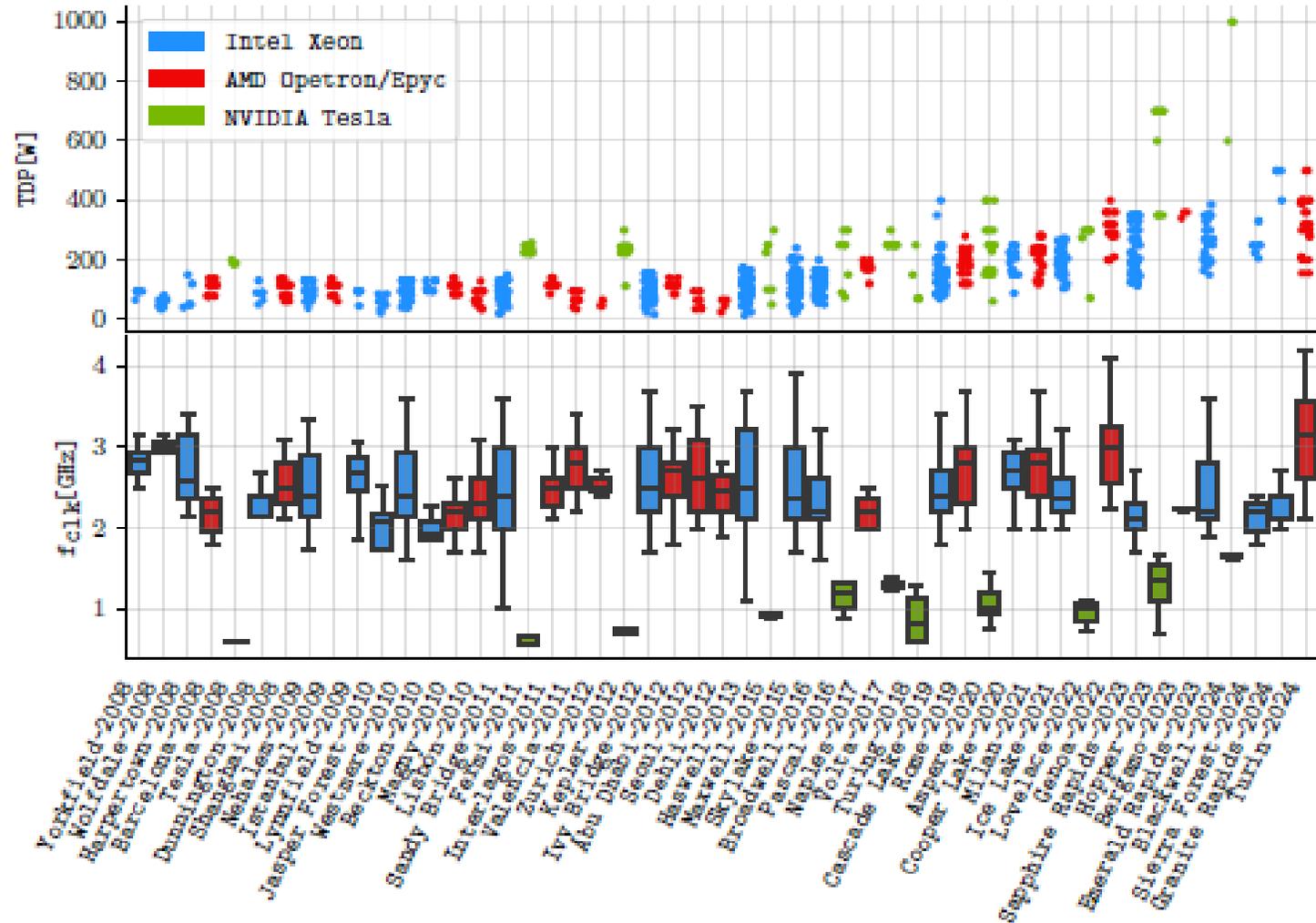# Energy efficiency and power management techniques for AI accelerators

Luca Benini
ETH Zurich, Univ. Bologna

ISSCC 2026 Forum 3.3

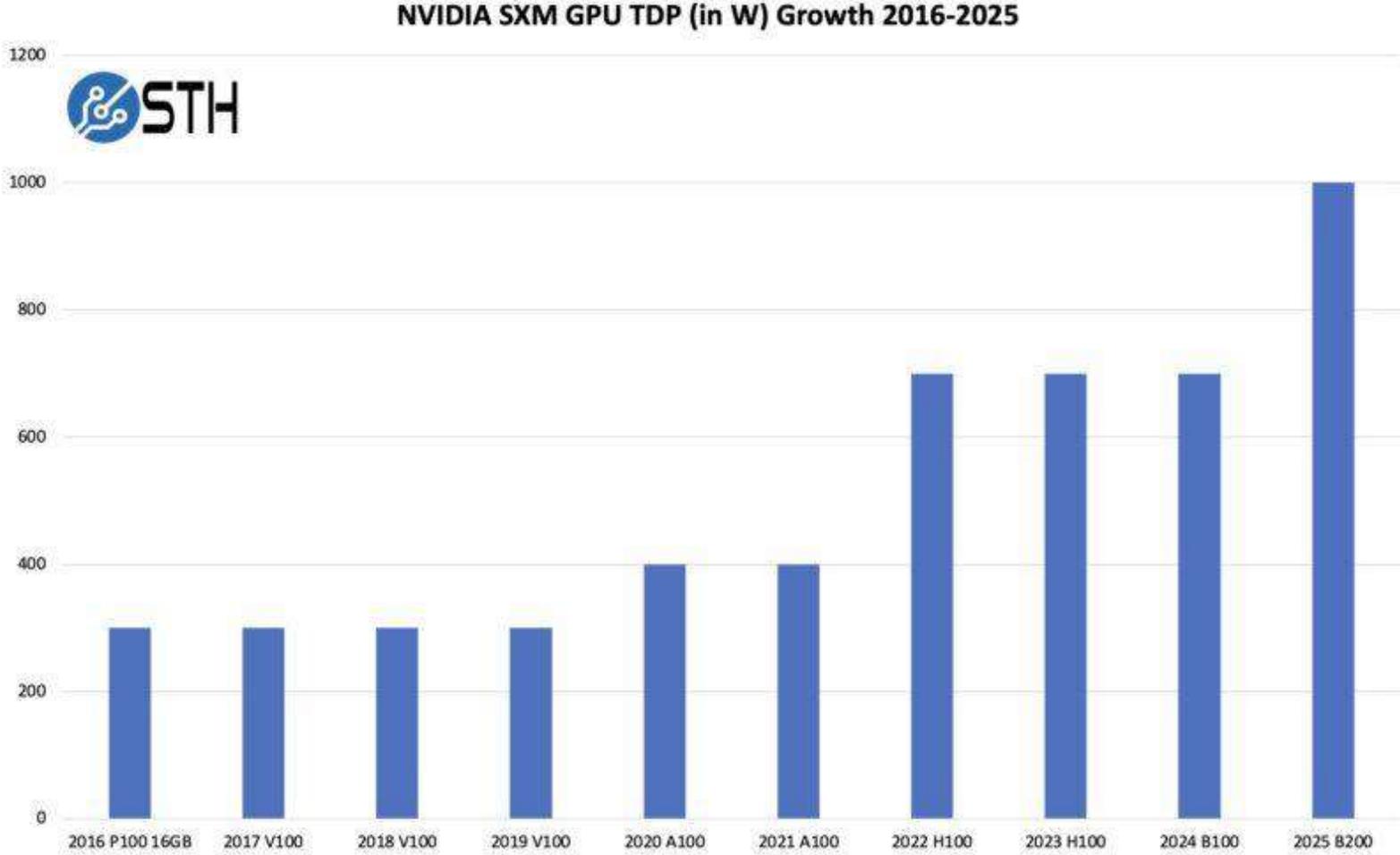# TDP of server Chips is growing

Clock speed and TDP of Intel, AMD, NVIDIA high-end CPUs/GPUs

# TDP of server AI Chips is growing Faster



NVIDIA SXM GPU TDP (in W) Growth 2016-2025

https://www.servethehome.com/why-servers-are-using-so-much-power-tdp-growth-over-time-supermicro-vertiv-intel-amd-nvidia

# Why? Ultra-high density, 2.5→3D, WS integration

# The MW AI Rack (in 2029)



https://wccftech.com/ai-servers-can-reach-rack-density-of-1000-kw-likely-with-nvidias-next-gen-rubin-ultra-architecture/?utm_source=chatgpt.com

# Why?

☐ GPU Density: AI systems rely on dense clusters of GPUs and TPUs with larger #MAC/Area  at high avg. utilization → more power, heat than CPUs

☐ Continuous Workloads: AI training and inference → sustained high-power for extended periods, unlike the more variable workloads of classical computing

☐ Chip Proximity: To achieve maximum performance and low latency, components are packed closer together → larger power and thermal demands within a single rack package → rack

# Outline

- ☐ <span style="color:red">Boosting efficiency for AI workloads</span>

- ☐ Managing idleness and heterogeneity in accelerated systems

- ☐ Using AI for managing AI

- ☐ Conclusions, future perspectives

# The Cambrian Explosion

**Domain-Specific Architectures**

# Domain-Specific Architecture Classes

|  | CPU Core | CPU Core Vector Engine | Parallel Thread Accelerator (GPUs) | Tensor Array Accelerator (TPU, Groq, TensorCore) | Microcore Mesh Accelerator (Cerebras, Tenstorrent) | Computational Block Accelerator (FPGA) | Custom Dataflow Accelerator (ASIC) |
|---|---|---|---|---|---|---|---|
| **Technology label** | CPU | Vector | GPU | Tensor | Manycore | FPGA | ASIC |
| **ALUs per core** | 1-4 | 8-32 | 32-64 | 8x8 to 256x256 | 1 to 4 | Various | App. specific |
| **Cores per processor** | 4-64 | 4-64 | 8-128 | 1 to 4 | 100s to 1M | Various | App. specific |
| **Parallel performance** | Low | Medium Low | Medium High | High | High | High | Very High |
| **Comp. efficiency (Ops/W)** | Low | Medium Low | Medium | High | High | High | Very High |
| **Comp. flexibility** | Very High | Medium Low | Medium | Medium Low | Medium | Low | Very Low |
| **Computation scheduling** | Dynamic by instruction | Dynamic by instruction | Dynamic by kernel | Static by kernel | Static by kernel | Fixed | Fixed |
| **Code redesign** | Seconds | Seconds | Seconds/Minutes | Minutes | Minutes/Hours | Hours | Months |

Increasing application specificity -> greater parallel performance -> narrower application/computational kernel enablement

