



Ten years of PULP: The Evolution of the Species from IoT to HPC

Davide Rossi

davide.rossi@unibo.it

PULP Platform

Open Source Hardware, the way it should be!



@pulp_platform

pulp-platform.org

youtube.com/pulp_platform



PULP: The Origins

- Project started in **2013**
- A collaboration between **University of Bologna** and **ETH Zürich**
 - We wanted to start with a clean slate, no need to remain compatible to legacy systems,
no dependencies with any commercial IP
- **Original Academic/Research goal**
 - Push energy efficiency of IoT computing systems as much as possible
 - Create a compute platform and an ecosystem used for research on Computer Architectures by our team as well as other groups in the World
- **Original Approach**
 - Exploit Open Source IPs as much as possible
 - “We’ll never do a processor”
 - “We’ll never do a compiler”
 - Team of 5-6 people at the beginning

Back to 2013 Open (And Close) Source Processors Landscape



- Reasonably Good Quality IPs
- Area around 50 kGates
- Open Source
- Compiler Support
- Community support
- We decided to go with OR1200

Architectural Requirements:

- DTCM interface
- PCACHE
- Blocking bus interface
- User-defined extension interface
- Floating point unit
- NOT obfuscated RTL
- Instruction Set Simulator

A12000

PROCESSOR	PCACHE	BUS	EXT IF	F.P. UNIT	RTL VISIBLE	ISS	AREA (Kg)
STxp70	YES	T3	YES	YES	YES	YES	42+
REISC4	YES	CUSTOM	NO	NO	?	?	?
ARM CORTEX M0+	NO	AMBA	NO	NO	(NO)	YES	~9
ARM CORTEX M3	NO	AMBA	NO	NO	(NO)	YES	~27
ARM CORTEX M4	NO	AMBA	NO	YES	(NO)	YES	~38 (FPU)
ARC EM4	NO	AMBA	YES	YES	(NO)	YES	10 - 20
ARC EM6	YES	AMBA	YES	YES	(NO)	YES	20 - 40
LATTICE Mico32	YES	WISHB.	NO	NO	NO	NO	~30
OPENRISC	YES	AMBA	NO	YES	YES	NO	~61
							~50

TARGET G-CC
LLVM

→ EXECUTORS SOFTWARE

COMPILATION BUONA QUALITÀ DEN

First Presence of PULP @ ORCONF in 2013

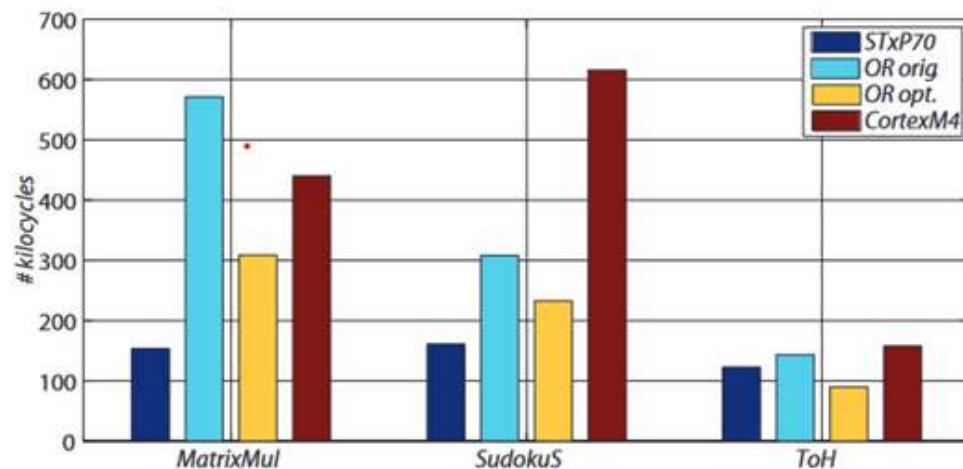
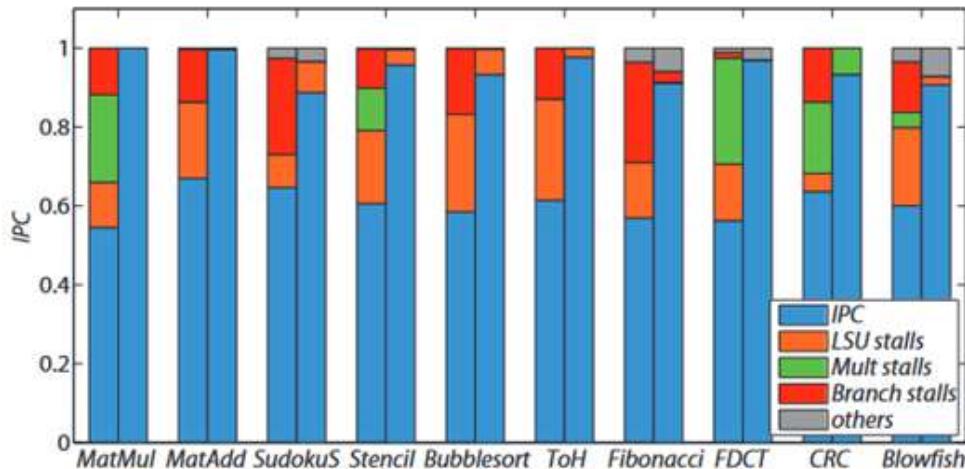
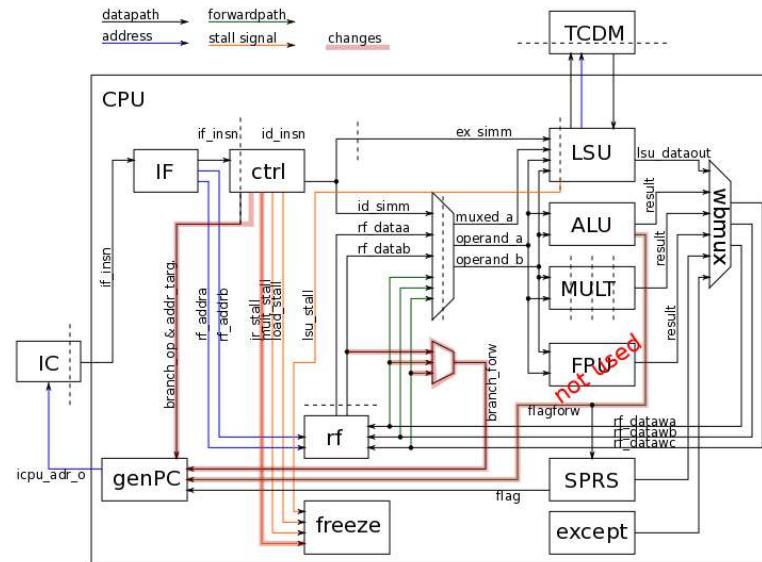


ORCONF 2013 5-6/10/2013

OR1200 Microarchitectural Analysis (And Optimizations)



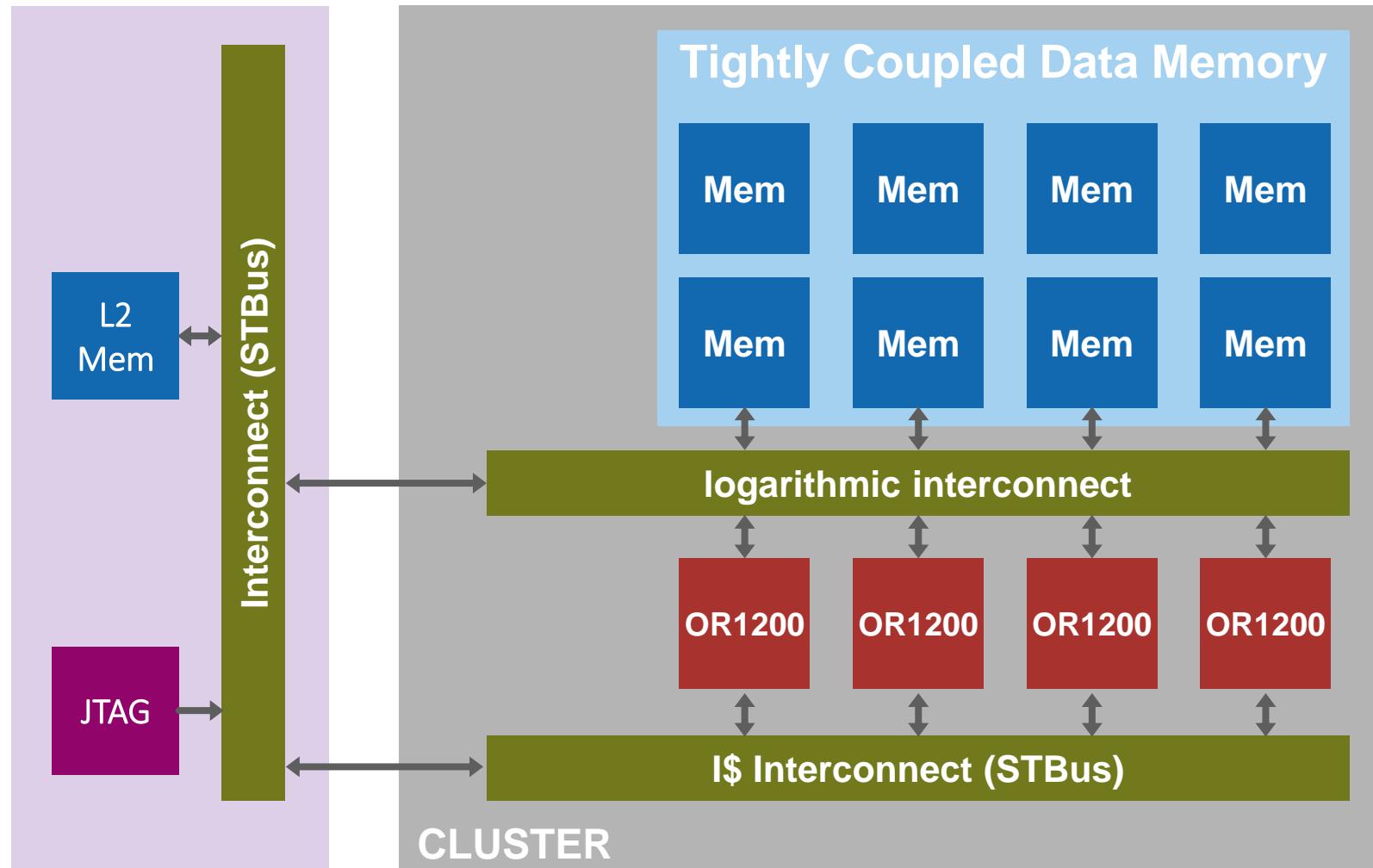
- Blocking LSU (2 cycles)
- Blocking Multiplier (3 cycles)
- I\$ → TCDM Combinational Path
- No Support for compressed instructions



The first “PULP cluster”: PULPv1 - TO Dec 2013



- 4-cores + I\$ cluster
- Logarithmic interco.
- I\$ + L2 Interco.
- L2 memory
- JTAG

