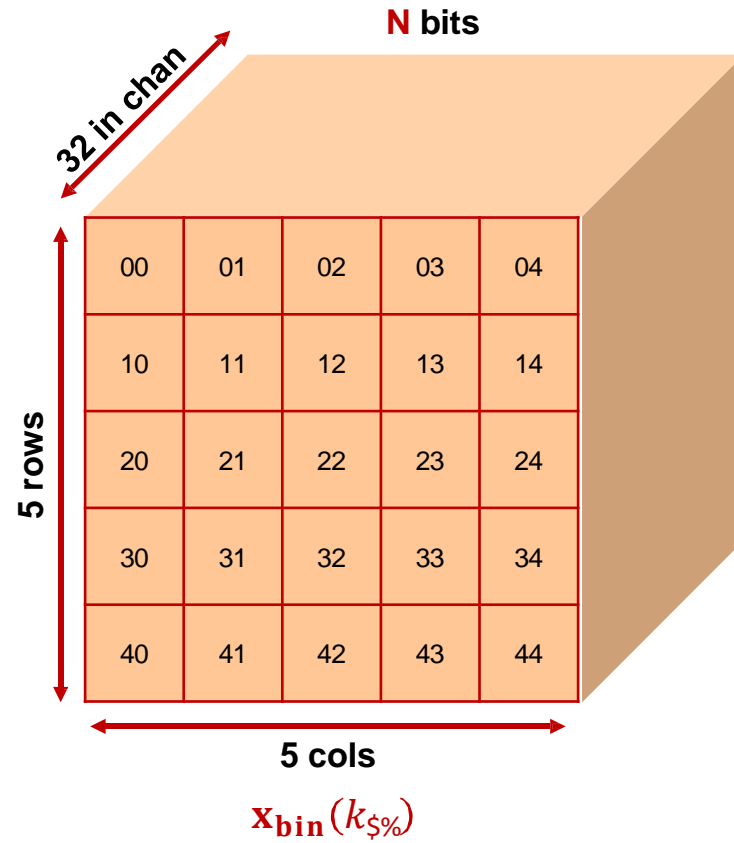


Multiplexer



Fence



Splitter



Merger



FIFO queues