# Heterogeneous Architecture

## Multiple Scales of acceleration
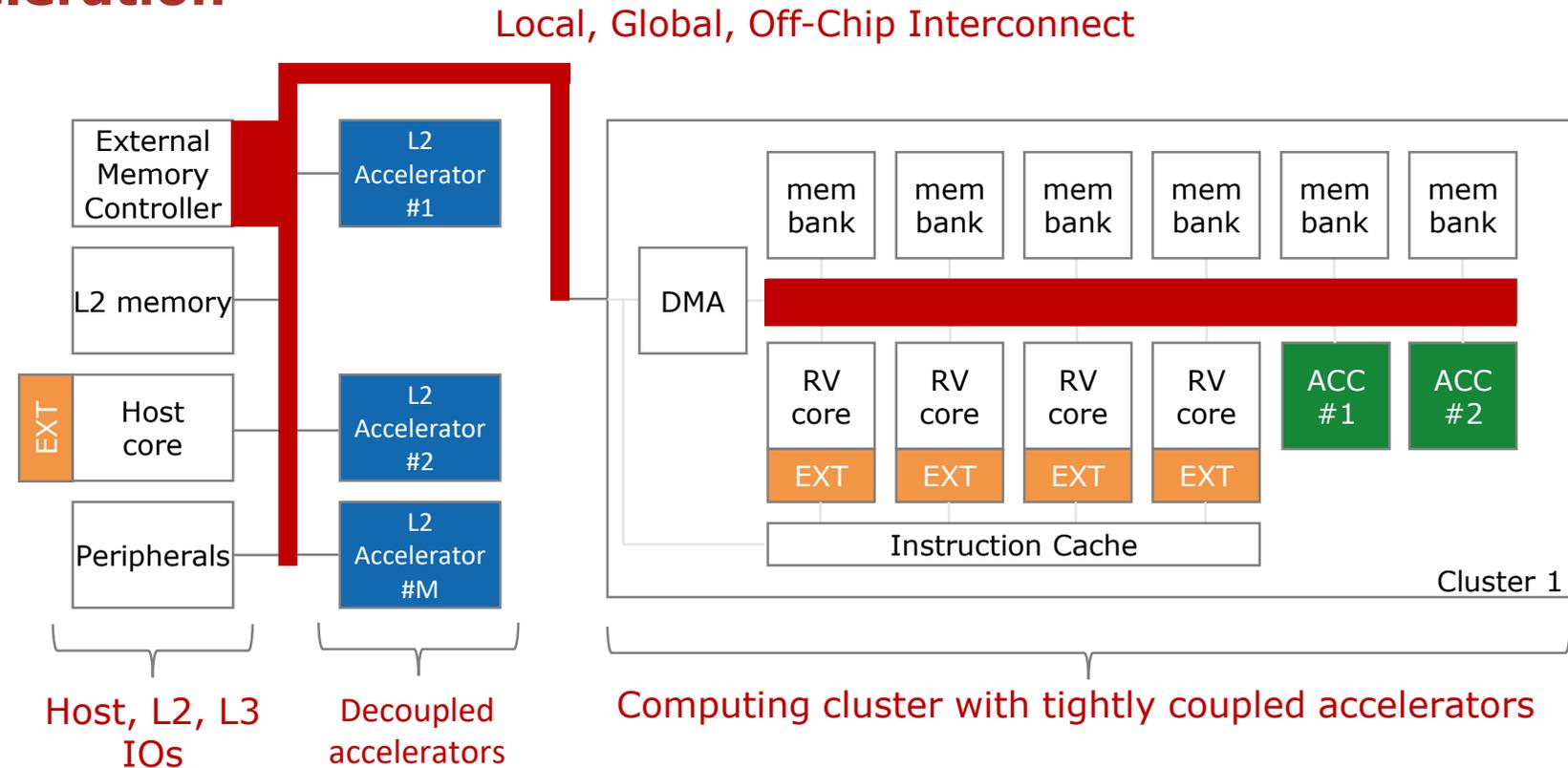
**Extensions to processor cores**
- Explore new extensions
- Efficient implementations

**Shared-memory Accelerators**
- Domain specific
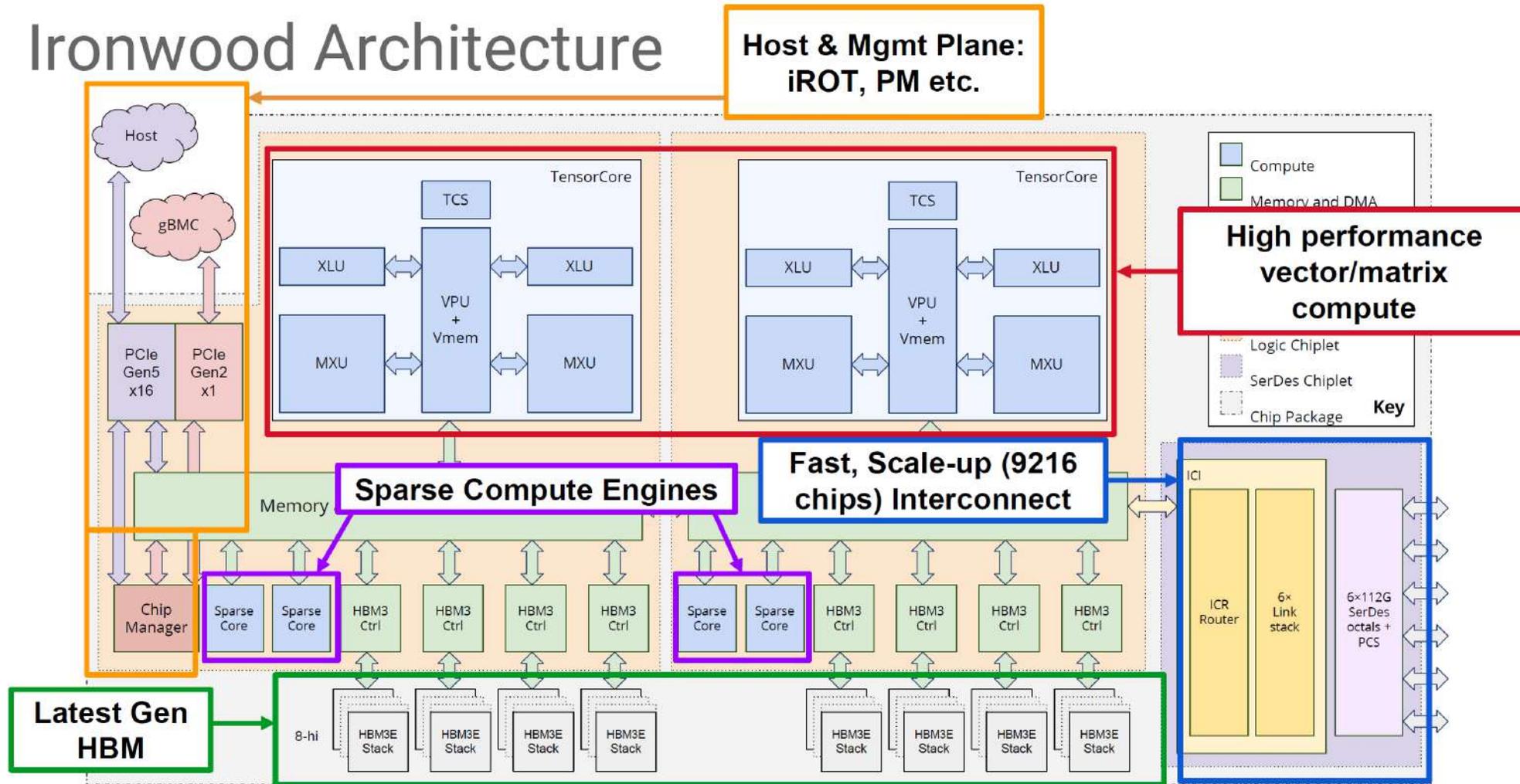- Local memory

**Multiple Decoupled Accelerators**
- Communication
- Synchronization

Local, Global, Off-Chip Interconnect

| External Memory Controller | L2 Accelerator #1 | | mem bank | mem bank | mem bank | mem bank | mem bank | mem bank |

DMA

| RV core | RV core | RV core | RV core | ACC #1 | ACC #2 |
| EXT | EXT | EXT | EXT | | |

Instruction Cache

Cluster 1

EXT | Host core | L2 Accelerator #2

Peripherals | L2 Accelerator #M

L2 memory

Host, L2, L3 IOs

Decoupled accelerators

Computing cluster with tightly coupled accelerators
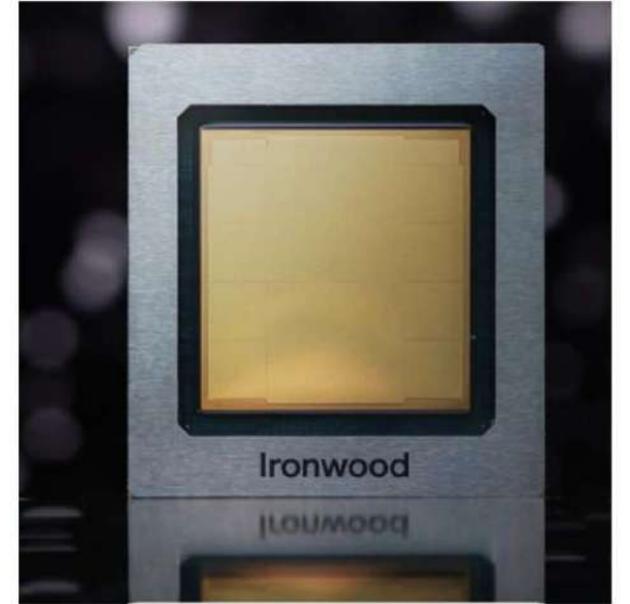
**Local, global, package, system**

# A notable Example: Google's Ironwood

# Ironwood Chip

- First dual compute die TPU, 4614 TFLOPS of FP8
  - >10x compute compared to TPU v5p
  - Capable of large scale pretraining of foundation models
- 8 stacks of HBM3E, peak BW 7.3 TB/s, capacity 192 GiB
  - Optimized for serving latest generation thinking models
- 1.2 TBps I/O to gluelessly scale-up to 9216 chips
- Industry leading cold plate thermal solution
- Integrated root-of-trust (iROT) for secure computing
- Functional BIST & silent data corruption (SDC) mitigation
- Logic repair to improve yield
- Dynamic voltage/frequency scaling for efficient perf/W
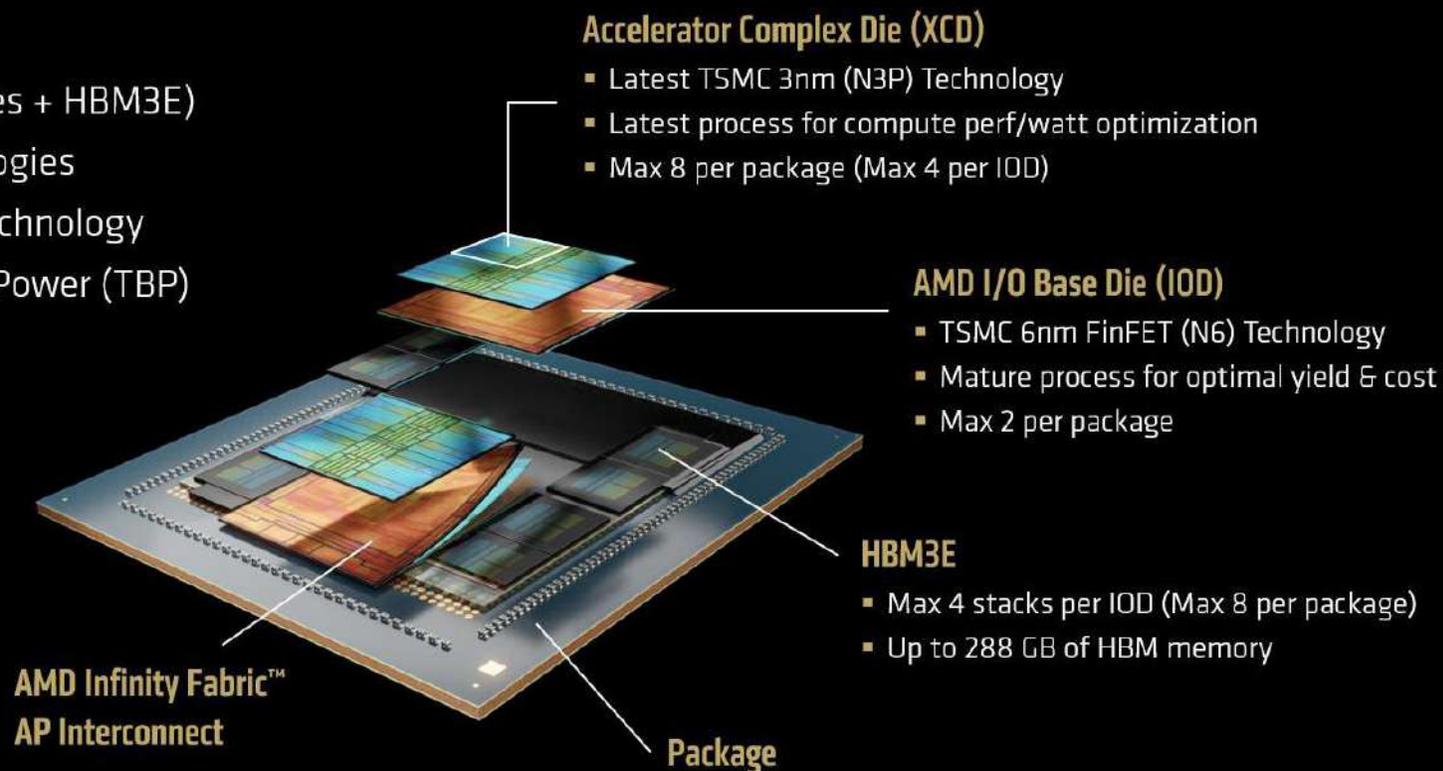- AI based ALU circuits, floor plan optimization