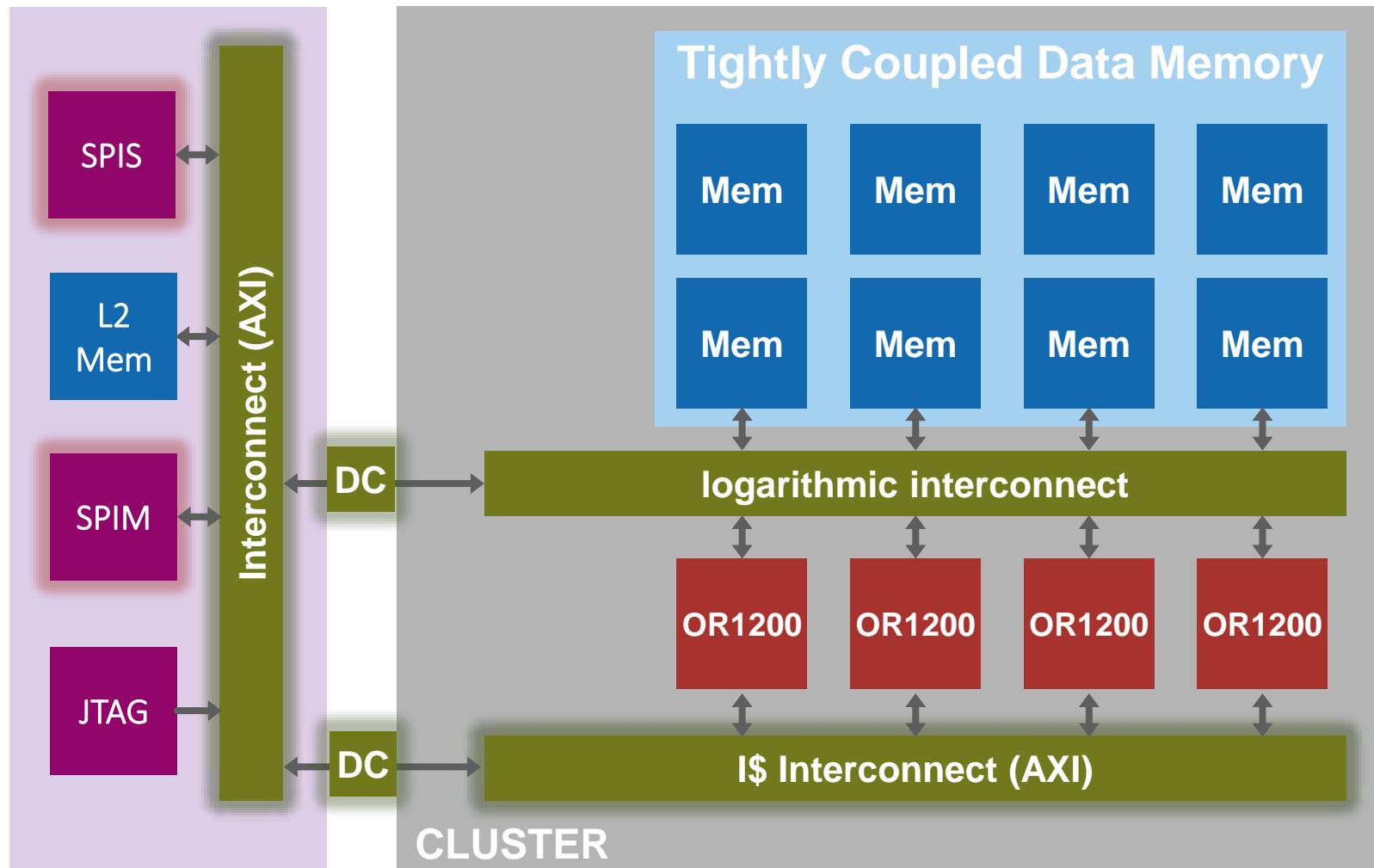


....not really “Open Source”...😊

The Evolution of the Species: PULPv2 - TO Dec 2014



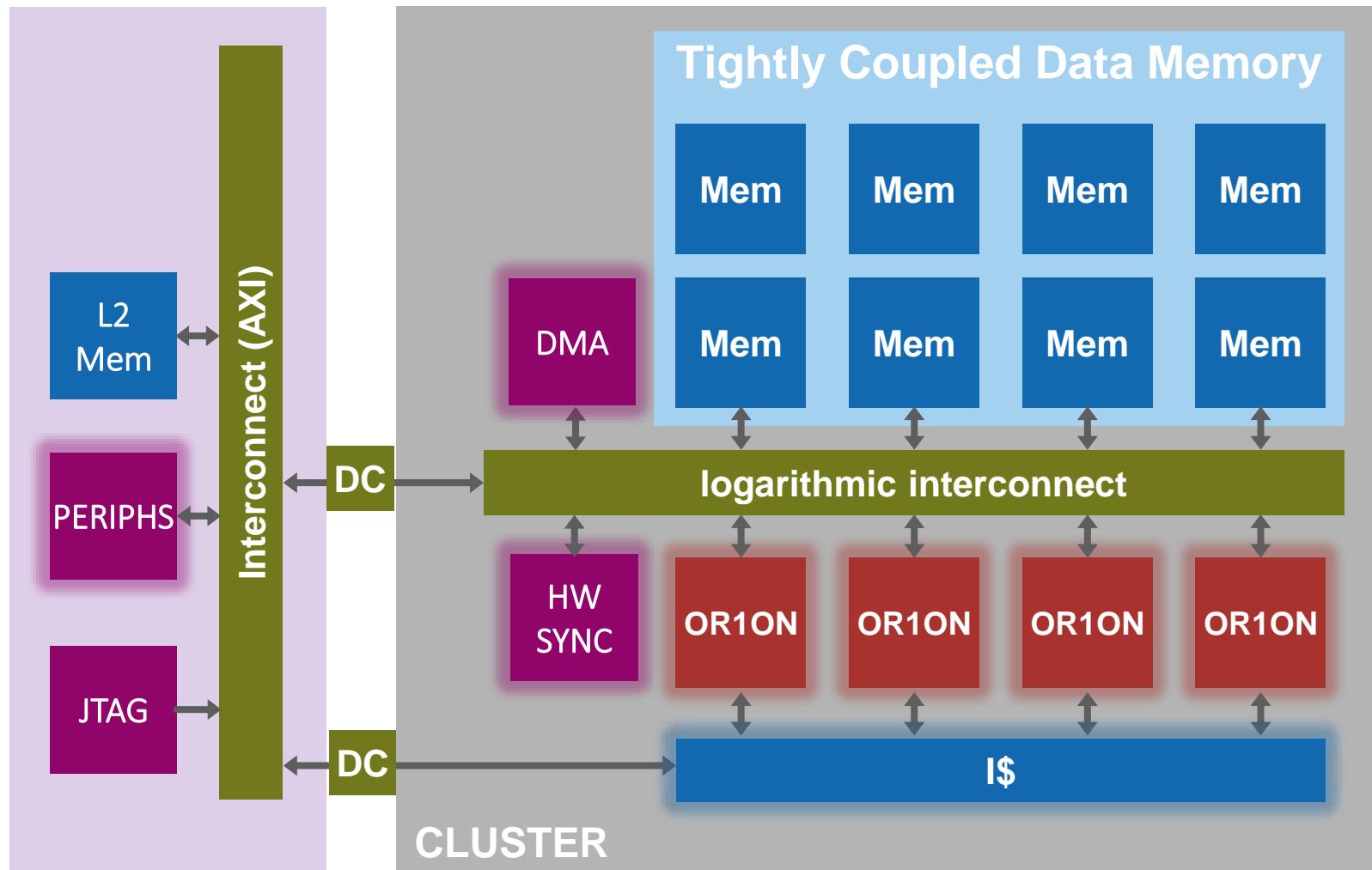
- PULPv1 +
- New (very rudimental) but Open Source AXI interconnect
- Some new peripherals (SPIM, SPIS)
- Instruction Cache
- DVFS-ready



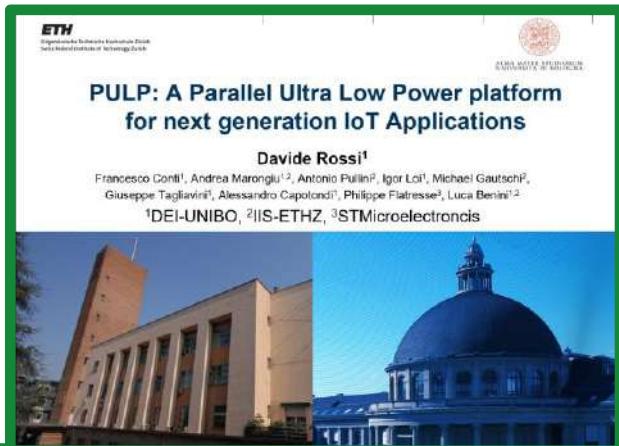
The Evolution of the Species: PULPv3 - TO Dec 2015



- PULPv2 +
- ORION cores with DSP Extensions (brand new core)
- Hardware Synchronizer
- Instruction Cache
- Standard Peripherals



PULP 1,2,3 at Hot Chips 27, 2015 + Other Open Source Contribs.



MIAOW: An Open Source GPGPU
www.miaowgpu.org

Vinay Gangadhar, Raghu Balasubramanian, Mario Drumond, Ziliang Guo,
Jai Menon, Cherin Joseph, Robin Prakash, Sharath Prasad, Pradip Vallathol,
Karu Sankaralingam

Vertical Research Group
University of Wisconsin - Madison



Raven: A 28nm RISC-V Vector Processor with
Integrated Switched-Capacitor DC-DC
Converters and Adaptive Clocking

Yunsup Lee, Brian Zimmer, Andrew Waterman,
Alberto Puggelli, Jaehwa Kwak, Ruzica Jevtic, Ben Keller,
Steve Bailey, Milovan Blagojevic, Pi-Feng Chiu,
Henry Cook, Rimas Avizienis, Brian Richards,
Elad Alon, Borivoje Nikolic, Krste Asanovic

University of California, Berkeley



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA



Wind of Change....First RISC-V Workshop, January 2015

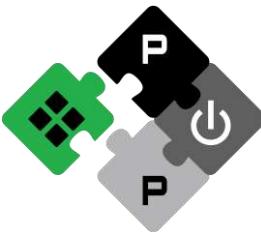


- We started to push the accelerator...
 - Switch to RISC-V (RTL, compilers, etc...)
 - size of the group increased (20-30 people)

January 14-15, 2015
Marriott Hotel, Monterey, CA



PULP - Open HW: RISC-V Cores,... and more



RISC-V Cores			
RI5CY	Ibex	Snitch	Ariane + Ara 64b
32b	32b	32b	64b

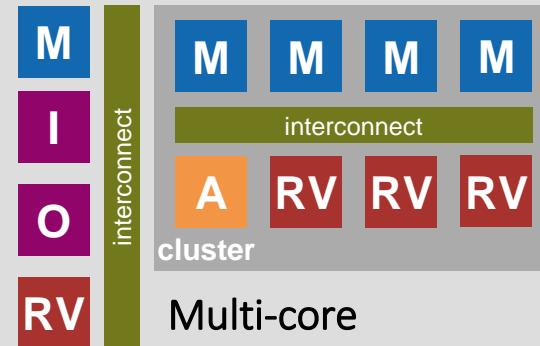
Peripherals	
JTAG	SPI
UART	I2S
DMA	GPIO

Interconnect
Logarithmic interconnect
APB – Peripheral Bus
AXI4 – Interconnect

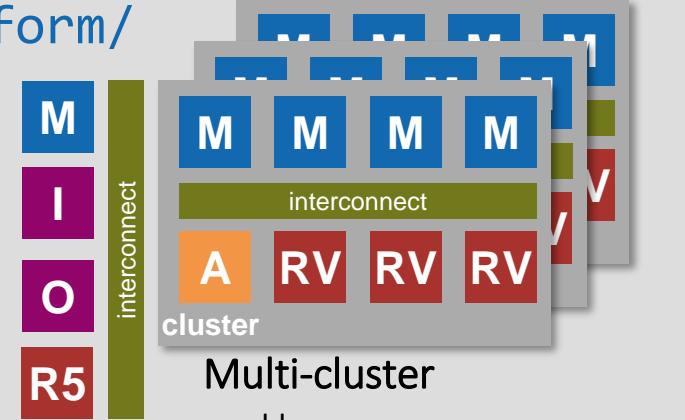
Platforms



- Single Core
 - PULPino
 - PULPissimo



- Multi-core
 - Control-PULP
 - Kraken



- Multi-cluster
 - Hero
 - Occamy

IOT

Accelerators

HWCE,1,2,3
(convolution)

Neurostream
(ML train)

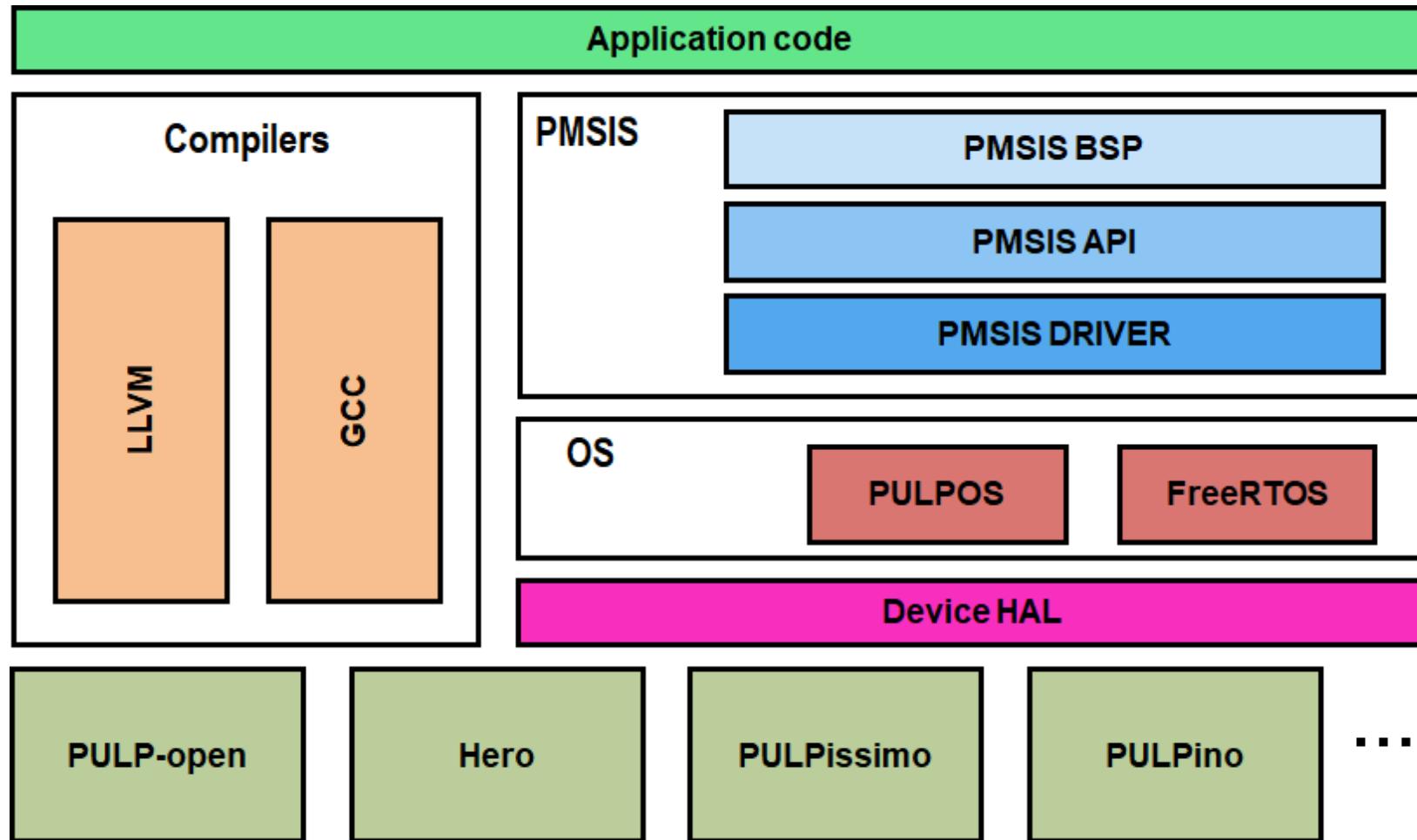
FPNEW
(64/32/16/8)

ARA
(DFP vector)