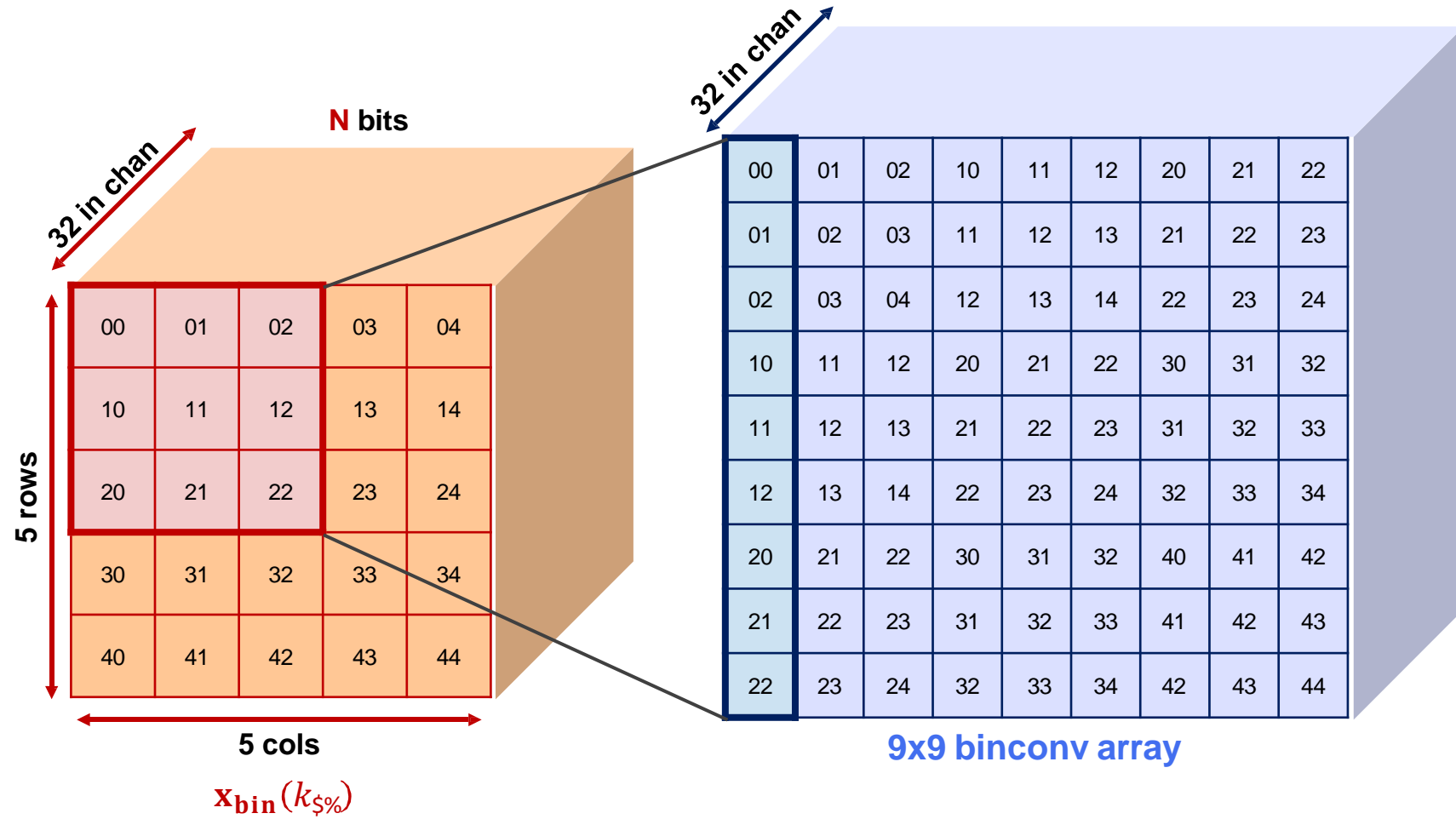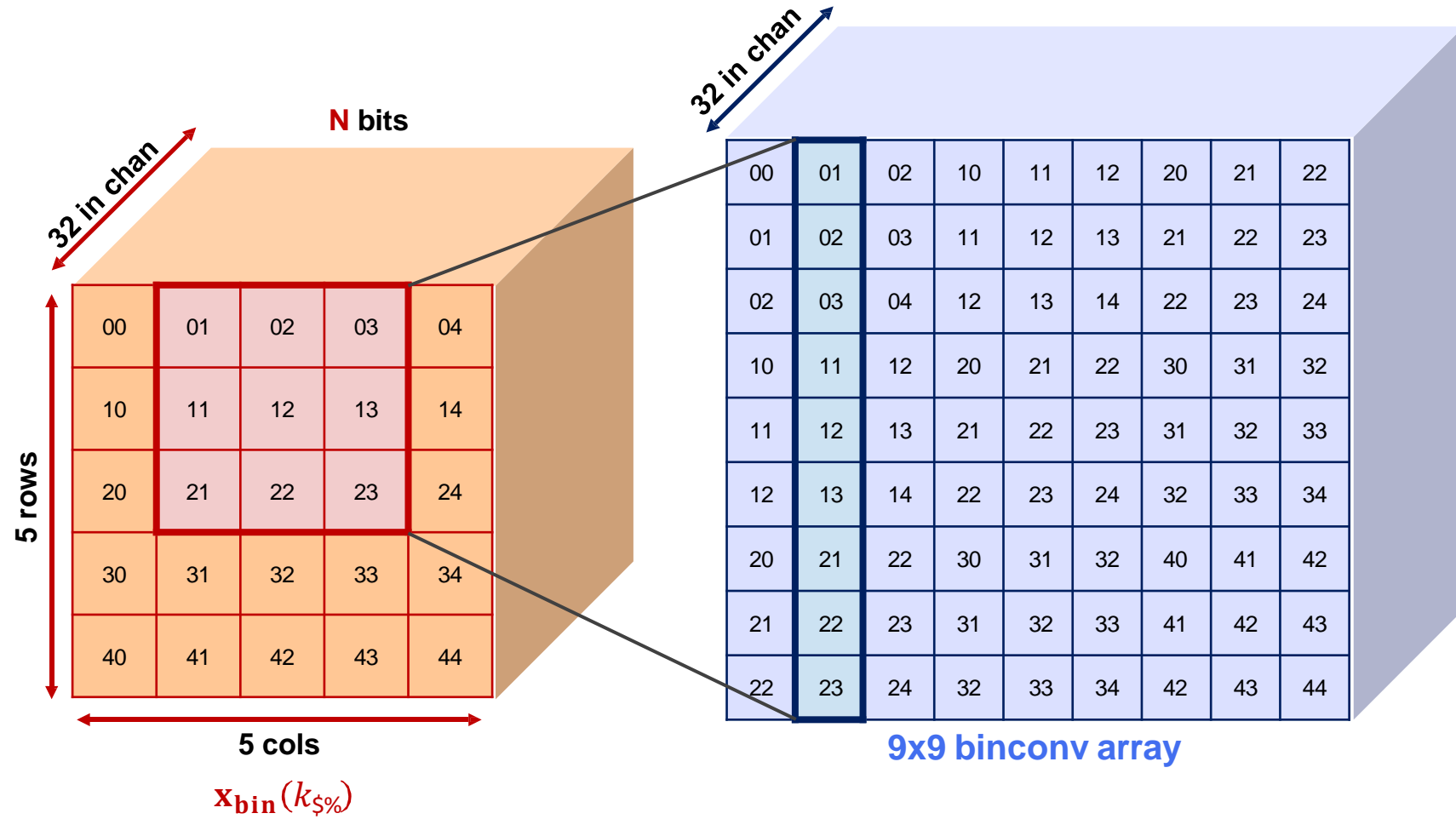
# RBE: 3x3 mapping