# Looking forward: more (3D) Heterogeneity



**AMD Instinct™ MI350 Series GPU**
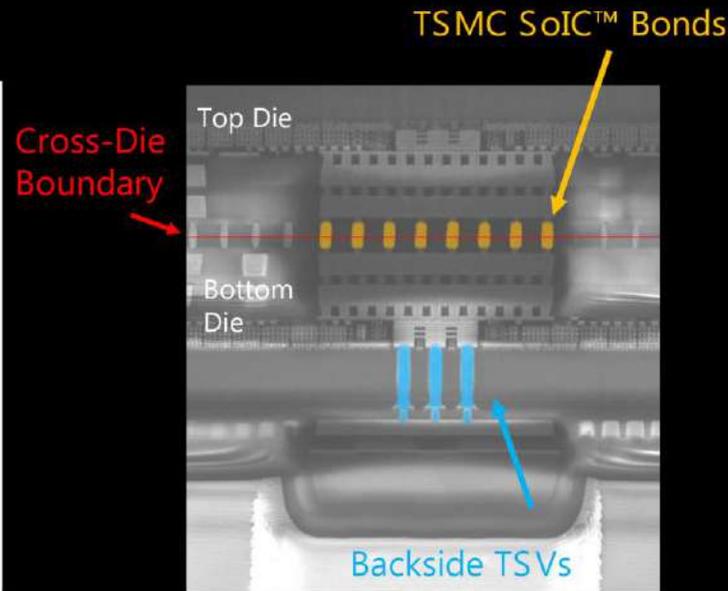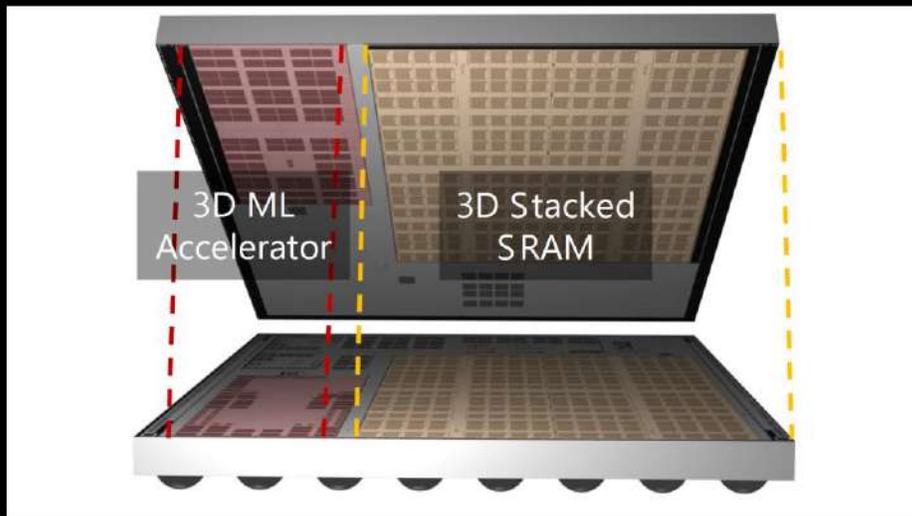
State-of-the-Art Construction

- 185 Billion transistors
- 3D Multi-Chiplet (2 chiplet types + HBM3E)
- Heterogenous process technologies
- Proven COWOS-S packaging technology
- Up to 1,400 watts Total Board Power (TBP)

**Accelerator Complex Die (XCD)**
- Latest TSMC 3nm (N3P) Technology
- Latest process for compute perf/watt optimization
- Max 8 per package (Max 4 per IOD)

**AMD I/O Base Die (IOD)**
- TSMC 6nm FinFET (N6) Technology
- Mature process for optimal yield & cost
- Max 2 per package

**HBM3E**
- Max 4 stacks per IOD (Max 8 per package)
- Up to 288 GB of HBM memory

**AMD Infinity Fabric™ AP Interconnect**

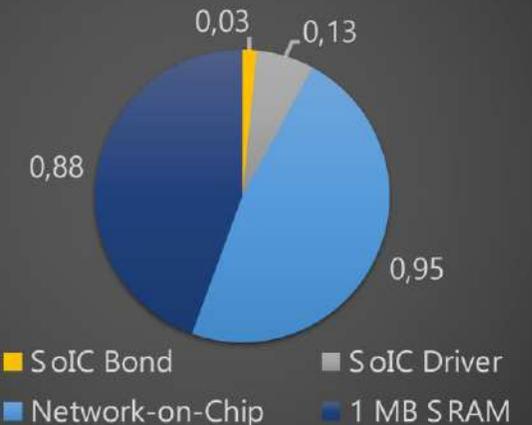**Package**

AMD – Hot Chips 2025

# Looking forward: more (3D) Heterogeneity



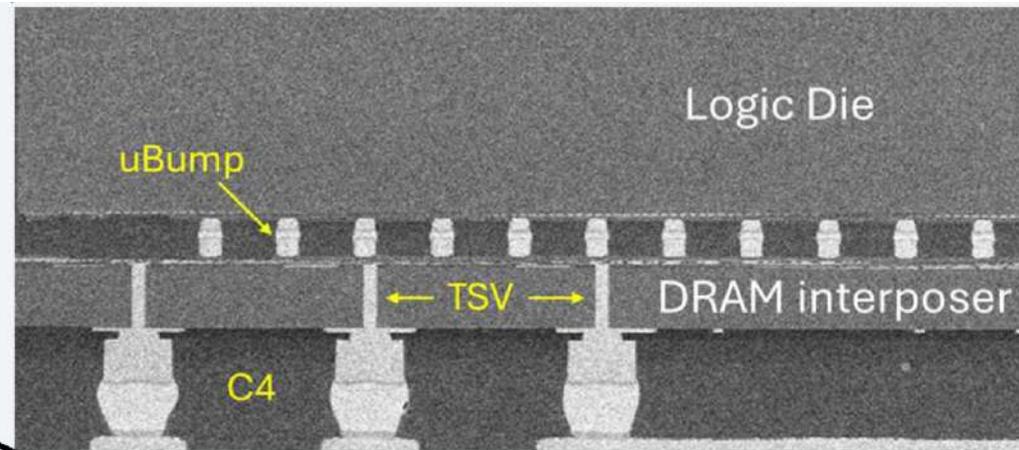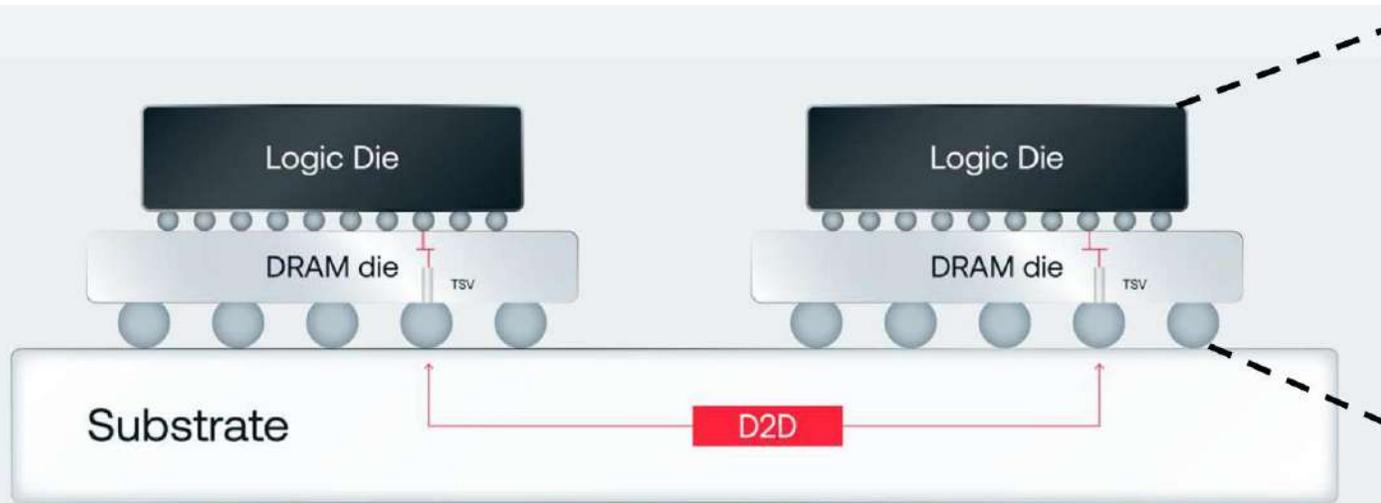## 3D Integration Enables a Path to AR Silicon

- 3D stacking for addressing the combined Power-Performance-Area challenges of AR SoCs within the same 2D area footprint
- Implemented a prototype 3D stacked AR SoC: Two 7nm dies, wafer-on-wafer face-to-face stacked at < 2um bonding pitch using TSMC SoIC™ bonding technology
- Integrates 3D ML Accelerator, 3D Stacked SRAM, CPU for realistic model deployment

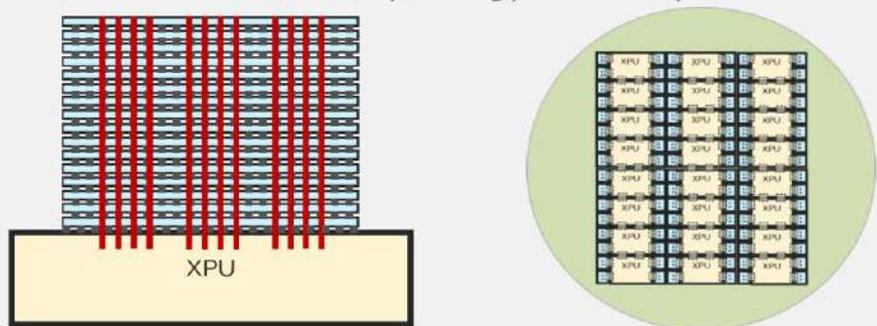META – ISSCC24 (Beigne et al.)

# Looking forward: more (3D) Heterogeneity



- Top die: TSMC N5 logic
- Bottom die: 3D DRAM
- Integration: 36µ Face to Face (F2F) stacking
- Proven low cost, high volume, high yield process
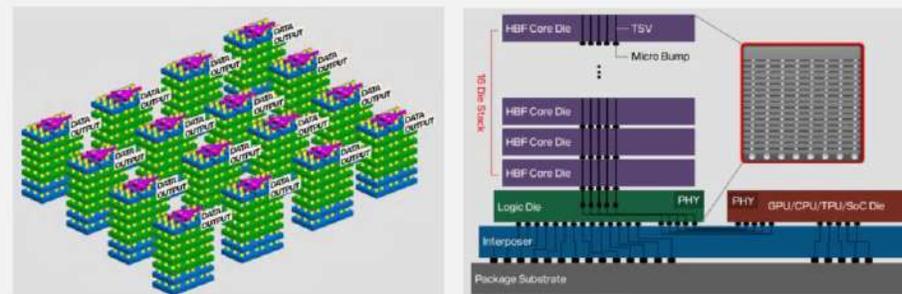
# Looking forward: more (3D) Heterogeneity



**Memory Integration Innovations**
(Ex: 3D DRAM + XPU, Wafer Scale Integration)
>4-5X bandwidth, density, energy-efficiency vs. 2.5D HBM

XPU

**Memory Technology Innovations**
(Example: SanDisk High-Bandwidth NAND Flash)
8-16X capacity at bandwidth/cost equivalent to HBM

Source: SanDisk Investor Day, 2025

**System Integration Innovations**
(Ex: Spatial 3D / Z-Axis Integration)

Planar (Pancake) Stack

L1 Cache | Cores
L2/L3 Cache Stack

Z-Axis (Loaf Bread) Stack

L2/L3 Cache Stack
L1 Cache | Cores

Scale-Up

>10X Dense System Integration

High-Bandwidth High-Capacity Memory Stack

Legacy I/O

Photonics

Source: J. Fryman, DARPA ERI Summit, 2023

INTEL CHISIC 2025

# Efficiency Quest – Summary

☐ Huge variety of domain-specific architectures → PM flexibility is crucial

☐ Heterogeneity is key → fine grained and heterogeneous power managers with lots of sensor data streams and multiple concurrent, heterogeneous PM tasks