HPC



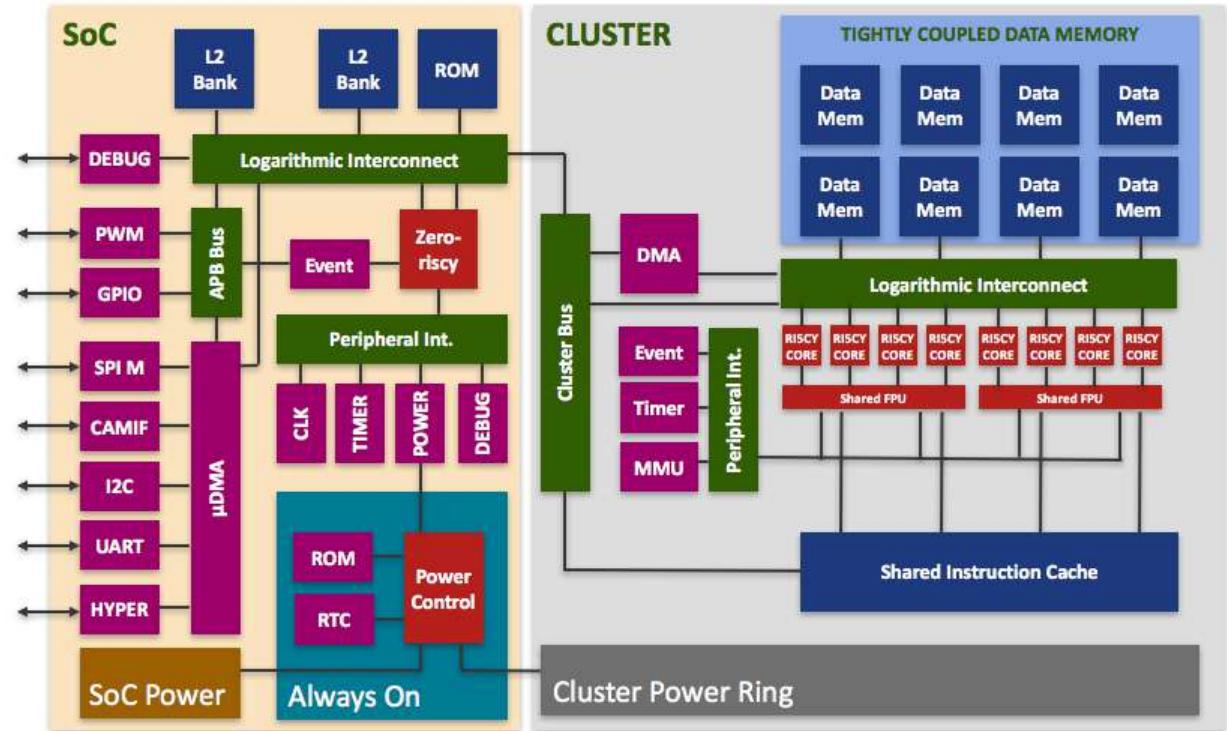
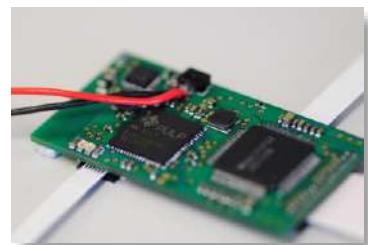
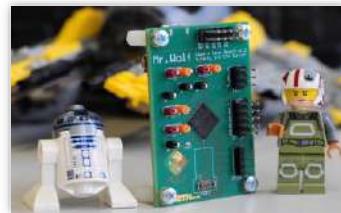
PULP Software Stack



Mr. Wolf SoC: First Board-Ready Chip (Dec 2018)



- First SoC with Fabric controller: 32-bit RISC-V processor (Zero-RISKY)
- Autonomous IO DMA Subsystem
- Rich set of peripherals
- Parallel Programmable Accelerator:
 - First pulp chip with RI5CY
 - First Floating-Point Capable PULP chip



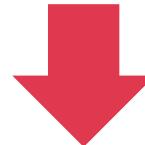
GREENWAVES
TECHNOLOGIES

bitcraze

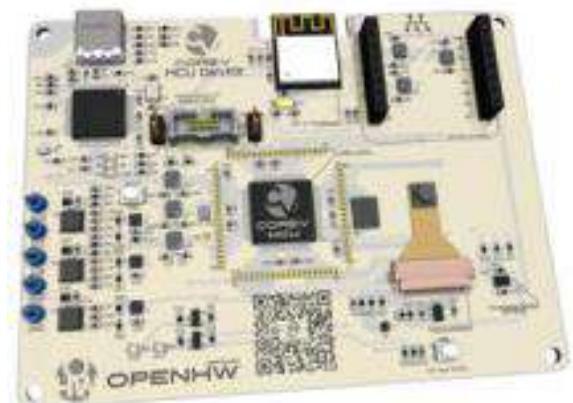
From Academia to Industry



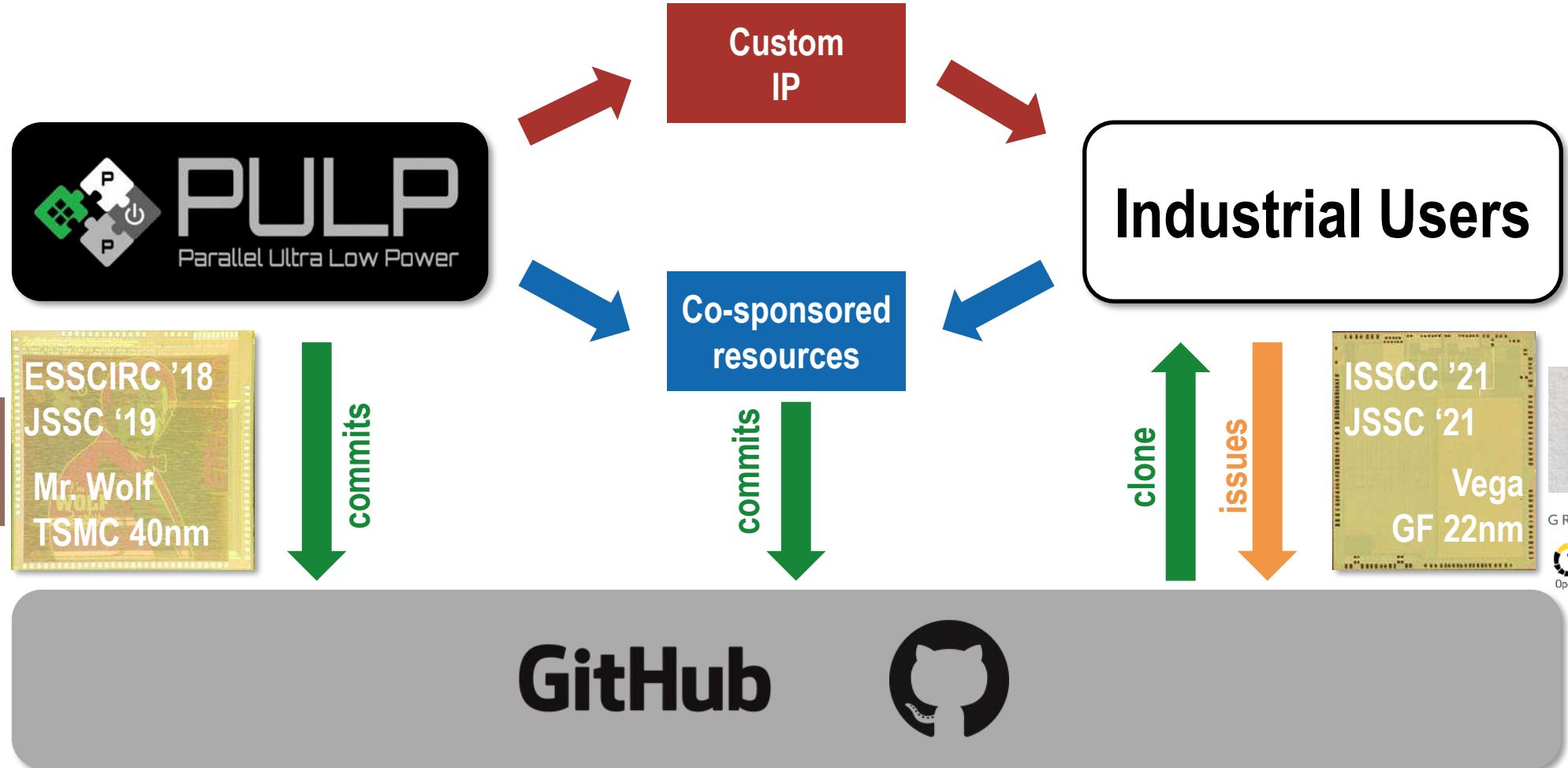
- **OpenHW Group** is a not-for-profit, global organization (EU, NA, Asia) where HW and SW designers collaborate in the development of open-source cores, related IP, tools and SW such as the **Core-V** family
- **OpenHW Group** provides an infrastructure for hosting high quality open-source HW in line with industry best practices.



RI5CY, ARIANE RISC-V cores,
FPU, AXI4 Components.



Industrial Paths for PULP IPs & Mutual Benefits





IoT & TinyML Applications and Chips

Vega: IoT Heterogeneous SoC



- Transprecision Floating-Point Units (16/32-bit float)¹
- 4 MB non-volatile MRAM (Weight Storage)²
- Programmable Cognitive Wake-up Unit based on HD Computing

Presented at ISSCC 2021

ISSCC 2021 / SESSION 4 / PROCESSORS / 4.4

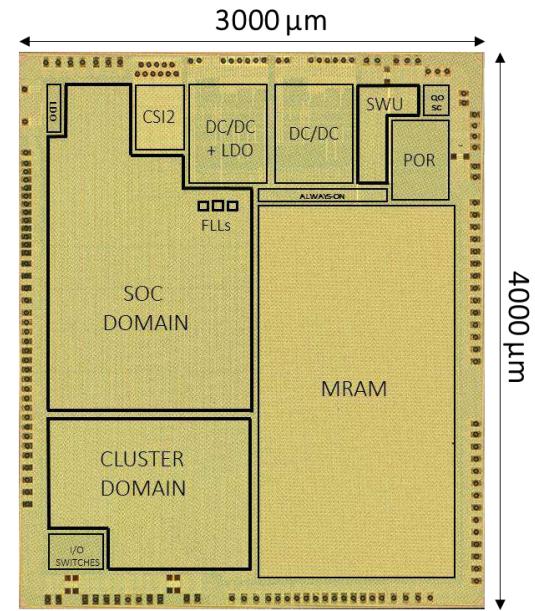
4.4 A 1.3TOPS/W @ 32GOPS Fully Integrated 10-Core SoC for IoT End-Nodes with 1.7 μ W Cognitive Wake-Up From MRAM-Based State-Retentive Sleep Mode

Davide Rossi¹, Francesco Conti¹, Manuel Eggimann², Stefan Mach², Alfio Di Mauro², Marco Guermandi^{1,3}, Giuseppe Tagliavini¹, Antonio Pullini^{2,3}, Igor Loi³, Jie Chen^{1,3}, Eric Flamand^{2,3}, Luca Benini^{1,2}

¹University of Bologna, Bologna, Italy

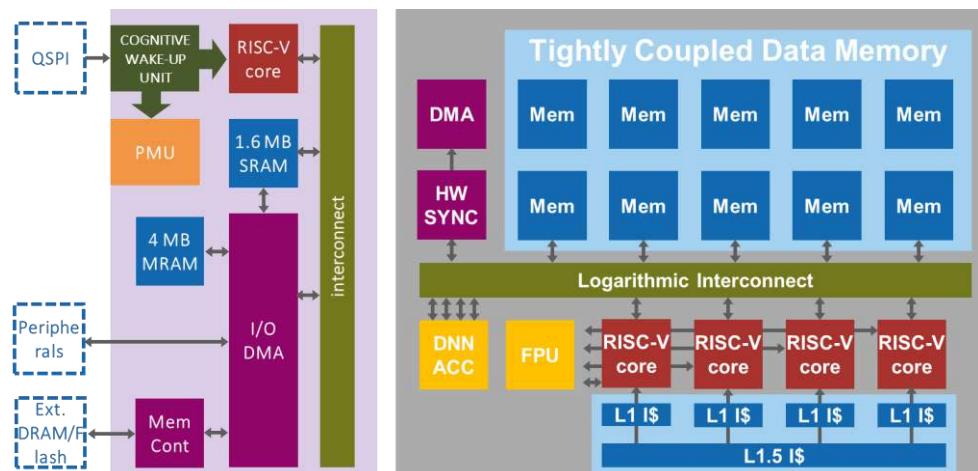
²ETH Zurich, Zurich, Switzerland

³Greenwaves Technologies, Grenoble, France



Prototype implemented in GF 22FDX

flip-well LVT & SLVT cells, 12mm² for Full SoC





State-of-the-Art Comparison

	SleepRunner [2]	Mr.Wolf [3]	SamurAI [4]	VEGA
Embedded NVM	-	-	SoA Int & FP Perf & Efficiency; Boost of 4.3x & 2.8x compared to Wolf	4 MB MRAM
Wake-up Sources	WiC	GPIO, RT		GPIO, RTC, Cognitive
Best Int Perf.	31 MOPS (32b)	12.1 GOPS	1.5 GOPS	15.6 GOPS
Best.Int Eff. @ Perf.	97 MOPS/mW (32b) @ 18.6 MOPS (32b)	190 GOPS/W @ 3.8 GOPS	230 GOPS/W @110 MOPS	614 GOPS/W @7.6 GOPS
Best FP Perf.		1 GFLOPS		4 GFLOPS
Best FP Eff. @Perf	-	18 GFLOPS/W @350 MFLOPS	-	158 GFLOPS/W @ 2 GFLOPS
Best ML Perf.			36 GOPS	32.2 GOPS
Best ML Eff. @Perf	-	-	1.3 TOPS/W @ 2.8 GOPS	1.3 TOPS/W @15.6 GOPS