9x9 binconv array

$\mathbf{x_{bin}}(k_{\$\%})$

# RBE: 3x3 mapping



$\mathbf{x_{bin}}(k_{\$\%})$

9x9 binconv array

# RBE: 3x3 mapping



$$\mathbf{x_{bin}}(k_{\$\%})$$

9x9 binconv array

# RBE: 3x3 mapping



$W_{\mathbf{bin}}(k_{\$\%}, k_{\&'(})$

**9x9 binconv array**

# RBE: 3x3 mapping



$W_{\mathbf{bin}}(k_{\$\%}, k_{\&'(})$

**9x9 binconv array**

**32 out chan × M bits in 32×M cycles**

# RBE: 3x3 mapping

# RBE: 1x1 mapping
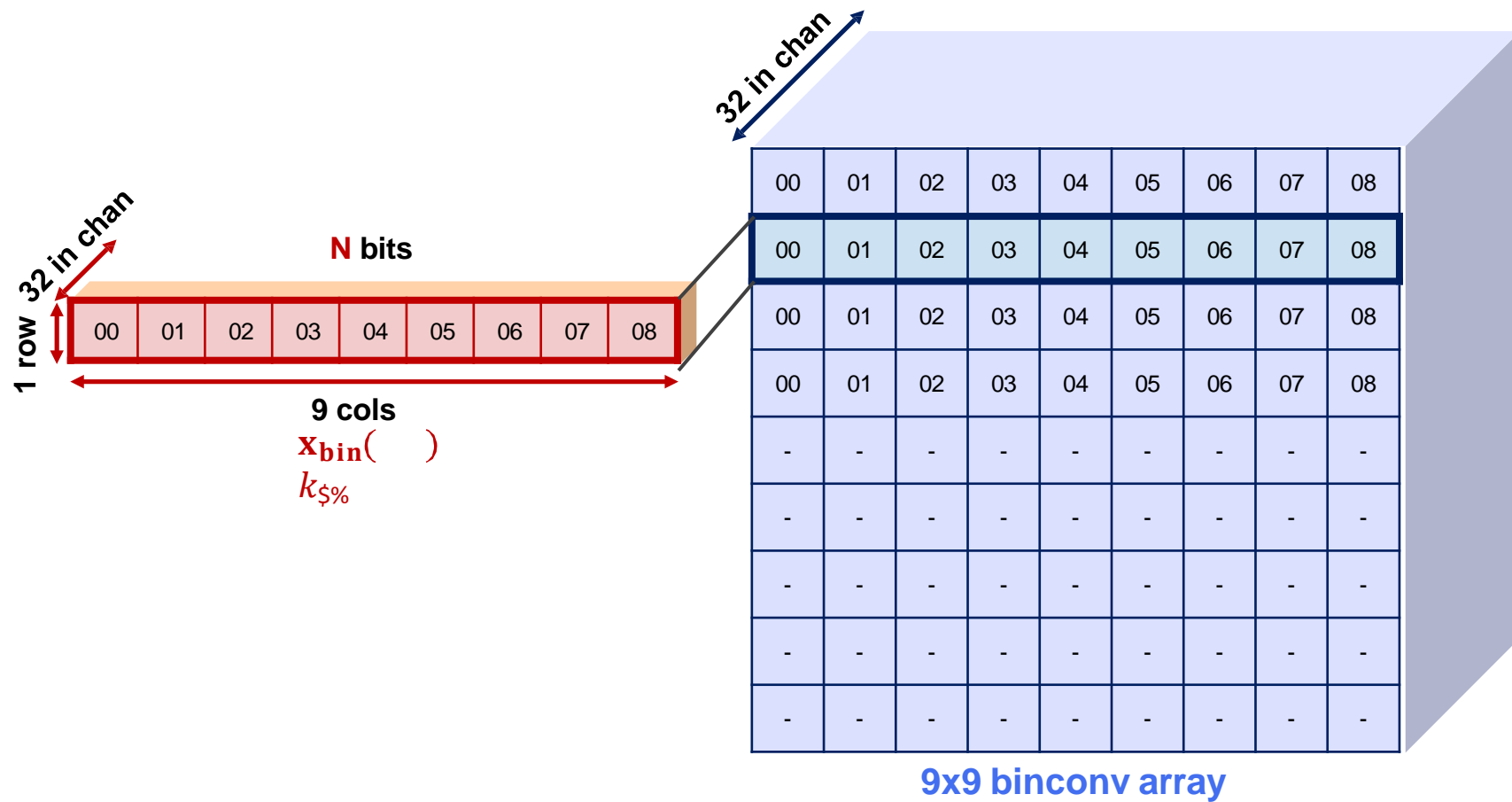


**32 in chan**

**N bits**

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |

**1 row**

**9 cols**

$x_{bin}(\quad)$

$k_{\$\%}$

**32 in chan**

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
|----|----|----|----|----|----|----|----|----|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| - | - | - | - | - | - | - | - | - |
| - | - | - | - | - | - | - | - | - |
| - | - | - | - | - | - | - | - | - |
| - | - | - | - | - | - | - | - | - |
| - | - | - | - | - | - | - | - | - |

**9x9 binconv array**

# RBE: 1x1 mapping



9x9 binconv array

# RBE: 1x1 mapping



**N bits**

**1 row** **32 in chan**

00 01 02 03 04 05 06 07 08

**9 cols**

$\mathbf{x_{bin}}(\quad)$
$k_{\$\%}$

**32 in chan**

| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
|----|----|----|----|----|----|----|----|----|
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 |
| -  | -  | -  | -  | -  | -  | -  | -  | -  |
| -  | -  | -  | -  | -  | -  | -  | -  | -  |
| -  | -  | -  | -  | -  | -  | -  | -  | -  |
| -  | -  | -  | -  | -  | -  | -  | -  | -  |
| -  | -  | -  | -  | -  | -  | -  | -  | -  |

**9x9 binconv array**

# RBE: 1x1 mapping



9x9 binconv array

# RBE: 1x1 mapping



9x9 binconv array

# RBE: 1x1 mapping



$W_{\mathbf{bin}}(k_{\$\%}, k_{\&'(})$

**9x9 binconv array**

**32 out chan in 32 cycles**

# RBE: 1x1 mapping