☐ 3D integration → closed-loop massive MIMO power and thermal management is essential – active cooling is likely going fine-grained

# Outline

☐ Boosting efficiency for AI workloads

☐ Managing idleness and heterogeneity in accelerated systems
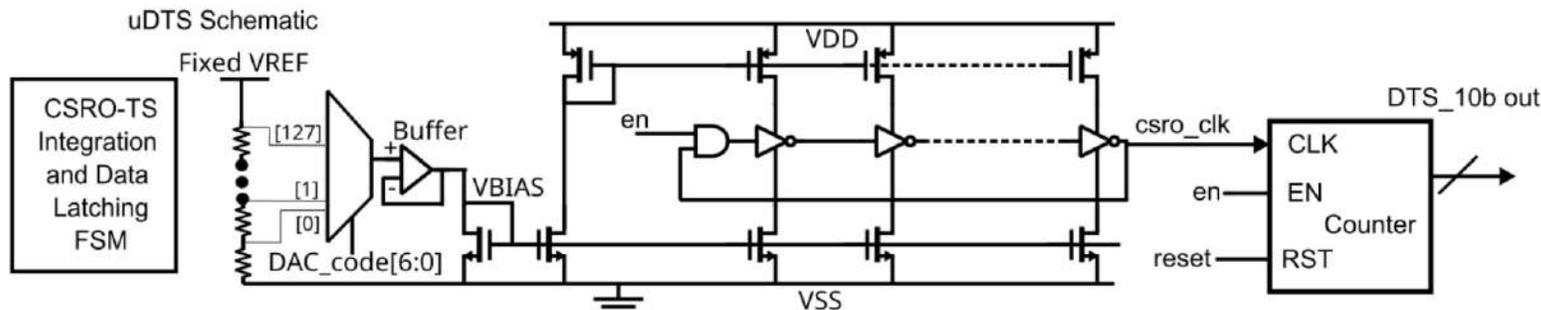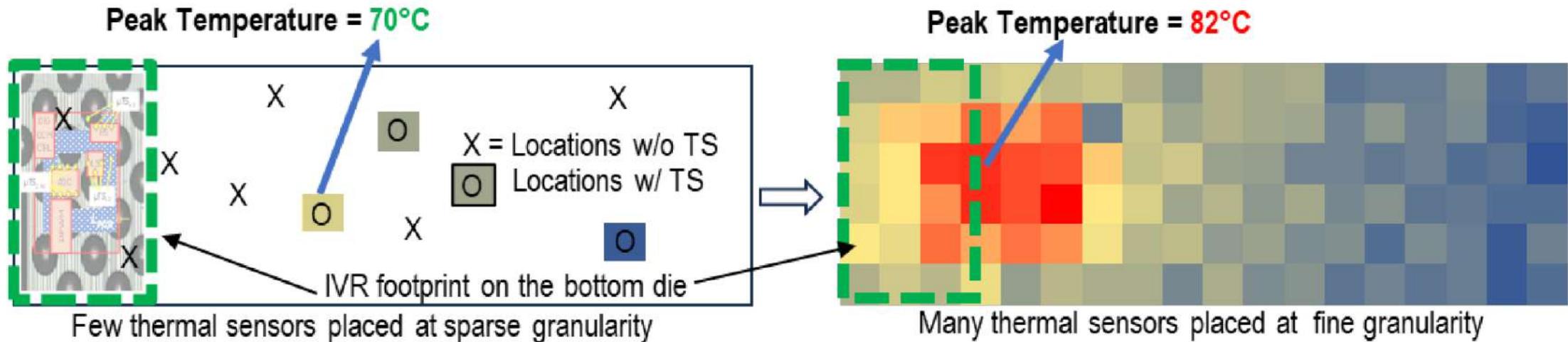
☐ Using AI for managing AI

☐ Conclusions, future perspectives

# Ultra-fine grained (digital) sensing (Temp)

☐ Micro-scale Sensors placed in a fine-grained array improve resolution and enable better thermal decisions



Peak Temperature = 70°C

X = Locations w/o TS
O Locations w/ TS

IVR footprint on the bottom die

Few thermal sensors placed at sparse granularity

Peak Temperature = 82°C

Many thermal sensors placed at fine granularity

uDTS Schematic

CSRO-TS Integration and Data Latching FSM

Fixed VREF

[127]
[1]
[0]
DAC_code[6:0]

Buffer
VBIAS

VDD

en

csro_clk

VSS

DTS_10b out

CLK
EN
Counter
RST

en
reset

- Linearity and sensitivity of the CSRO-TS are programmed through VBIAS
- Programmable integration time enables noise robustness

INTEL ISSCC25 8.8

# Ultra-fine grained (mixed-signal) sensing (Vdd droop)



- DDC: Variation-tolerant reference generation ⟶ High Accuracy
- RDD: Self-calibrated comparator
  - High-speed comparator & Scaler ⟶ High-Speed
- Multiple RDDs share a single DDC. ⟶ Area Efficiency

Samsung ISSCC25 8.4

# Fine-grained Regulation (Vdd)



- **Minimize <u>total</u> energy dissipation**
  - Voltage Regulator (VR)
  - Load (SoC):

- **Fine-grained DVFS**
  - Improved $\eta_{soc}$
  - Needs multiple VRs

- **Power Distribution Network ➜ IVR**
  - Scalability vs $\eta_{VR}$ tradeoffs

Gatech ISSCC25 8.3

# Fine-grained Regulation (Vdd)



- **Partial Crossbar**
  - Near-optimal domain assignment

- **Runtime Reconfigurable connectivity**

- Sample State Variables from all LDOs
- Establish the **lowest $V_{buck}$** to satisfy all $V_{dd}$ targets

Gatech ISSCC25 8.3

# Fine-grained Regulation (Vdd)



- Domain runtime re-configurability ($M_0$, $M_1$)

- Controller for VR and hot-swapping domain transition

- LDO state variable broadcasting to each buck tile

- Phase selection for $1\phi$ or $2\phi$ Buck

Gatech ISSCC25 8.3

# MIMO Control: IBM Telum2 (Z-systems)



**IBM Telum**

**IBM Telum II**

Voltage

Timing & Performance Protection

Voltage Guard Band

Other Guard Band

Performance Protection

Timing Protection

**Timing & Performance Protection**

- lumped together
- static $V_{SET}$ and droop mitigation

**Timing Protection**

- Droop mitigation only

**Performance Protection**

- Dynamic
- Workload dependent

IBM ISSCC25 8.1

# MIMO Control: IBM Telum2 (Z-systems)



- ☐ Multiple interacting control loops
  - ■ Performance Loop
  - ■ Voltage droop loop
  - ■ Recoverable error loop (exploiting RAS feature – incl. robust core recovery)

IBM ISSCC25 8.1

# PM: System-Level View



- Out-of-band – zero overhead telemetry
- Node Pcap – Max perf @ Pnode<Pmax
- RAS – error and conditions reporting
  - cpufreq/ cpuidle
  - Based on O.S. metrics
  - Slow & often unused

System Management / RM

Power Cap

Energy vs. Throuput

System Management / RM

Application

Hints/Prescription

Operating System

In band

Governors

DIMM

PowerOn/Reset

RAS

Node Power Cap

Out of band

VRM

RJ45 | BMC

Power Controller

PE

S

# PM: System-Level View

**Power Management standard HW/SW interfaces:**
**In-band:**
- The ACPI (Advanced Configuration and Power Interface) for power states.
- The **SCMI (System Control and Management Interface)** for OS communication.

System Management / RM

System Management / RM

Power Cap

Energy vs. Througput

Application

Hints/Prescription

Operating System

In band

Governors

DIMM

PowerOn/Reset

RAS

Node Power Cap

Out of band

VRM

S

RJ45

BMC

Power Controller

PE

# PM: System-Level View

**Power Management standard HW/SW interfaces:**
**In-band:**
- The ACPI for power states
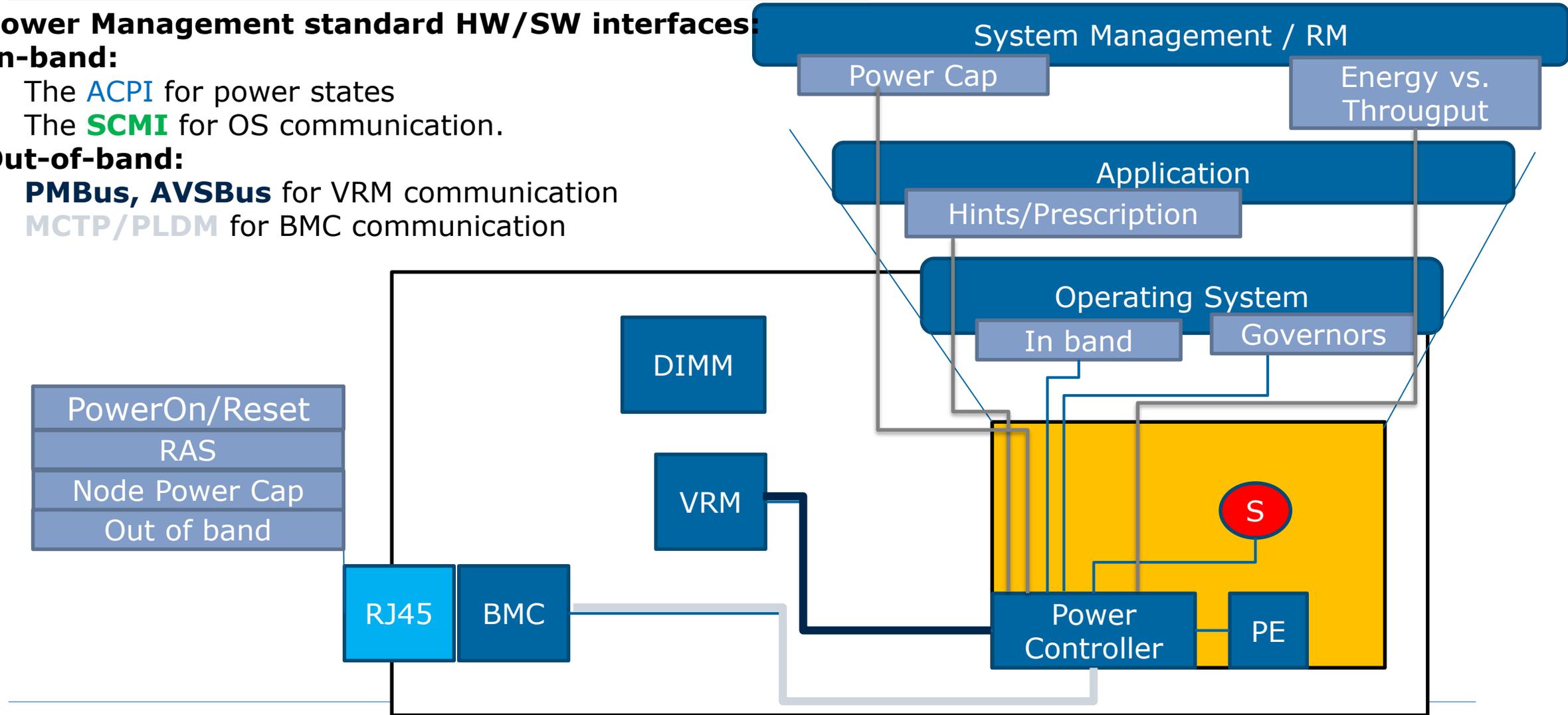- The **SCMI** for OS communication.