Open Software Stack for VEGA: *QuantLab*, *Nemo*, *Dory*

QuantLab

Quantization Laboratory

NEMO

NEural Minimization for pytOrch

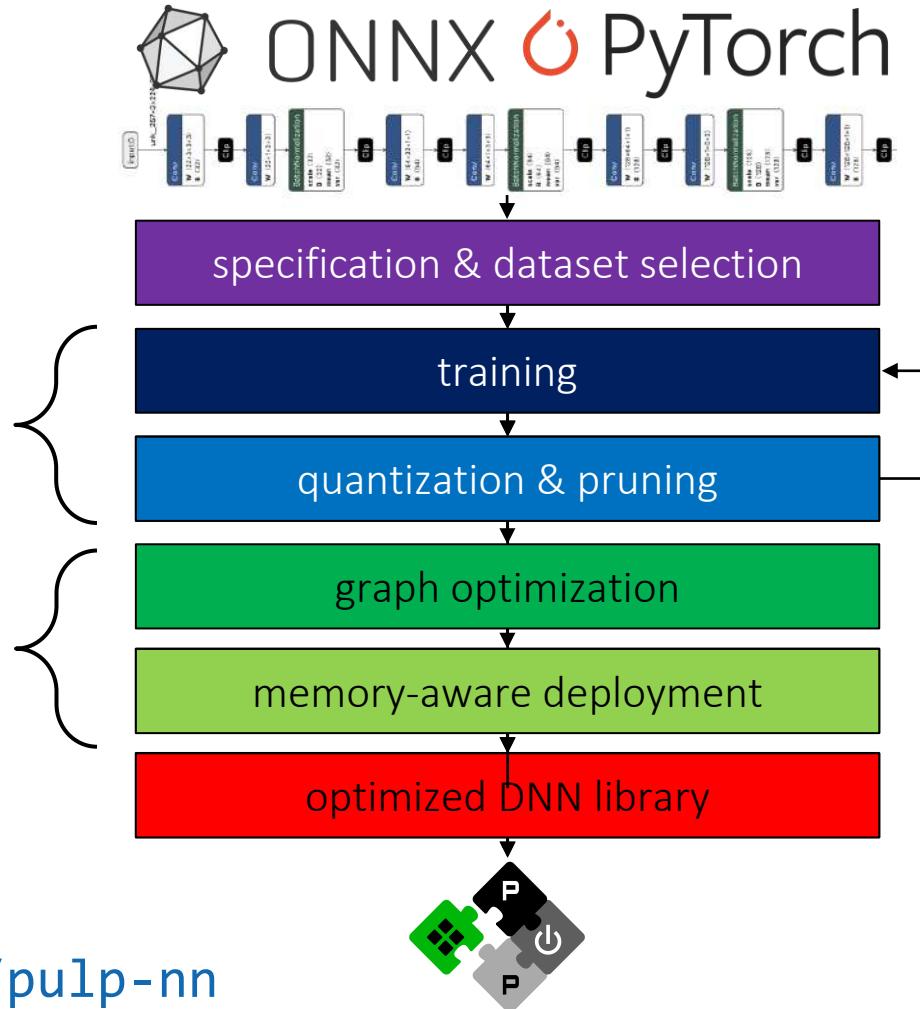
DORY

Deployment Oriented to memoRY

PULP-NN

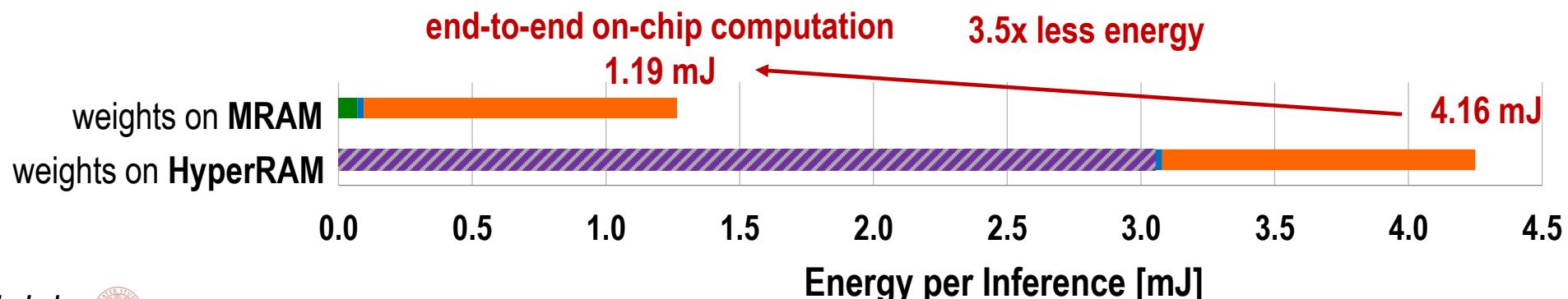
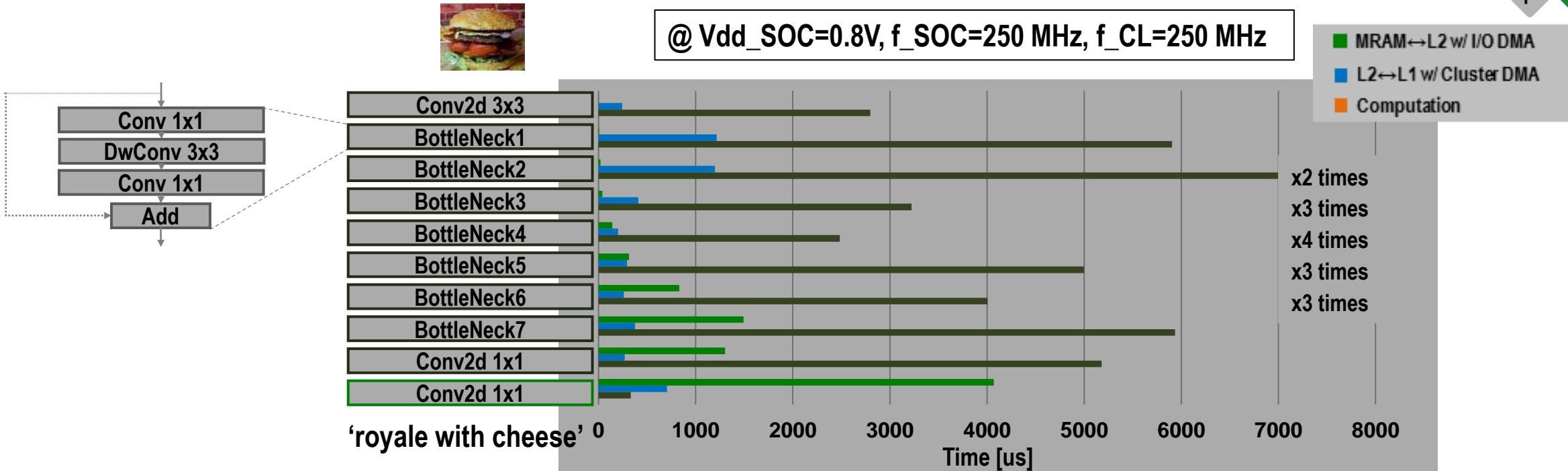
PULP Neural Network backend

github.com/pulp-platform/nemo,/dory,/pulp-nn





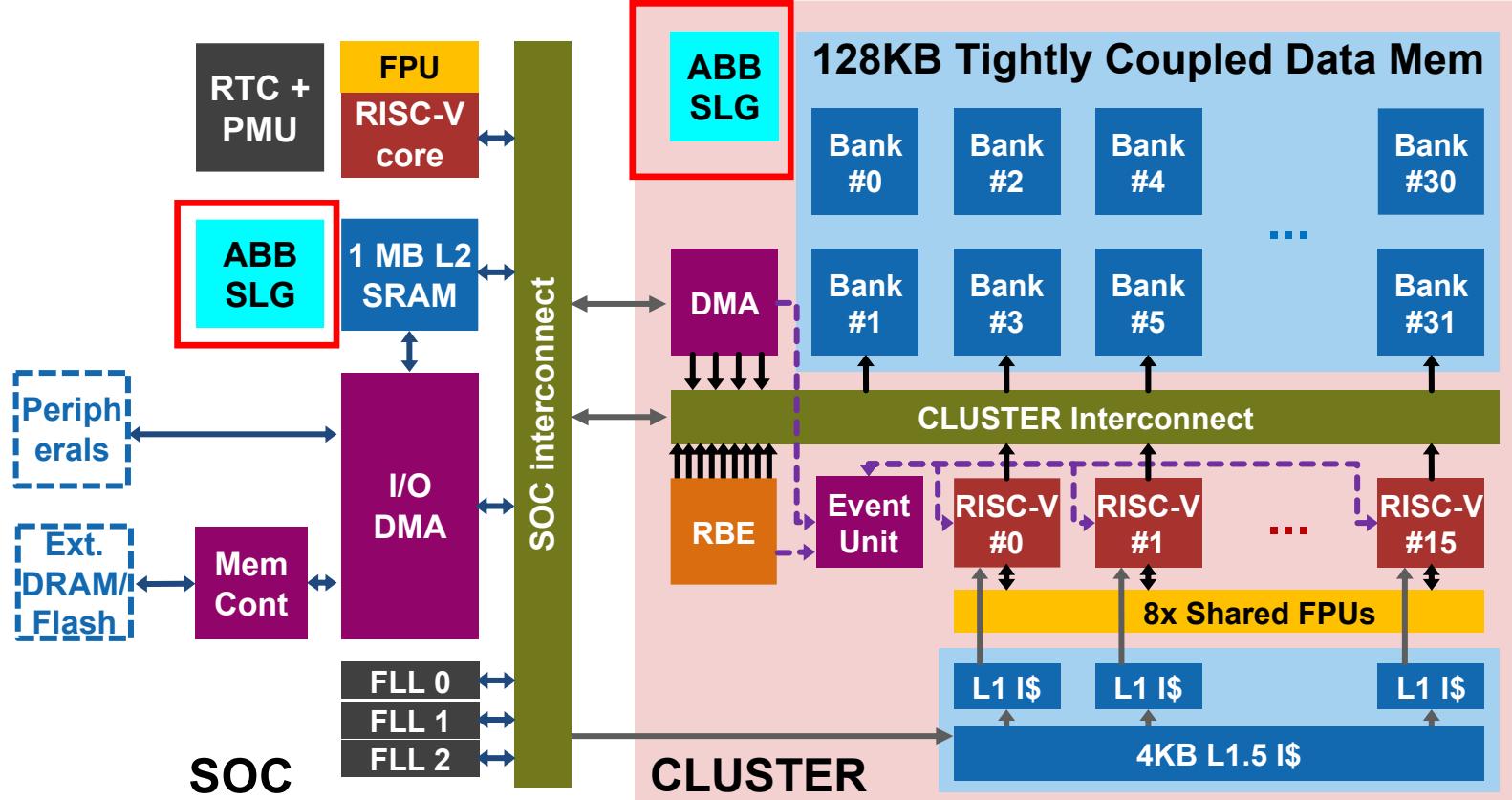
First Fully Integrated MobileNetV2 on an IoT SoC



MARSELLUS: AI-IoT Heterogeneous SoC



- Machine Learning ISA Extensions
- 2-8b Reconfigurable Binary Engine for 3x3, 1x1 DNN kernels
- Adaptive Body Biasing with on-the-fly control



22.1 A 12.4TOPS/W @ 136GOPS AI-IoT System-on-Chip with 16 RISC-V, 2-to-8b Precision-Scalable DNN Acceleration and 30%-Boost Adaptive Body Biasing
Francesco Conti¹, Davide Rossi¹, Gianna Paulin², Angelo Garofalo¹, Alfio Di Mauro², Georg Ruetishauer², Gianmarco Ottavi¹, Manuel Eggimann², Hayate Okuhara¹, Vincent Huard³, Olivier Montfort³, Lionel Jure³, Nils Exibard³, Pascal Gouedo³, Mathieu Louvat³, Emmanuel Botte³, Luca Benini^{1,2}

Presented at ISSCC 2023

State-of-the-Art Comparison (HW Accelerated)

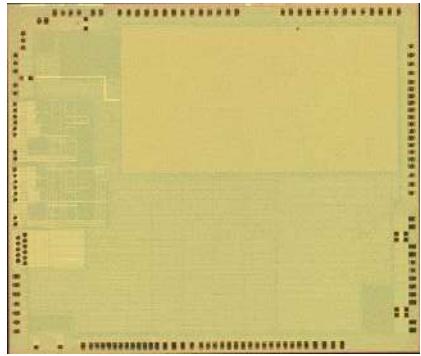


	VEGA [1] ISSCC'21	SAMURAI [2] VLSIC'20	DIANA [3] ISSCC'22	MARSELLUS <i>(this work)</i>
Technology	22nm FDX	28nm FD-SOI	22nm FDX + AIMC	22nm FDX
Die Area	10 mm ²	4.5 mm ²	10.24 mm ²	18.7 mm ² (cluster 1.9 mm ²)
Cores	10x RV32IMCFXpulp + HWCE	1x RV32IMCFXpulp + Digital Accel.	1x RV32CIMFXpulp + Digital Accel. + AIMC	1x RV32IMCFXpulp + 16x RV32IMCFXpulpnn + RBE
Max Frequency	450 MHz	350 MHz	320 MHz	420 MHz
Power range	1.7 uW - 49.4 mW	6.4 uW - 96 mW	10-129 mW (digital)	12.8 mW - 123 mW
Best SW (INT) Perf	15.6 GOPS (8 RISC-V)	1.5 GOPS (1 RISC-V)	-	90 GOPS (16 RISC-V M&L 2x2b, 0.8V+ABB)
Best SW (INT) Eff	614 GOPS/W @ 7.6 GOPS (8 RISC-V)	230 GOPS/W @ 110 MOPS (1 RISC-V)	-	1.66 TOPS/W @ 19 GOPS (16 RISC-V M&L 2x2b)
Best HW-Accel Perf	32.2 GOPS (HWCE)	36 GOPS (Dig)	180 GOPS (Digital), DNN Acceleration AIMC	637 GOPS (RBE 2x2b, 0.8V+ABB)
Best HW-Accel Eff	1.3 TOPS/W @ 15.6 GOPS (HWCE)	1.3 TOPS/W @ 2 TOPS/W (Digital Accel.)	>10x w.r.t. VEGA 500 TOPS/W (AIMC)	12.4 TOPS/W @ 136 GOPS (RBE 2x2b)

IoT SoC Playground: Success Stories



*Fixed weight
precision to 8/16-bit*



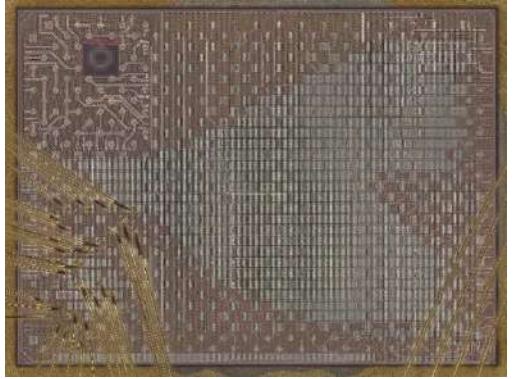
GREENWAVES
TECHNOLOGIES

**Vega 22nm, ISSCC'21
(UNIBO + GreenWaves)**



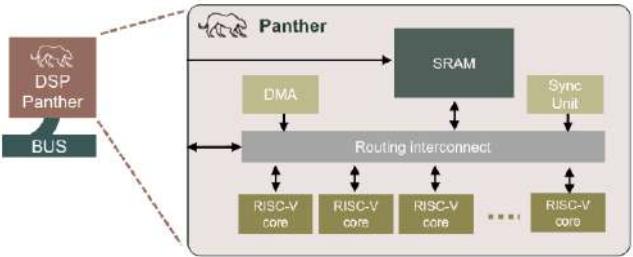
*Fixed weight
precision to 8/16-bit*

**GAP9 SoC (commercial)
65% of GAP9 is based on
PULP open-source IPs**



DOLPHIN
DESIGN

**Marsellus 22nm, ISSCC'23
(UNIBO + Dolphin Design)**



**Panther DSP IP (commercial)
90% HW and 20% SW of Phanter based on
PULP open-source IPs**

Higher-Performance PULP?