**Out-of-band:**
- **PMBus, AVSBus** for VRM communication
- **MCTP/PLDM** for BMC communication

System Management / RM

System Management / RM

Power Cap

Energy vs. Throuput

Application

Hints/Prescription

Operating System

In band

Governors

PowerOn/Reset

RAS

Node Power Cap

Out of band

DIMM

VRM

RJ45

BMC

Power Controller

PE

S

# PM: System-Level View



- ☐ **High level controller (HLC) governors**
  - ■ Operating System (OS): e.g., OSPM in the Linux kernel
  - ■ Baseboard Management Controller (BMC): e.g., Intel Aspeed off-chip regulator
- ☐ **Low level controller (LLC) embedded platform**
  - ■ Power control subsystem (PCS): typically, a single-core MCU with < 1MiB memory footprint
  - ■ Power control firmware (PCF): FW routine on top of the PCS

# Focus on LLC: Function



**Control algorithm in a periodic loop**

Plant, Ctrl stats, Events

Constraints & Targets
e.g. $T_{MAX}$, $Perf_{MIN}$

Control Task

HLC (SW)

Sensors

Actuatos

**RT Constraint: control task ends within period, with small jitter**

**Controlled System (massive MIMO)**

# ControlPULP: Open-source Controller for PM



**PFCT:** Periodic Frequency Control Task
**PVCT:** Periodic Voltage Control Task

**HLC:** high level controller
**LLC:** low level controller

# ControlPULP: Open-source controller for PM



Governor OS - **on-chip HLC**

PCF + FreeRTOS

HPC processing elements

PLLs

Secure boot subsystem

P1

P2

P3

AXI Slave          Boot

**ControlPULP PCS - on-chip LLC**

**PFCT**

Power comp

SCMI doorbells

Alpha comp.

PID comp.

Control Action

PVCT

F/V comp.

I/O

**Two challenges**

1. Number of controlled computing cores is increasing
2. Advanced DTPM is expensive (**computation size** and **feedback time-scale**)

CPU silicon die and socket

Motherboard

Off-die out-of-band

Voltage IN/OUT

DVS (MCTP)

PMBUS AVSBUS

I2C

VRMs

BMC off-chip HLC

**PFCT:** Periodic Frequency Control Task
**PVCT:** Periodic Voltage Control Task

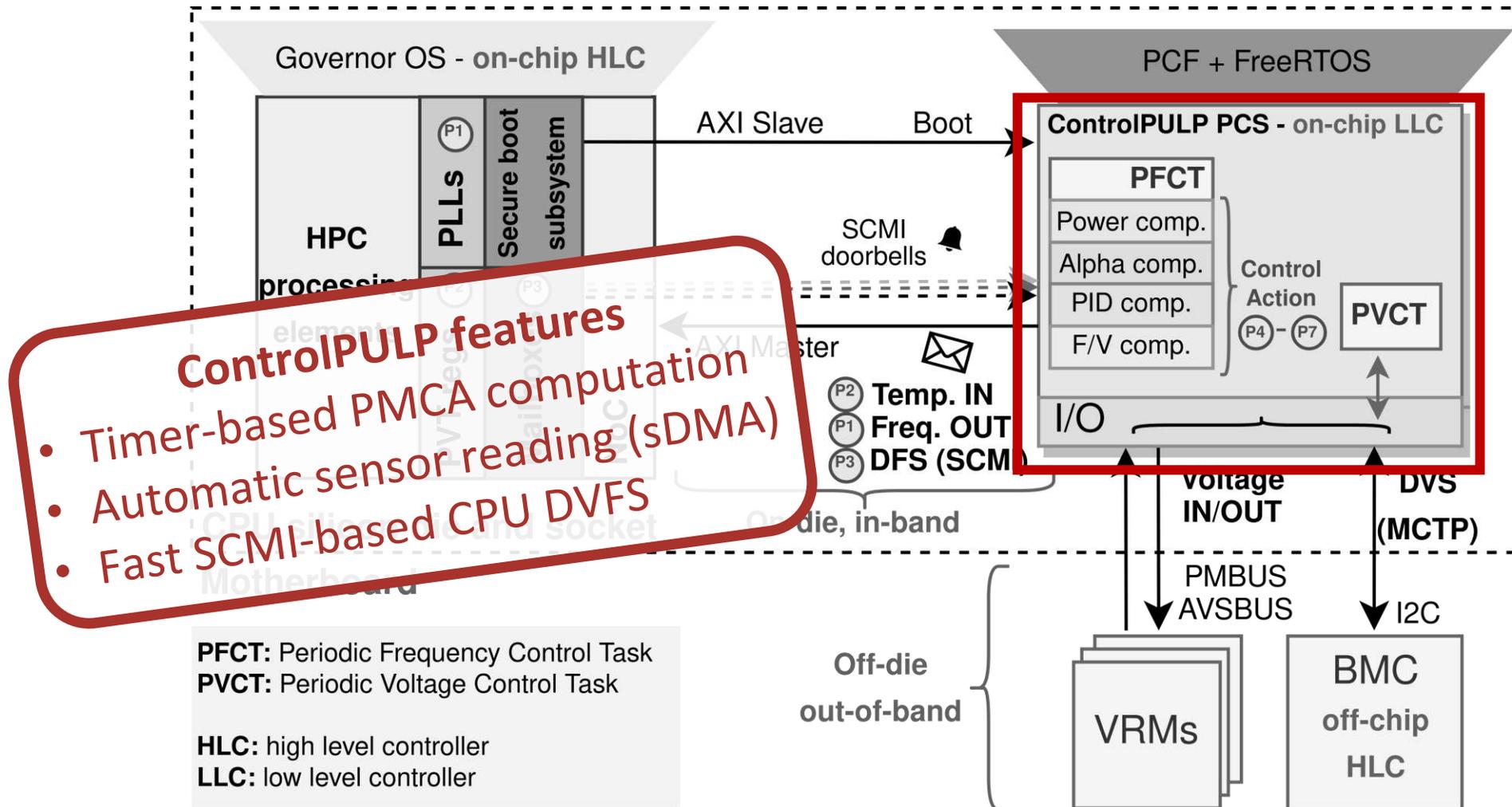**HLC:** high level controller
**LLC:** low level controller

# ControlPULP: Open-source controller for PM



Two requirements for an embedded PCS
1. High energy efficiency
2. Soft/firm real-time capabilities

PFCT: Periodic Frequency Control Task
PVCT: Periodic Voltage Control Task

HLC: high level controller
LLC: low level controller

# ControlPULP: Open-source controller for PM



**ControlPULP features**
- Timer-based PMCA computation
- Automatic sensor reading (sDMA)
- Fast SCMI-based CPU DVFS

**PFCT:** Periodic Frequency Control Task
**PVCT:** Periodic Voltage Control Task

**HLC:** high level controller
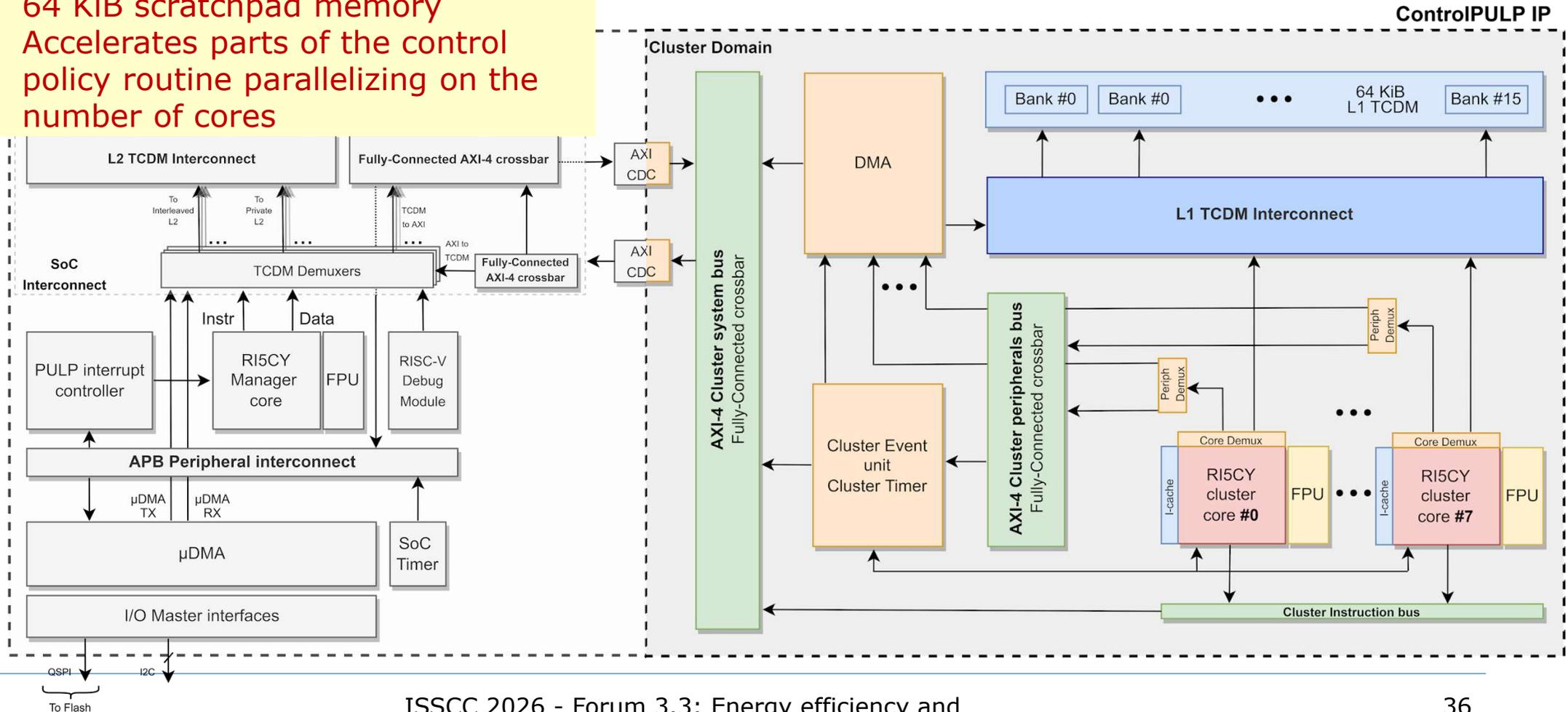**LLC:** low level controller

# ControlPULP Architecture



**Single-core subsystem**
- 32-bit CV32E40P
- 512 KiB scratchpad memory
- Executes the main control policy routine
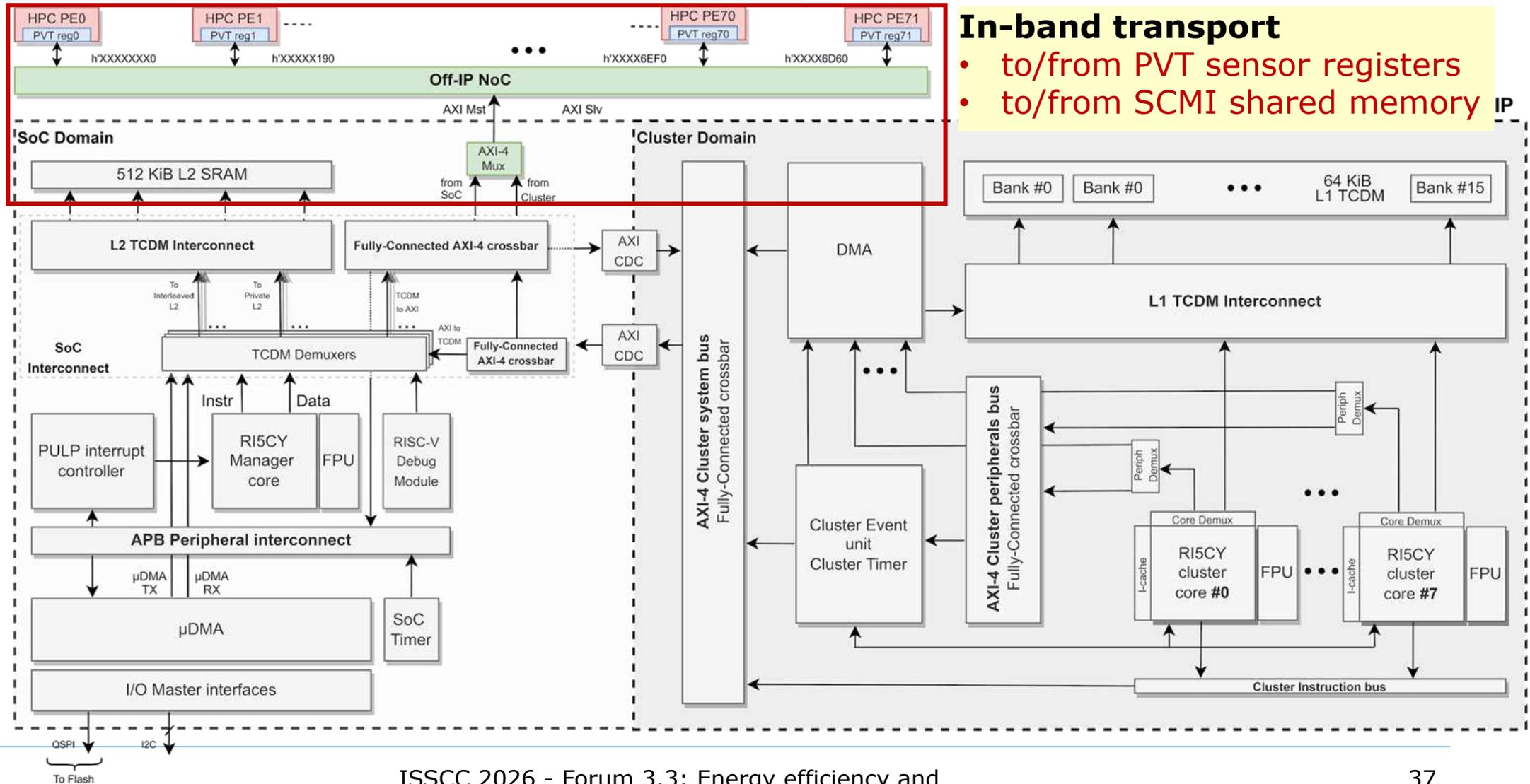- Offloads tasks to accelerator cluster

# ControlPULP Architecture

**Multi-core cluster subsystem**
- 8 32-bit CV32E40P
- 64 KiB scratchpad memory
- Accelerates parts of the control policy routine parallelizing on the number of cores
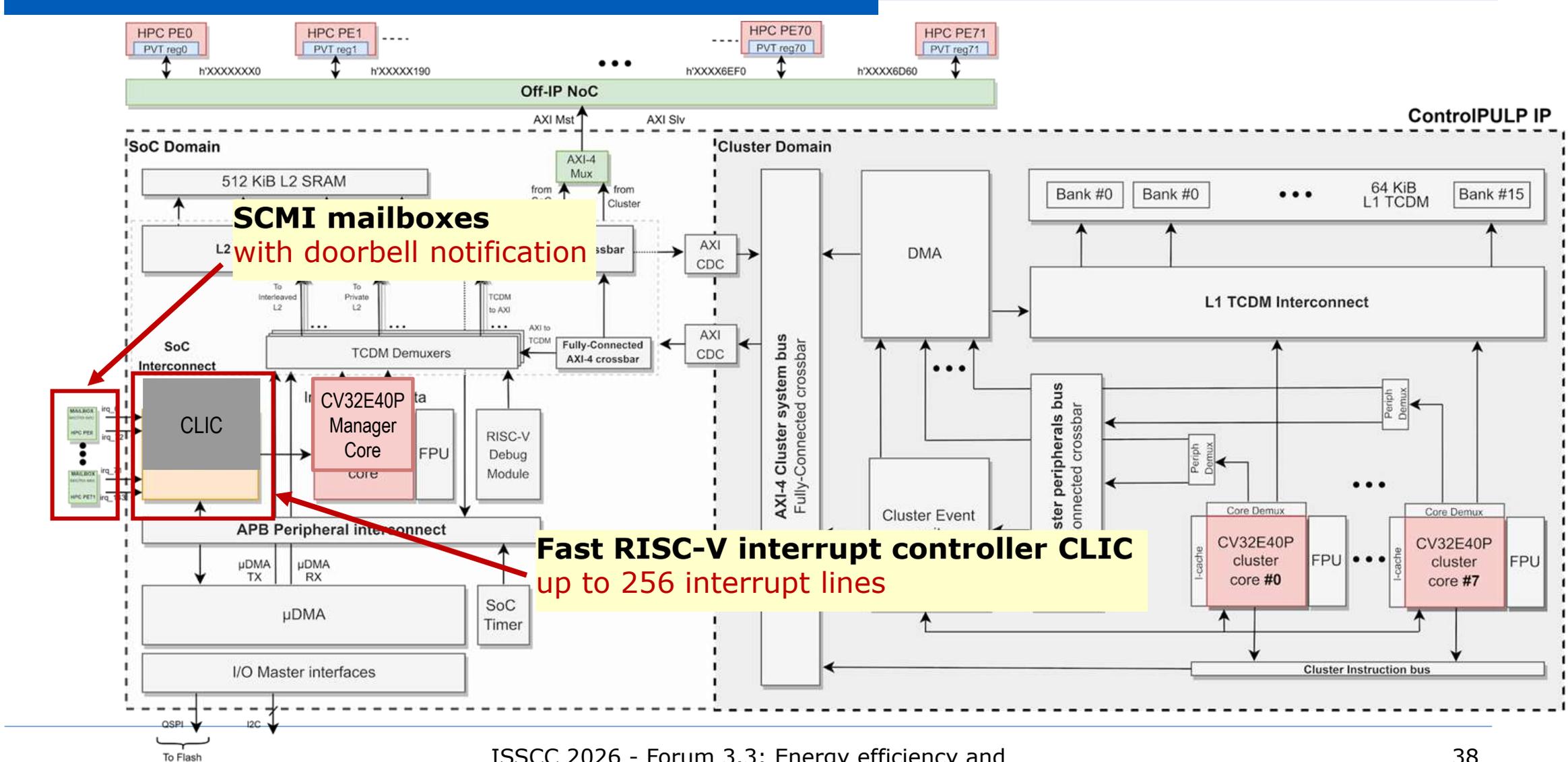
# ControlPULP Architecture



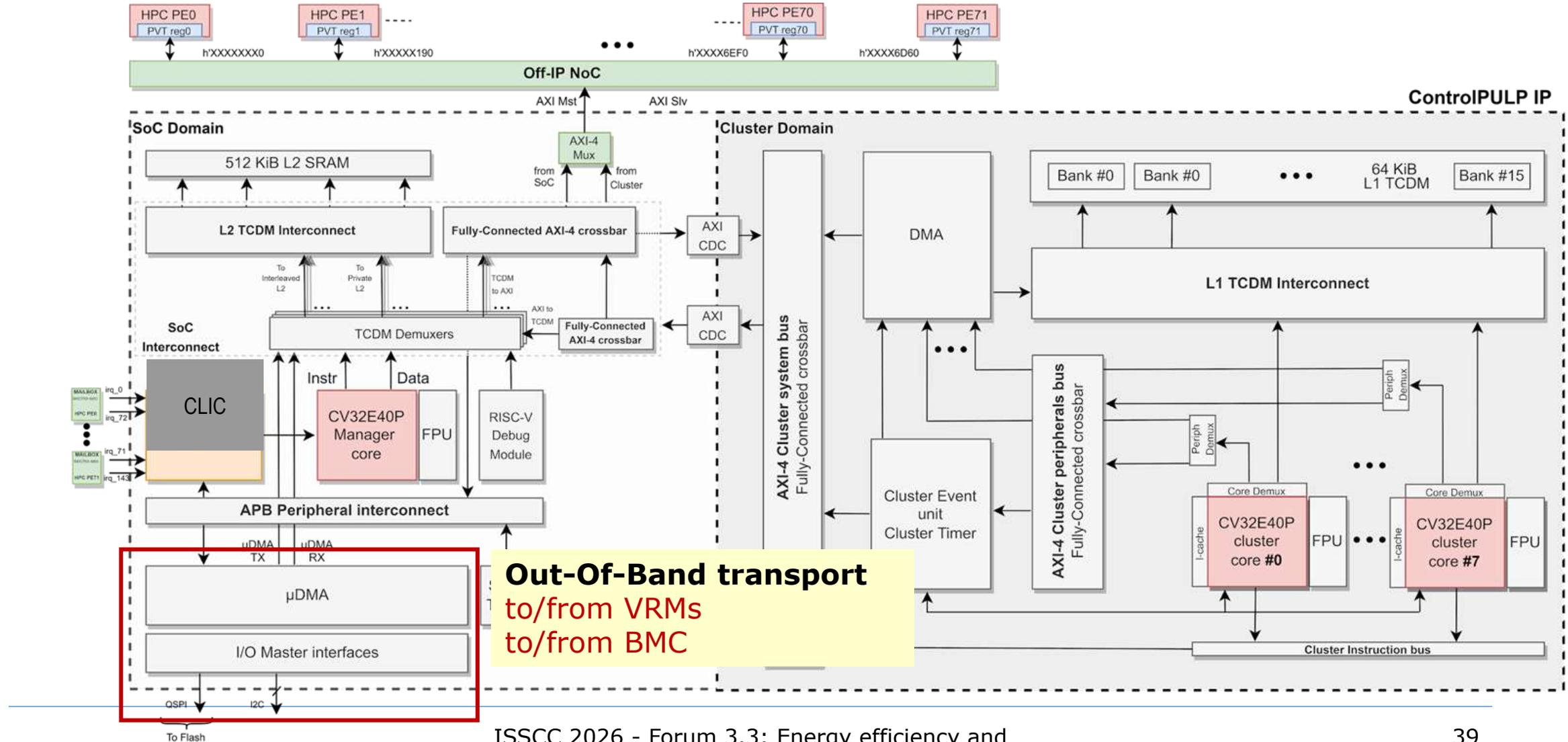**In-band transport**
- to/from PVT sensor registers
- to/from SCMI shared memory

# ControlPULP Architecture



**SCMI mailboxes**
with doorbell notification

**Fast RISC-V interrupt controller CLIC**
up to 256 interrupt lines

# ControlPULP Architecture

# ControlPULP: Software Stack

- Complete (application to hardware) software stack
- With a Real-time operation system, FreeRTOS

| Power Control Firmware |
| :---: |
| PULP FreeRTOS OS |
| PMSIS (BSP, API, DRIVER) |
| Device HAL |
| ControlPULP Hardware |

- Real-time OS (FreeRTOS) schedules tasks with periodic SysTick
- PCF (control policy) example tasks:
  - Periodic Frequency Control task: 2kHz, reads temperatures, computes DVFS, applies frequency on a per-PE basis
  - Periodic Voltage Control task: 8kHz, reads power rails consumption, applies voltage to groups of PEs
  - Communication Control task: SCMI, ACPI, MCTP communication
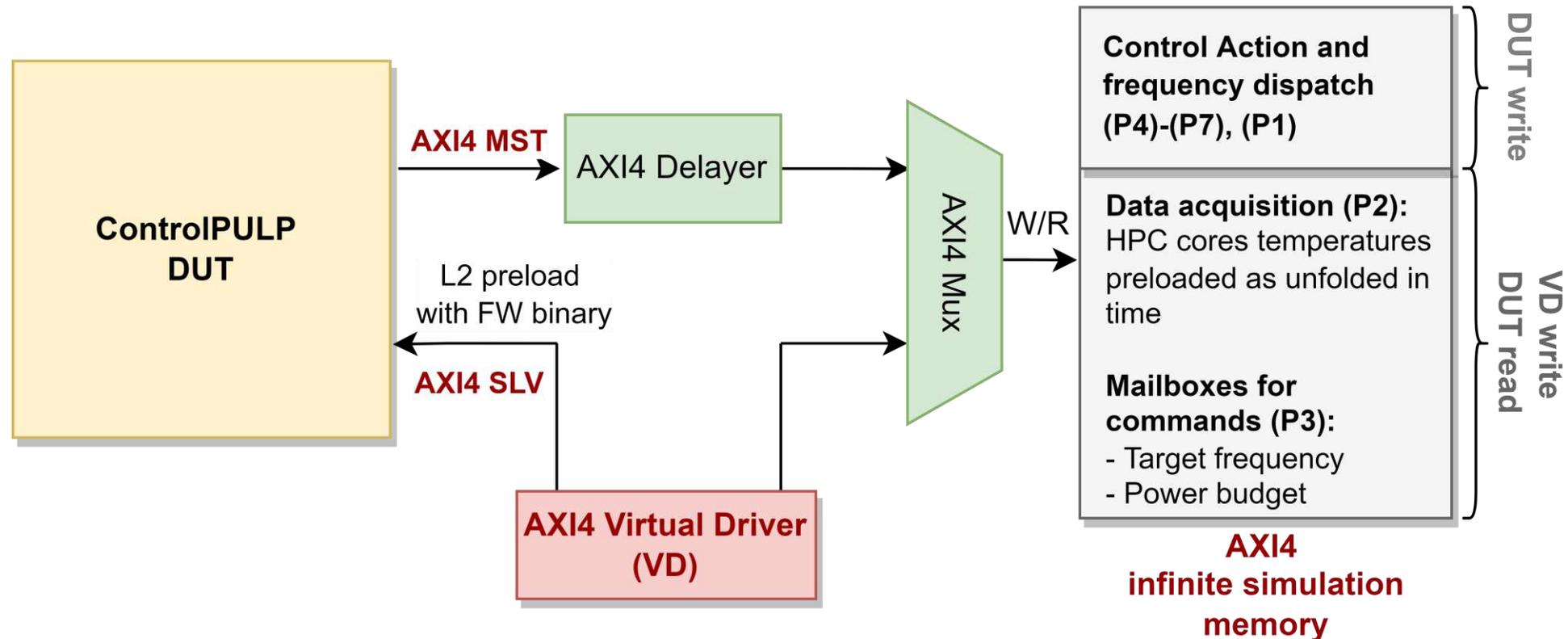
# ControlPULP Implementation