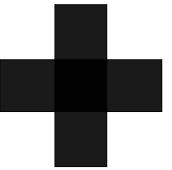


- So far, focus on **low-power** applications: near-sensor processing, nano-UAVs, etc...
- Two niches left out: **Linux-capable low-power MPUs**...
 - can exploit clusters for performance boost
- ... and what about **high-performance computing**?
 - energy efficiency of capital importance (**power = \$\$\$** spent for **cooling, energy bill**)
- **Hardware accelerators** provide a key technology for **HPC**
 - compute-dominated workloads
 - highly parallel workloads
 - efficiency in Joules/op and power envelope in kW are important metrics
 - flexibility is also of primary importance



HOST

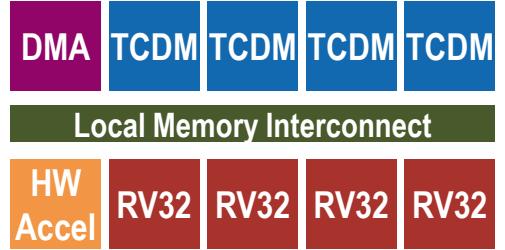
- General-purpose
- Linux-capable
- Versatility
- Programmability



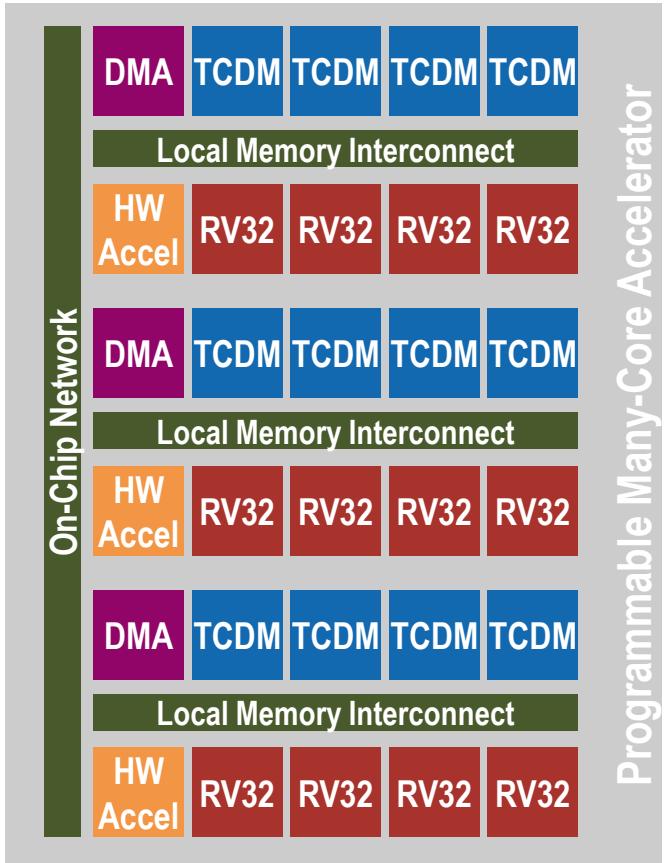
PMCA

- Parallel Manycore Accelerator (PMCA)
- Domain-specialized
- Energy-efficient

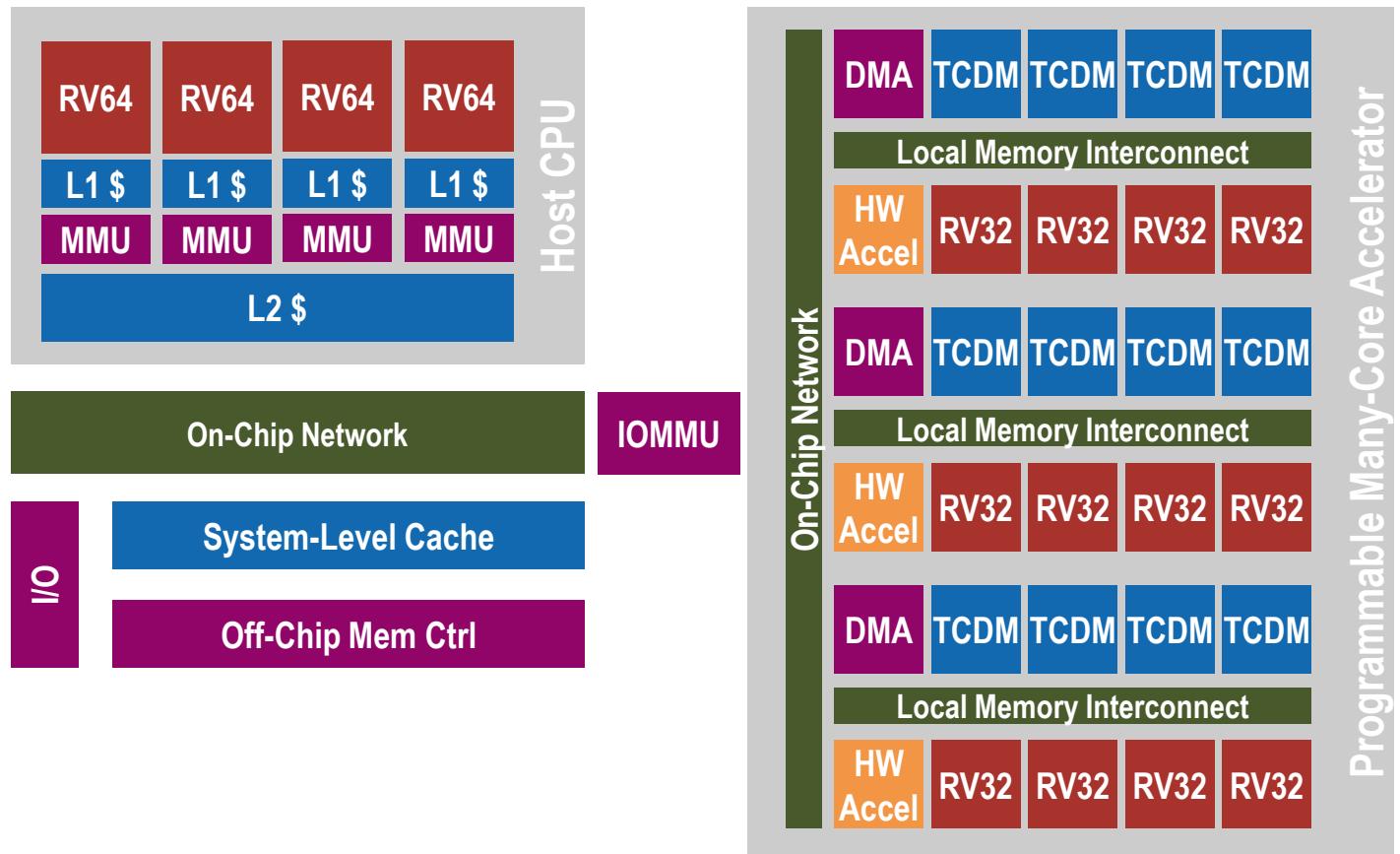
HERO: Hardware Architecture



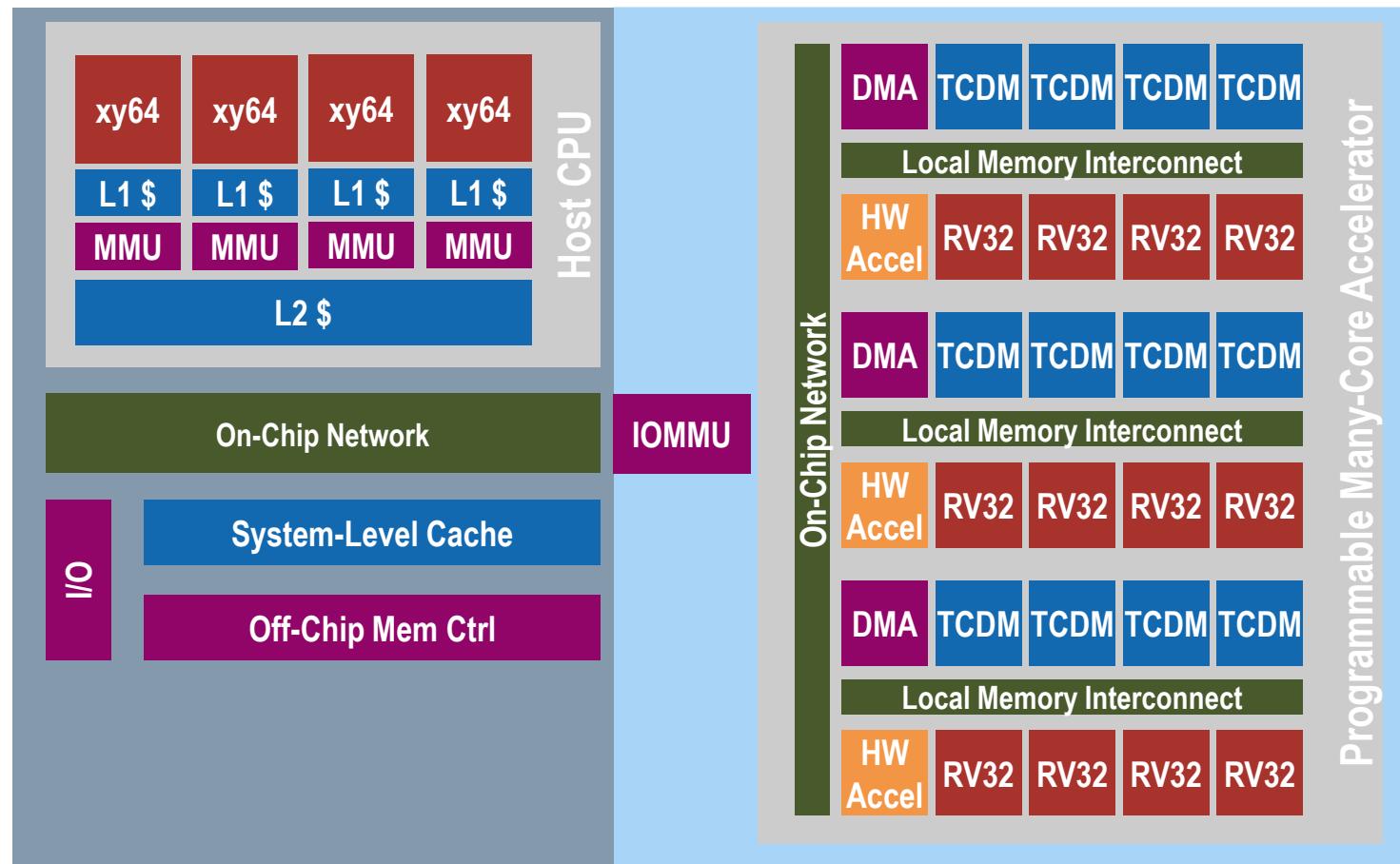
HERO: Hardware Architecture



HERO: Hardware Architecture



HERO: Hardware Architecture



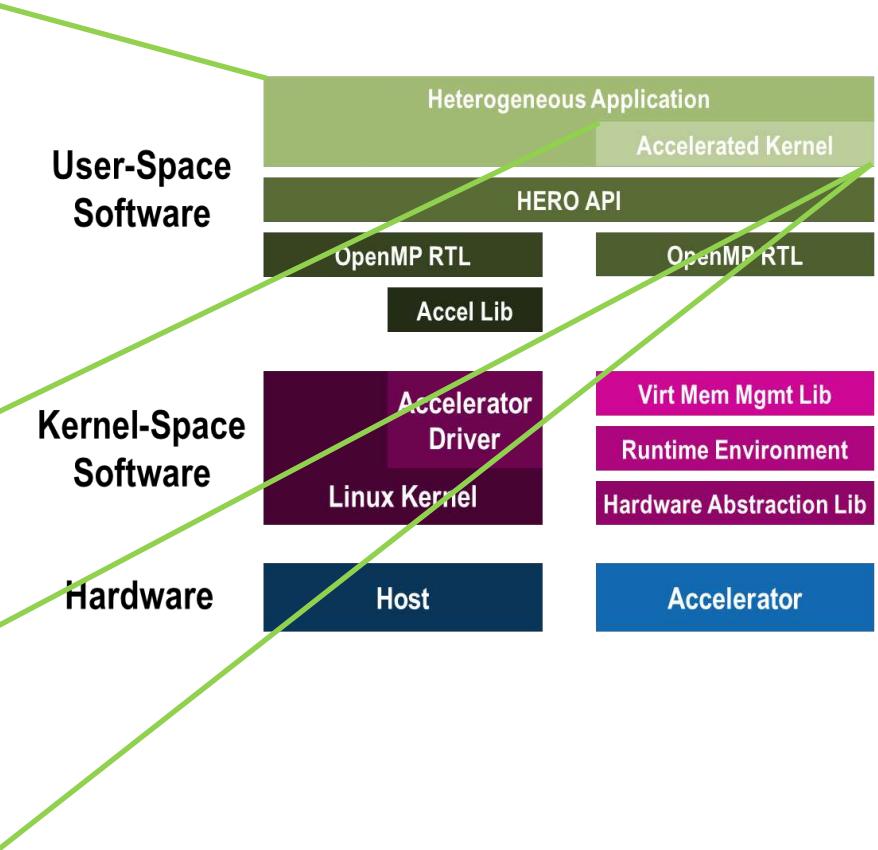
L. Valente et. al., "HULK-V: a Heterogeneous Ultra-low-power Linux capable RISC-V SoC",
DATE 2023



HERO: Software Architecture

Principle: **single-source** heterogeneous programming. Offload with OpenMP 4.5 target semantics.
Example:

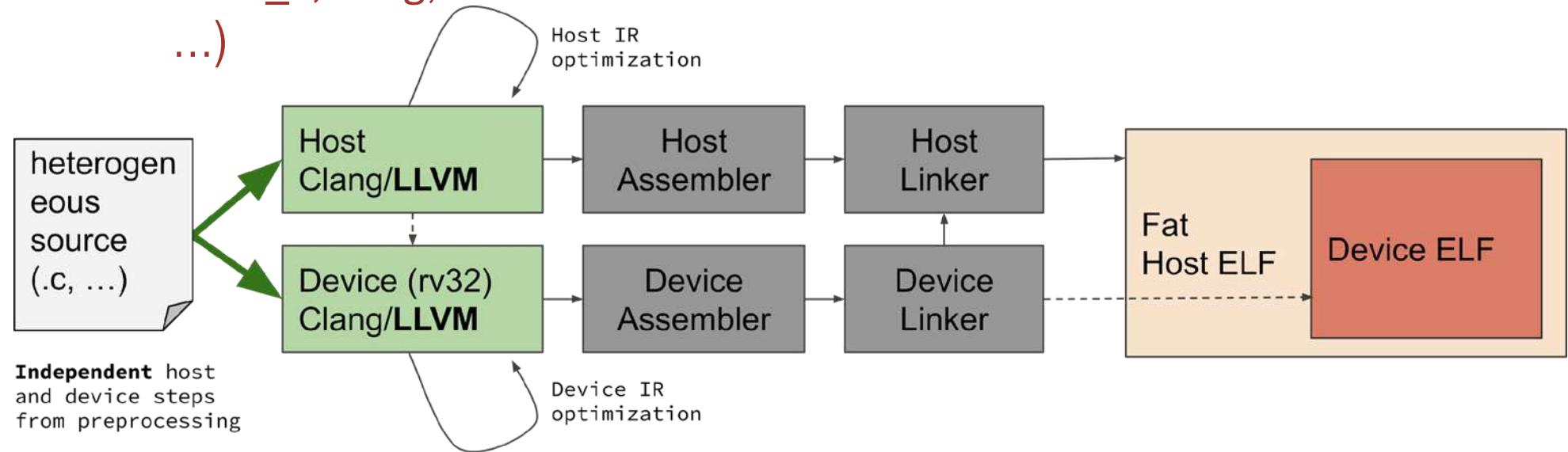
```
int main ()  
{  
    vertex vertices[N];  
    load(&vertices, N);  
    #pragma omp target map(tofrom: vertices)  
    {  
        #pragma omp parallel  
        for (i = 0; i < N; ++i)  
            process(vertices[i]);  
    }  
}
```



HERO: Heterogeneous Compilation

Single-source, single-binary heterogeneous compilation

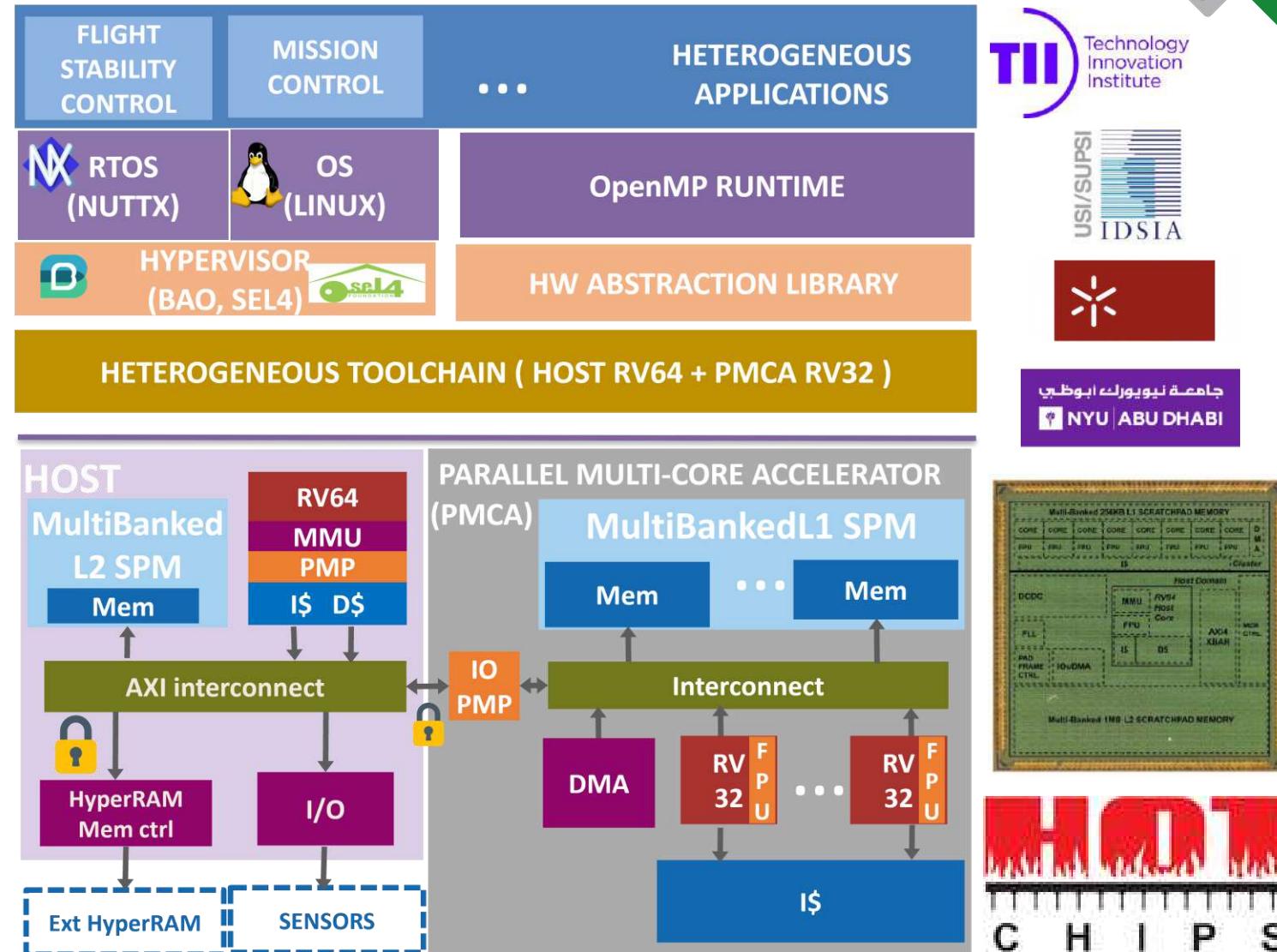
Single Source, Common API, Specialized Code	Separate data models (width of pointers, size_t, long, ...)	Separate ISAs	Separate Libraries	Single Binary, Nested Executables
--	--	----------------------	-------------------------------	--



OpenMP offloading requires a **host compiler plus one target compiler for each PMCA ISA** in the system.

Shaheen: Heterogeneous SoC for Nano-UAVs

- Secure Heterogeneous Application Processor
- Host Subsystem
 - CVA6
 - H-Extensions → BAO, Sel4
 - Timing channel attack protection
- PULP Cluster with Mixed-Precision Extension
- HyperRAM Memory Controller
 - Up to 512 Mbit
 - Up to 1.6 Gbit/s





PULP in HPC

Occamy – Massive Scaling



Dual Chiplet System Occamy:

- Technology: GF12LP+
- Area: 73mm²

Interposer Hedwig:

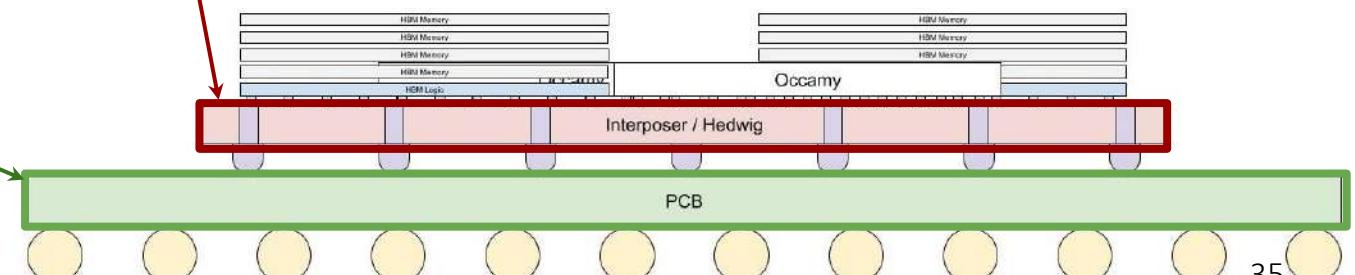
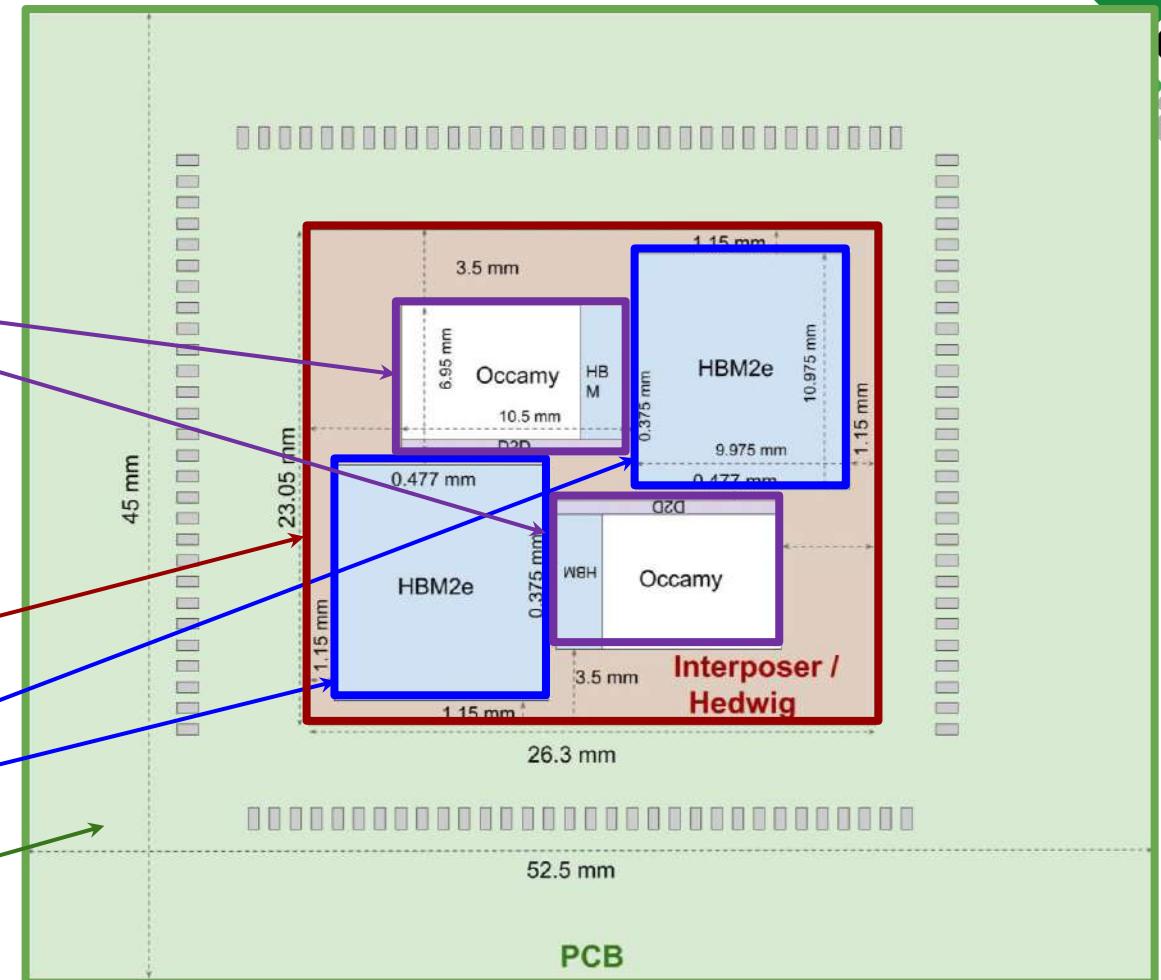
- Technology: 65nm, passive (only BEOL)
- Area: 26.3mm x 23.05mm

HBM2e:

- 16GB HBM2e (Micron)