- ☐ GF22 synthesis: one manager core, one cluster, 512KiB + 64KiB @500 MHz, TT
- ☐ Total Area of 9.1 MGE
- ☐ Estimated < 1% of a HPC processor die in modern technology node

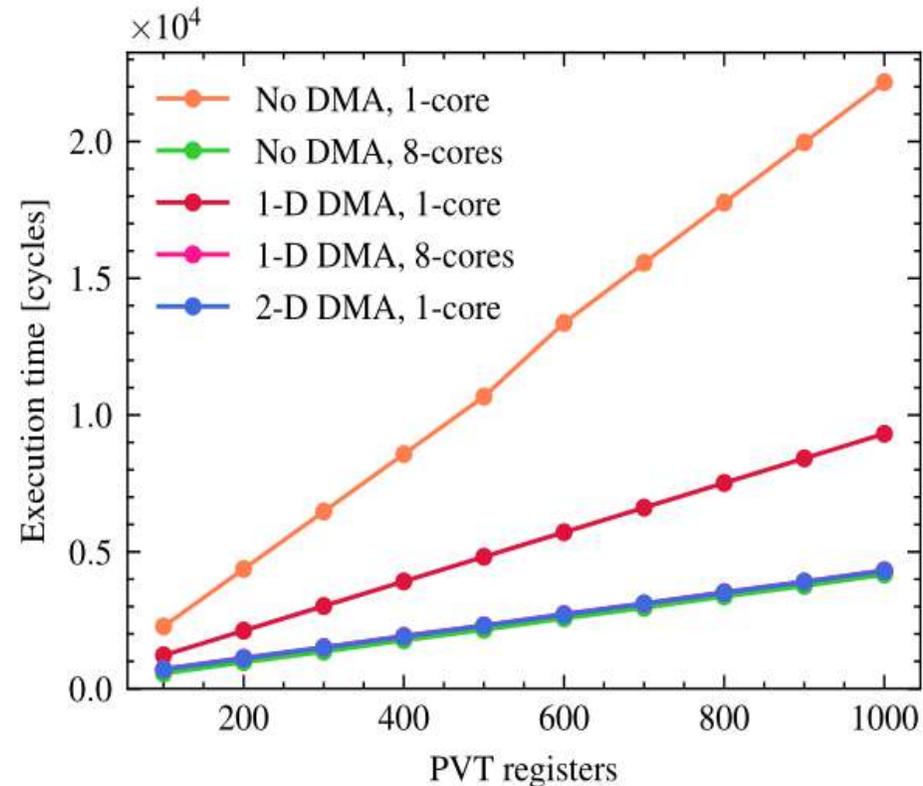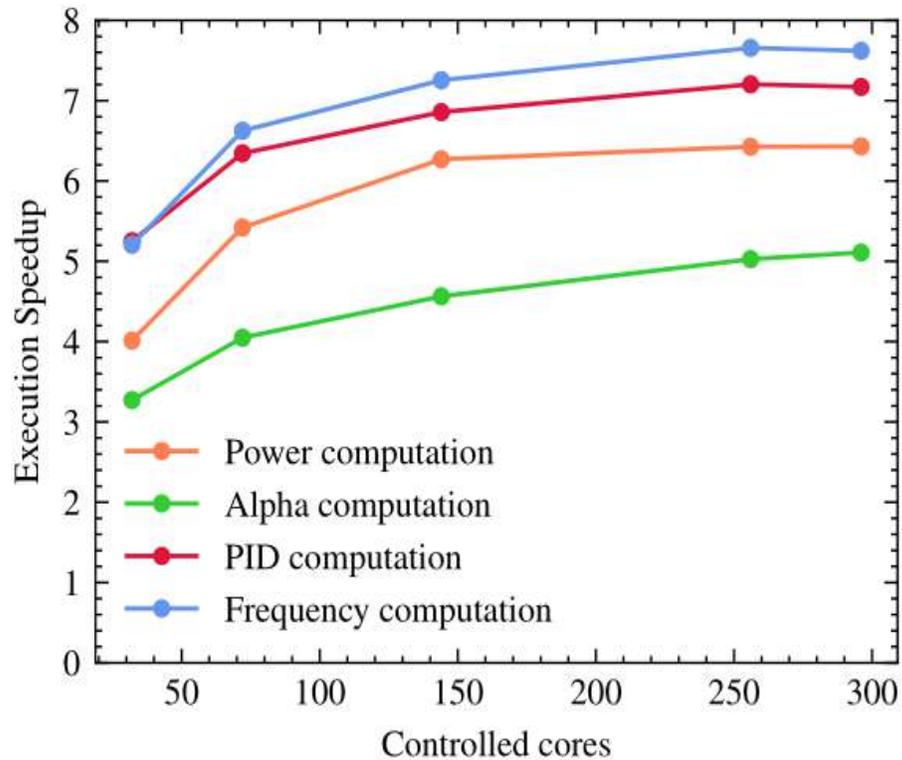| Unit | Area $[\text{mm}^2]$ | Area [kGE] | Percentage [%] |
|---|---|---|---|
| Cluster unit | 0.467 | 2336.7 | 25.5 |
| SoC unit | 0.135 | 675.9 | 7.39 |
| L1 SRAM | 0.119 | 595.7 | 6.51 |
| L2 SRAM | 1.108 | 5542.1 | 60.6 |
| Total | 1.830 | 9150.3 | 100 |

# ControlPULP Performance

☐ NoC latency and controlled system are modeled in a testbench environment
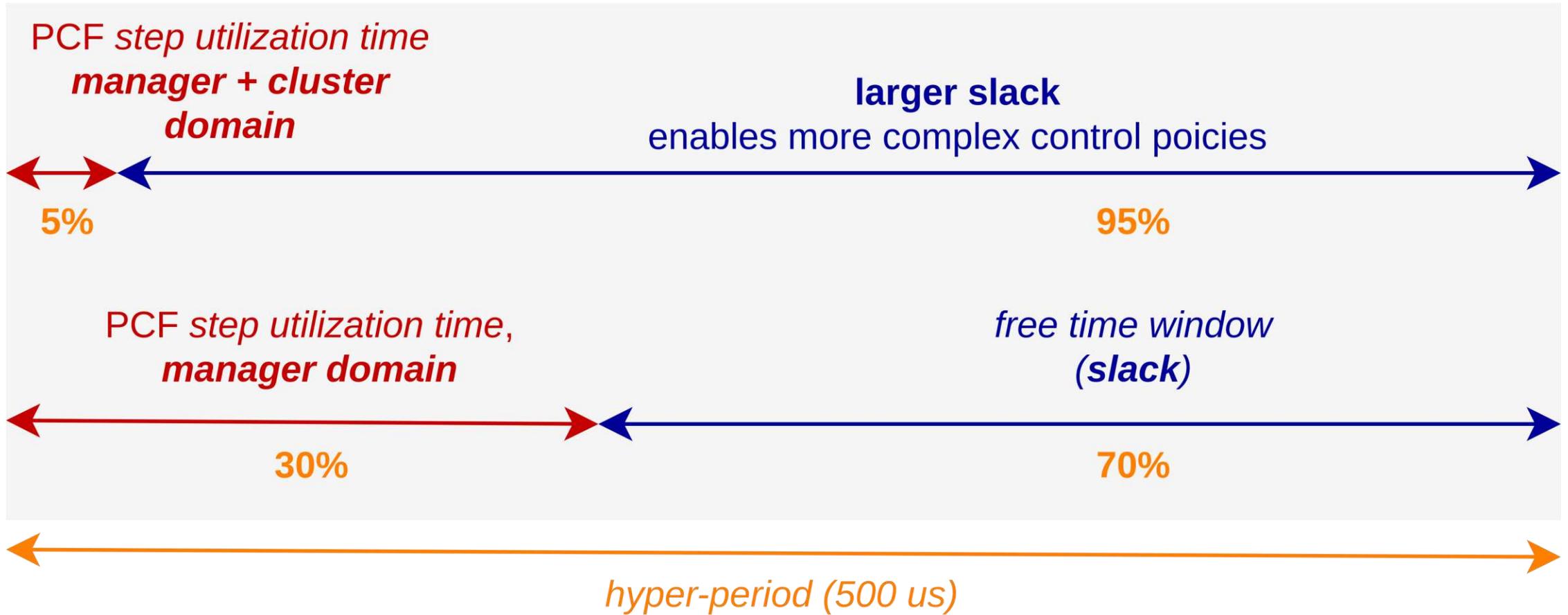☐ Evaluate: multi-core speedup (performance) and interrupt handling reactiveness (latency)

# ControlPULP Performance

☐ Programmable accelerator and DMA driven PCF speedup:
about 5x faster than single-core execution for 72 controlled cores (massive MIMO)

# Why Programmable PM Accelerarator?

☐ Programmable accelerator (cluster) advantage for real-time scenarios



PCF *step utilization time*
**manager + cluster domain**

**larger slack**
enables more complex control poicies

5%                    95%

PCF *step utilization time,*
**manager domain**

*free time window* (**slack**)
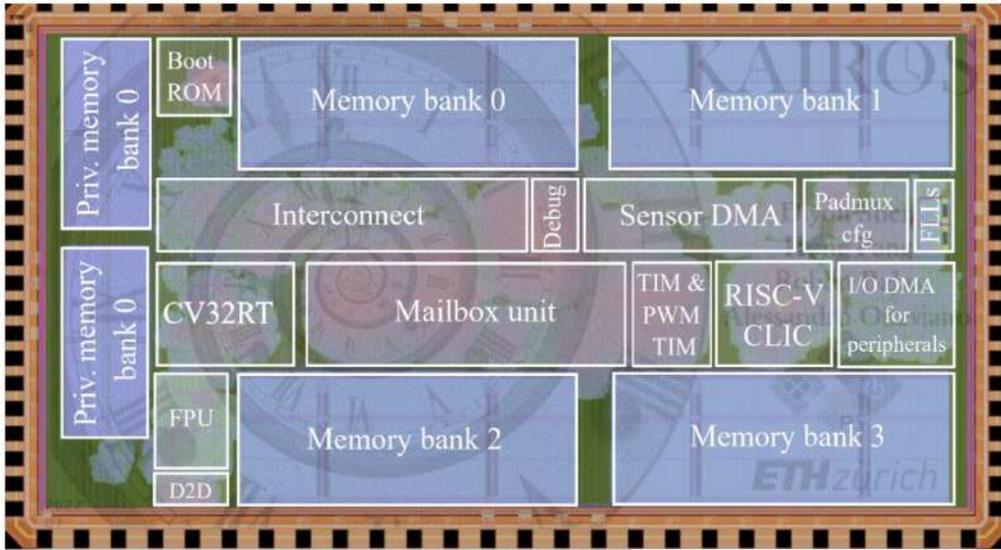
30%                    70%

*hyper-period (500 us)*

# From single to multi Die: ControlPULPlet

- CV32RT: fast interrupt handling in RISC-V
- D2D (die-to-die) connection for chiplet-based power control
- Scatter-gather DMA for N-dimensional periodic sensor reading
- 64x mailboxes for fast message exchange (e.g., SCMI)

# Kairos, a Proxy in TSMC65 CMOS



| | |
|---|---|
| ISA | rv32imafcxpulpv3xfastirq |
| Core | CV32RT [12] |
| On-chip SPM | 448 KiB |
| Technology | TSMC65 |
| Chip area | $7.2\,mm^2$ |
| $V_{DD}$ core/$V_{DD}$ IO | 1.2 V/2.5 V |
| Frequency range | 20 MHz to 380 MHz |
| Power envelope | < 45 mW |



A. Ottaviano et al.  TVLSI25

# Outline

- ☐ Boosting efficiency for AI workloads
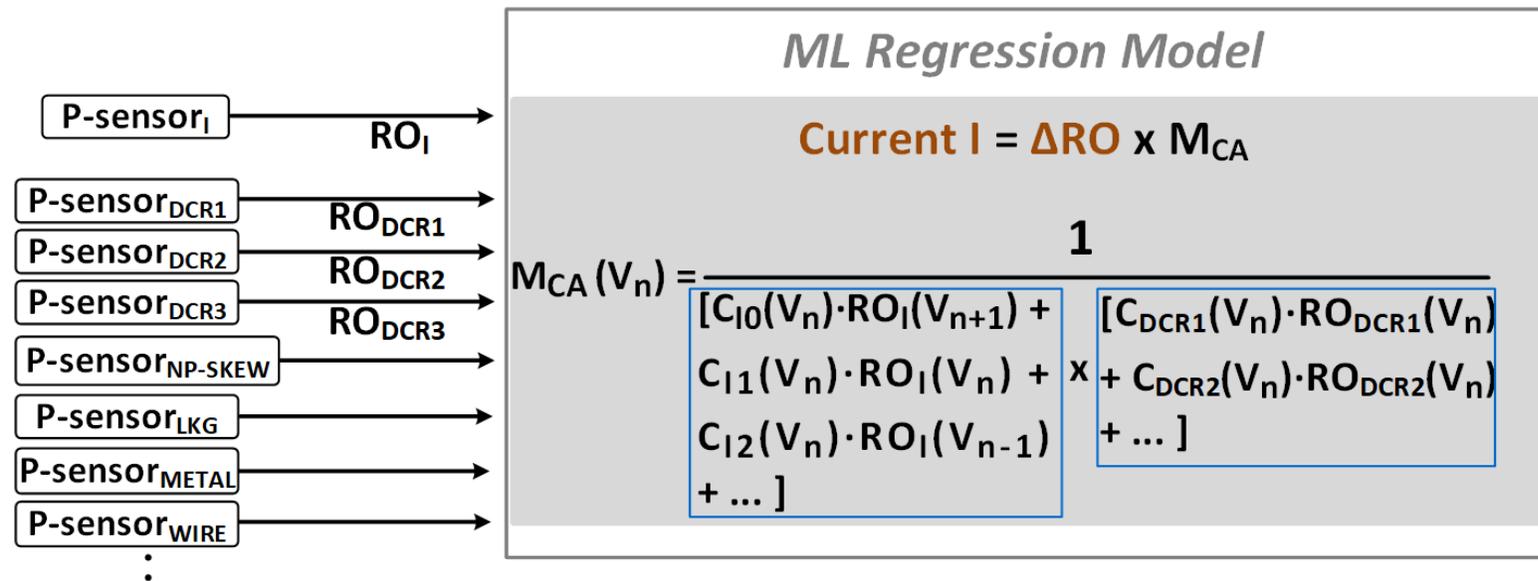
- ☐ Managing idleness and heterogeneity in accelerated systems

- ☐ Using AI for managing AI

- ☐ Conclusions, future perspectives
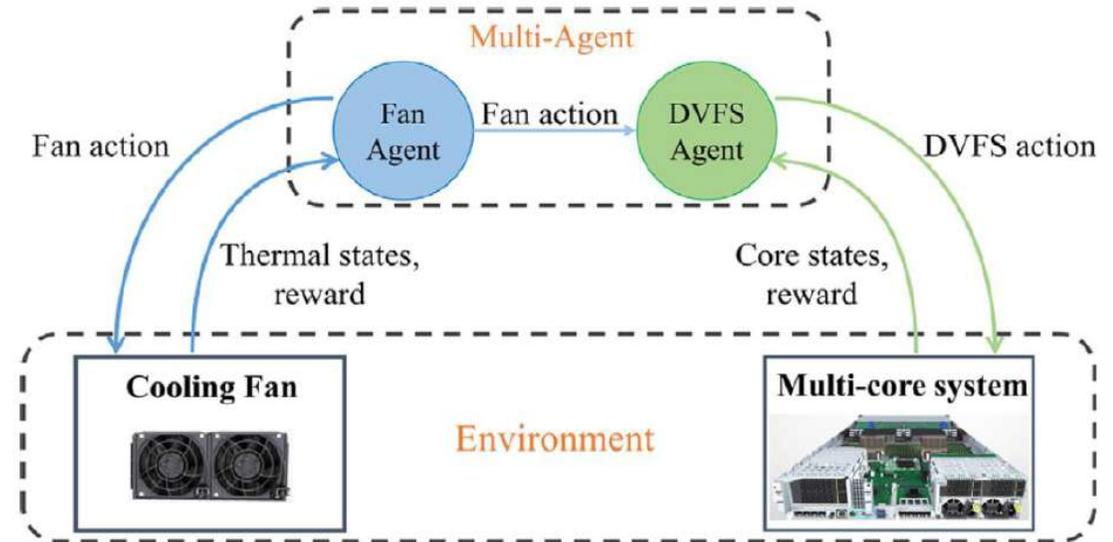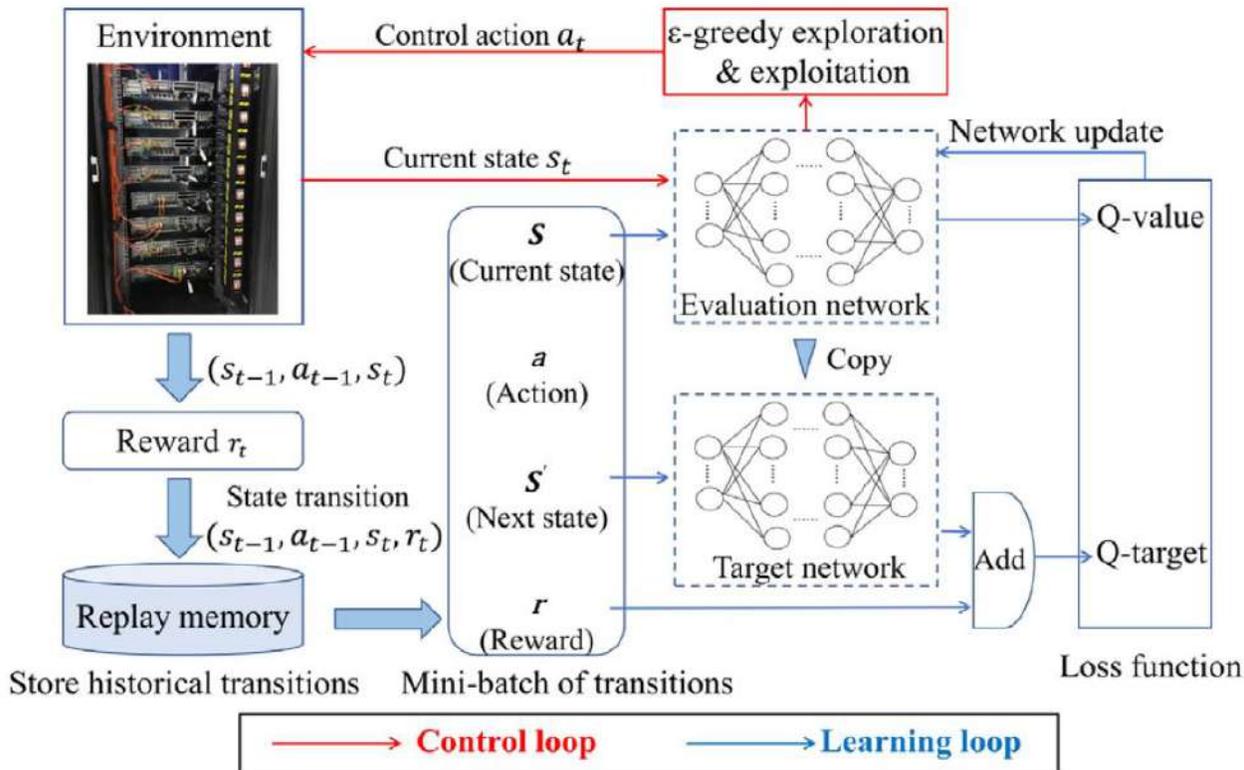
# ML Based learning for Non-linear Sensors

**Current sensing using Count-to-Ampere Models (MCA) based on Ring-oscillator sensors**

**ML Regression Model**

$$\text{Current I} = \Delta RO \times M_{CA}$$

P-sensor$_I$ → RO$_I$

P-sensor$_{DCR1}$ → RO$_{DCR1}$

P-sensor$_{DCR2}$ → RO$_{DCR2}$

P-sensor$_{DCR3}$ → RO$_{DCR3}$

P-sensor$_{NP-SKEW}$

P-sensor$_{LKG}$

P-sensor$_{METAL}$

P-sensor$_{WIRE}$

$$M_{CA}(V_n) = \frac{1}{[C_{I0}(V_n) \cdot RO_I(V_{n+1}) + C_{I1}(V_n) \cdot RO_I(V_n) + C_{I2}(V_n) \cdot RO_I(V_{n-1}) + \ldots] \times [C_{DCR1}(V_n) \cdot RO_{DCR1}(V_n) + C_{DCR2}(V_n) \cdot RO_{DCR2}(V_n) + \ldots]}$$

- ☐ ML calibration model leverages silicon data to refine the MCA
  - ■ Scope-measured current I as golden to train the coefficients of ΔRO (RO) from the I-sensor (P-sensors)
  - ■ Multi-scenarios, including generic benchmarks like Geekbenchv6, Antutu, and SpecInt2k6 etc.
- ☐ Good tolerances against voltage and temperature fluctuations
- ☐ Statistical ML used to eelect the most relevant RO terms (i.e., P-sensors) to simplify the coefficient terms of the model
- ☐ Device-specific fine tuning also possible

Mediatek ISSCC25 8.1

# Advanced AI-based Policies

☐ Control law can be learned (e.g. RL) offline and fine-tuned online



www.sciencedirect.com/science/article/abs/pii/S2210537924000222

# Outline

- ☐ Boosting efficiency for AI workloads

- ☐ Managing idleness and heterogeneity in accelerated systems

- ☐ Using AI for managing AI

- ☐ Conclusions, future perspectives

# Summing Up & Future Perspectives

- ☐ PM is more critical than ever for AI hardware
  - ■ Thermal bottleneck
  - ■ Power delivery bottleneck
  - ■ Impacts the AI bottom line (cost per token)
- ☐ Key requirements for PM in AI hardware
  - ■ Massive MIMO control (fine-grained & scaled-up)
  - ■ Heterogeneous (processors, accelerators, PIM, PIN…)
  - ■ Multi-die (2.5D → 3.5D)
- ☐ <span style="color:red">Real-time</span>, (multi-die) <span style="color:red">scalable</span> and <span style="color:red">programmable</span> PM controller architecture is needed (open platform is advisable)
- ☐ AI for power managing AI is a major trend
  - ■ Learn complex non-linear functions for sensing and actuation
  - ■ Learn optimal policies (e.g. RL)
  - ■ Online learning for fine-tuning device-specific calibration