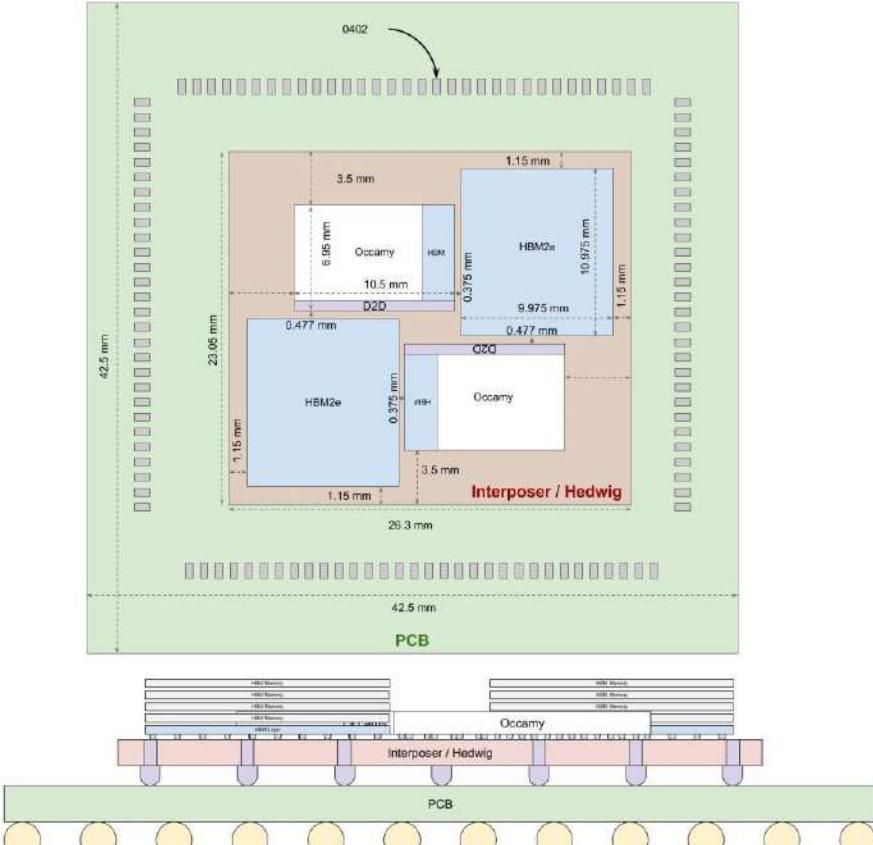
Fan-out PCB:

- RO4350B
- 52.5mm x 45mm

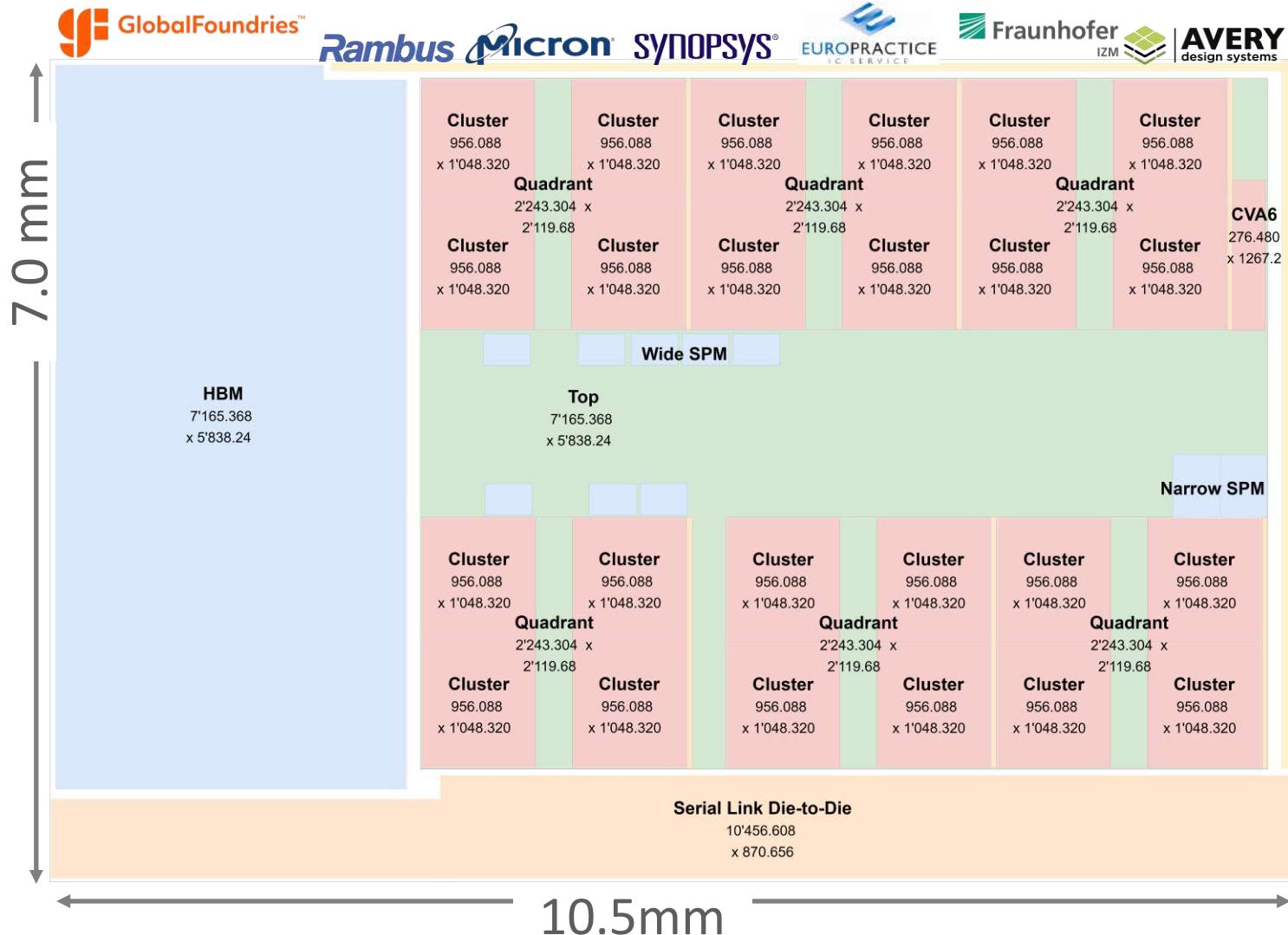




Occamy: Next-Generation HPC SoC



2.5D 2xOccamy+2HBM2 on interposer



RV-Based Snitch Cluster



8 Snitch compute cores

- Single-stage, small Integer control core

9th Core: DMA

- 512 bit data interface
- Efficient data movement

128 kB TCDM

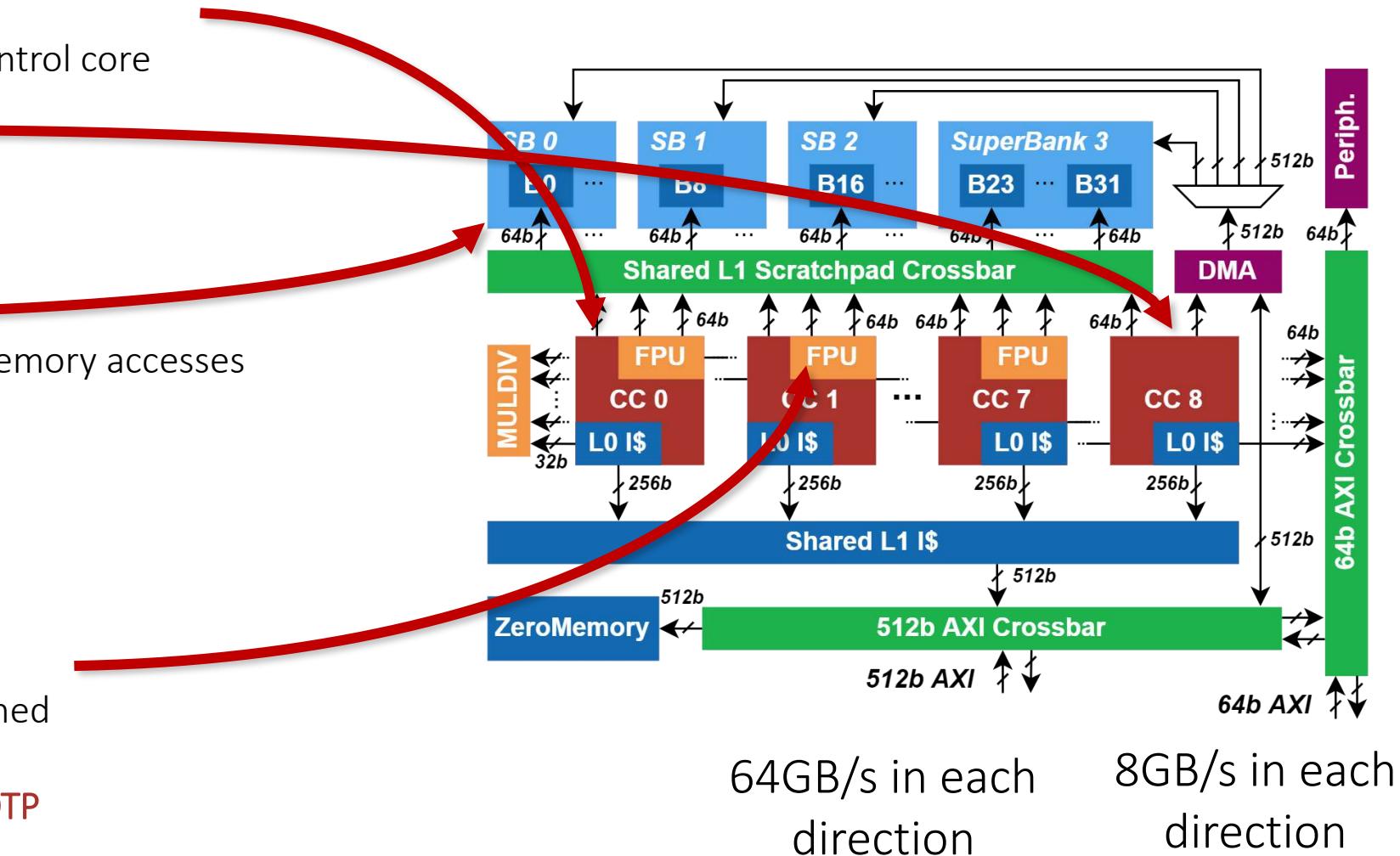
- Scratchpad for predictable memory accesses
- 32 Banks

Custom ISA extensions

- Xfrep, XSSR
- New: XSSR sparsity support

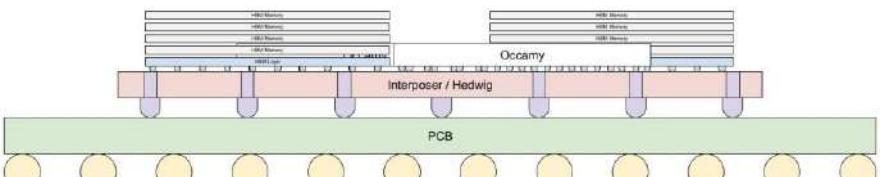
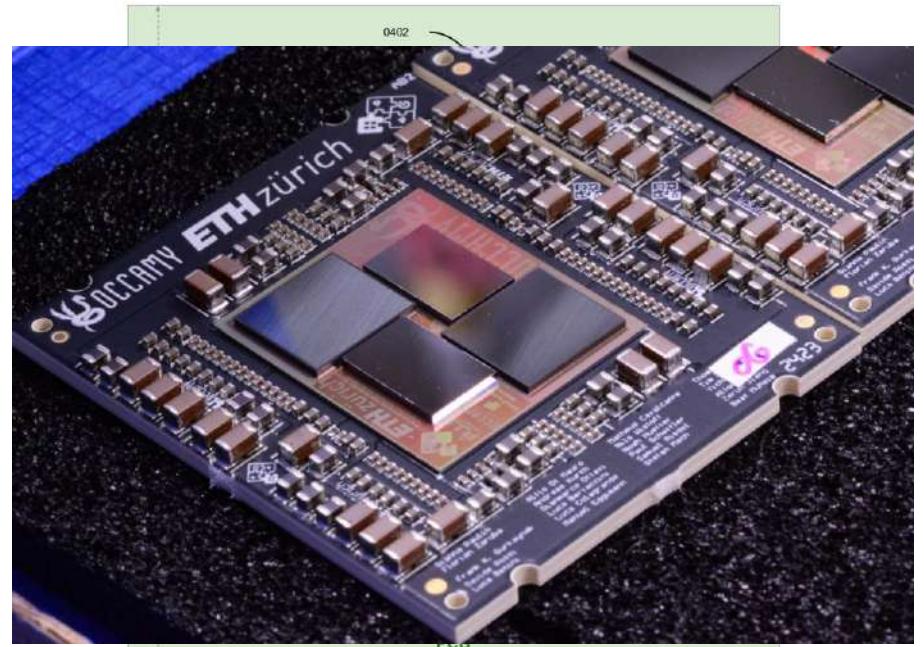
1 FPU per Snitch core

- Decoupled and heavily pipelined
- Multi-format FPU (+SIMD)
- New: Minifloat support + SDOTP





Occamy is On The Tester Now



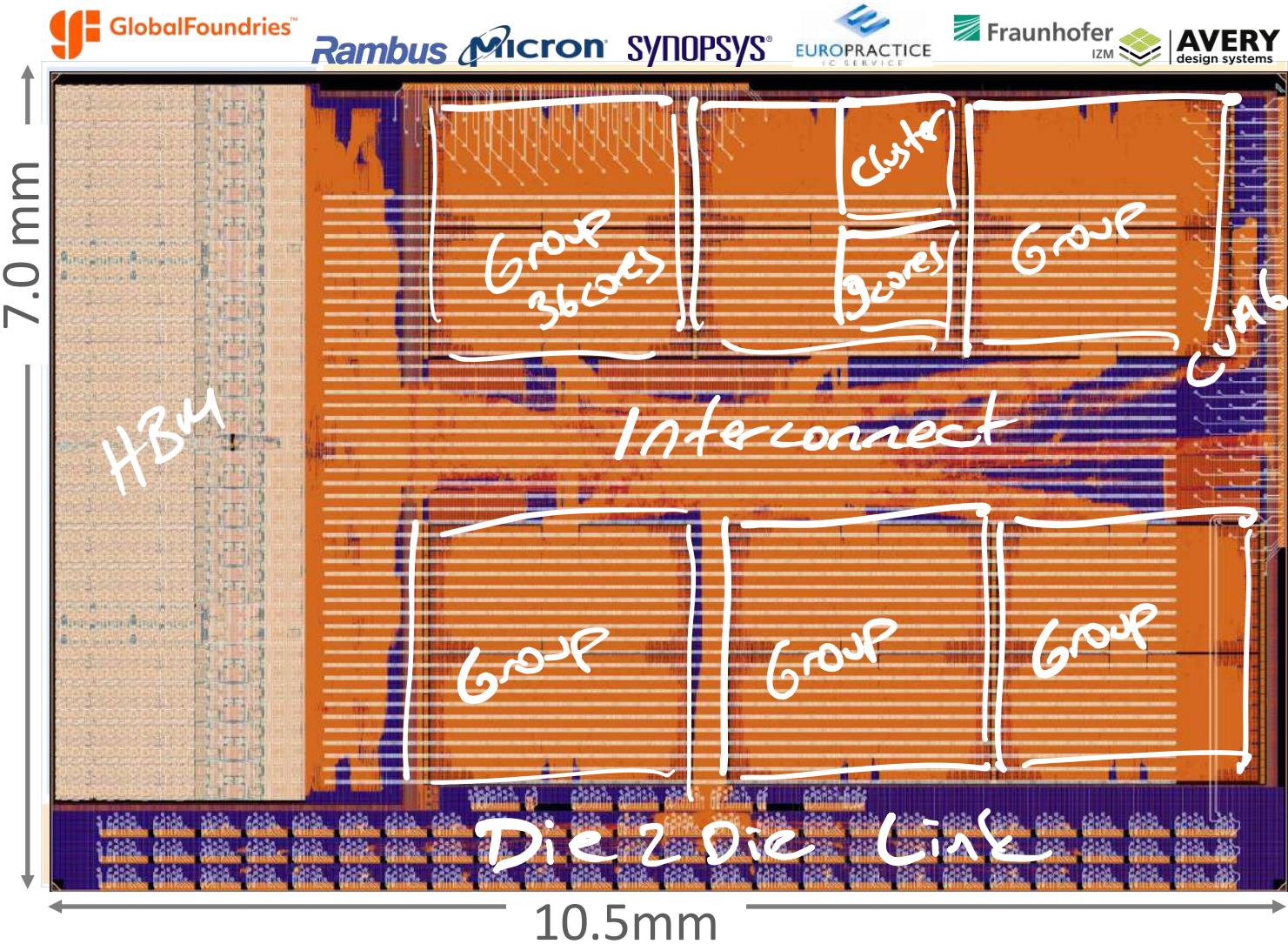
2.5D 2xOccamy+2HBM2 on interposer

GF12 @1GHz: FP64 384GFLOP/s, FP32,16,8

ETH zürich



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA





What's Next?

AI Saqr: Open-Research Platform for Secure UAVs



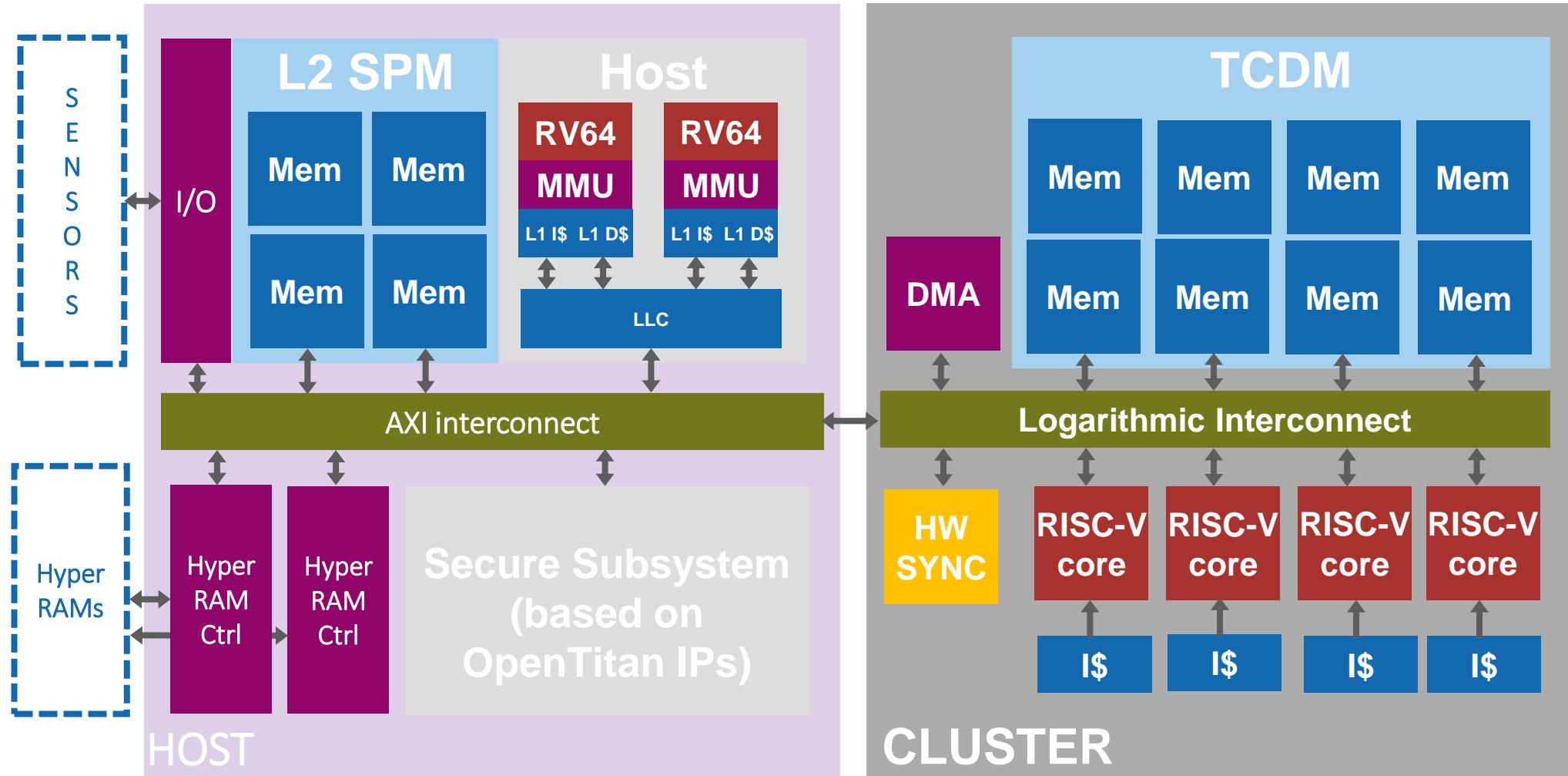
جامعة خليفة
Khalifa University

TII
Technology Innovation Institute

USI/SUPSI
IDSIA

Universidade do Minho

جامعة نيويورك أبوظبي
NYU ABU DHABI



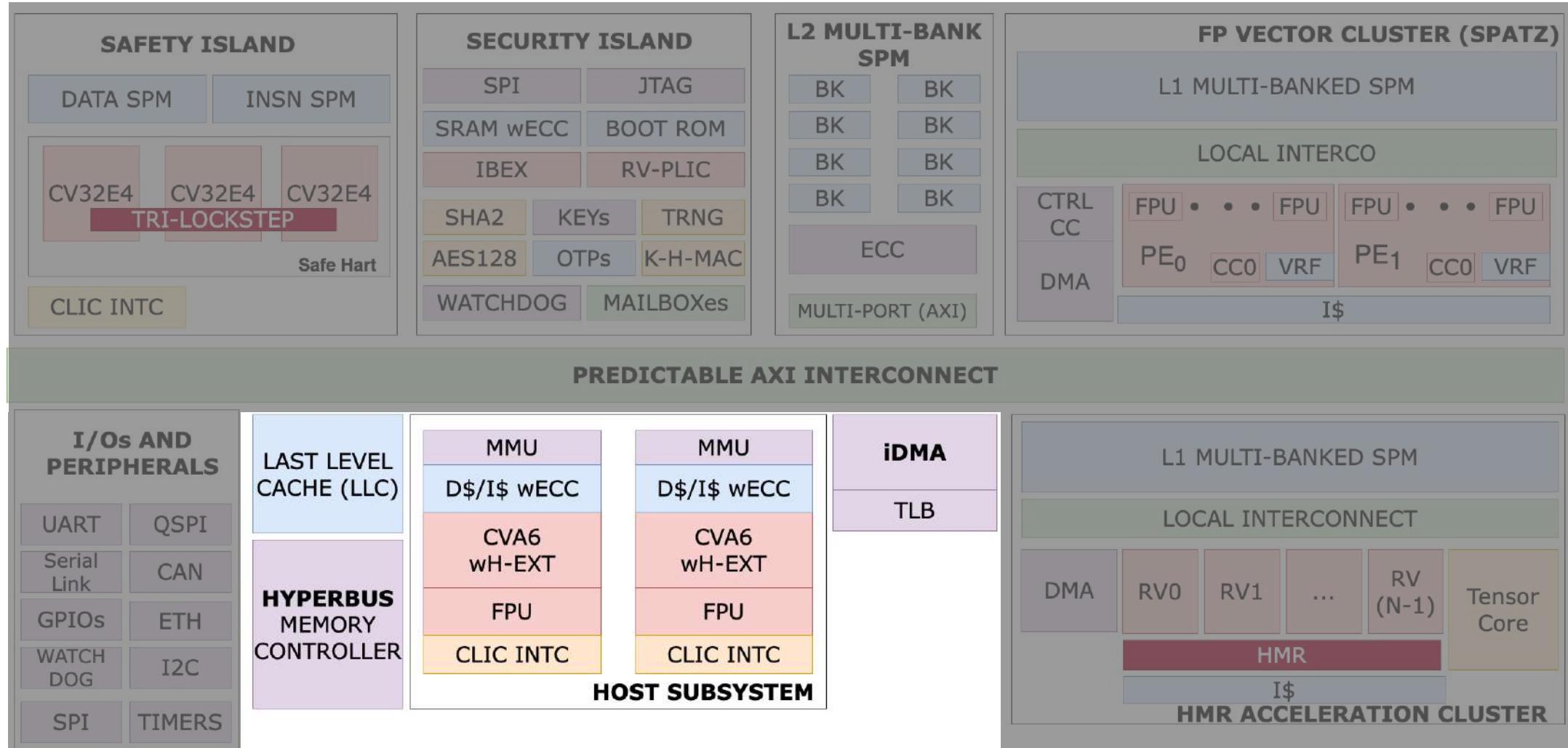
Carfield: Open-Research Platform for Automotive Systems



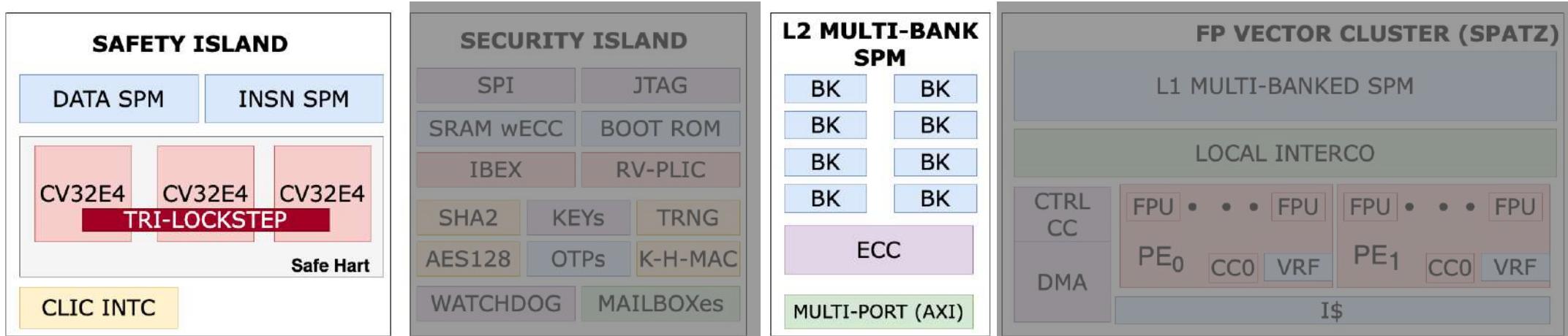
- Develop a **precompetitive RISC-V Automotive SoC**
 - Based on fully open-source HW and SW IPs
 - **Scalable** and **configurable** architectural template based on our PULP architectural ball-park
 - To adapt to different computing requirements
 - **Full SW stack** to address requirements of RISC-V based automotive applications
- Collaborative research roadmap for automotive-driven computing architectures
 - Functional safety, Hardware/Software Acceleration, Time-Predictability, Fast-Interrupts, HW-based Virtualization
- Benchmark RISC-V ecosystem and architectural solutions for automotive
 - Contributing to European interest to build an automotive reference platform around RISC-V
- Close the gap between RISC-V and ARM-powered solutions for automotive SoCs
- Fully aligned with EU Initiative (DG-connect) on RV for automotive (HW and SDV)!



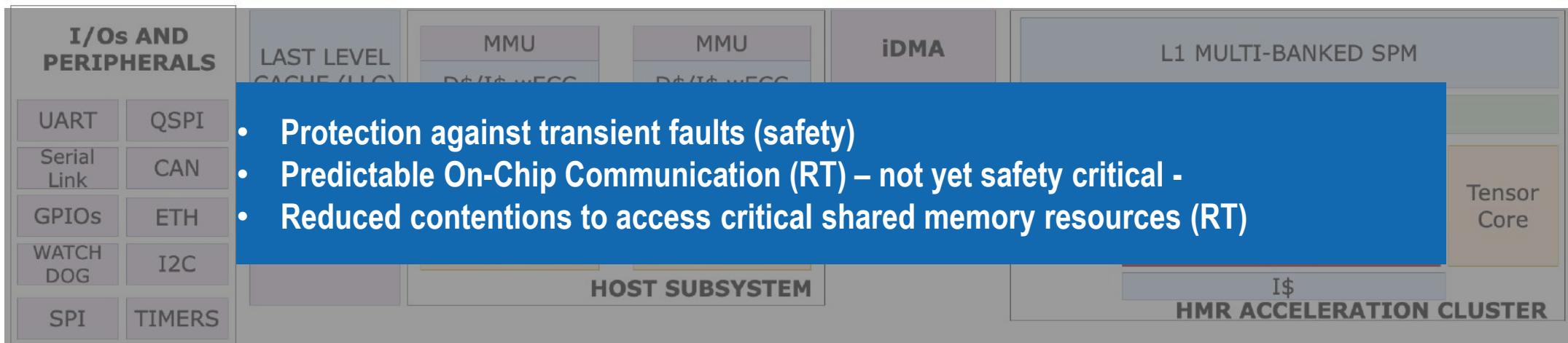
Host-Domain for Low-Criticality Linux-Based Applications



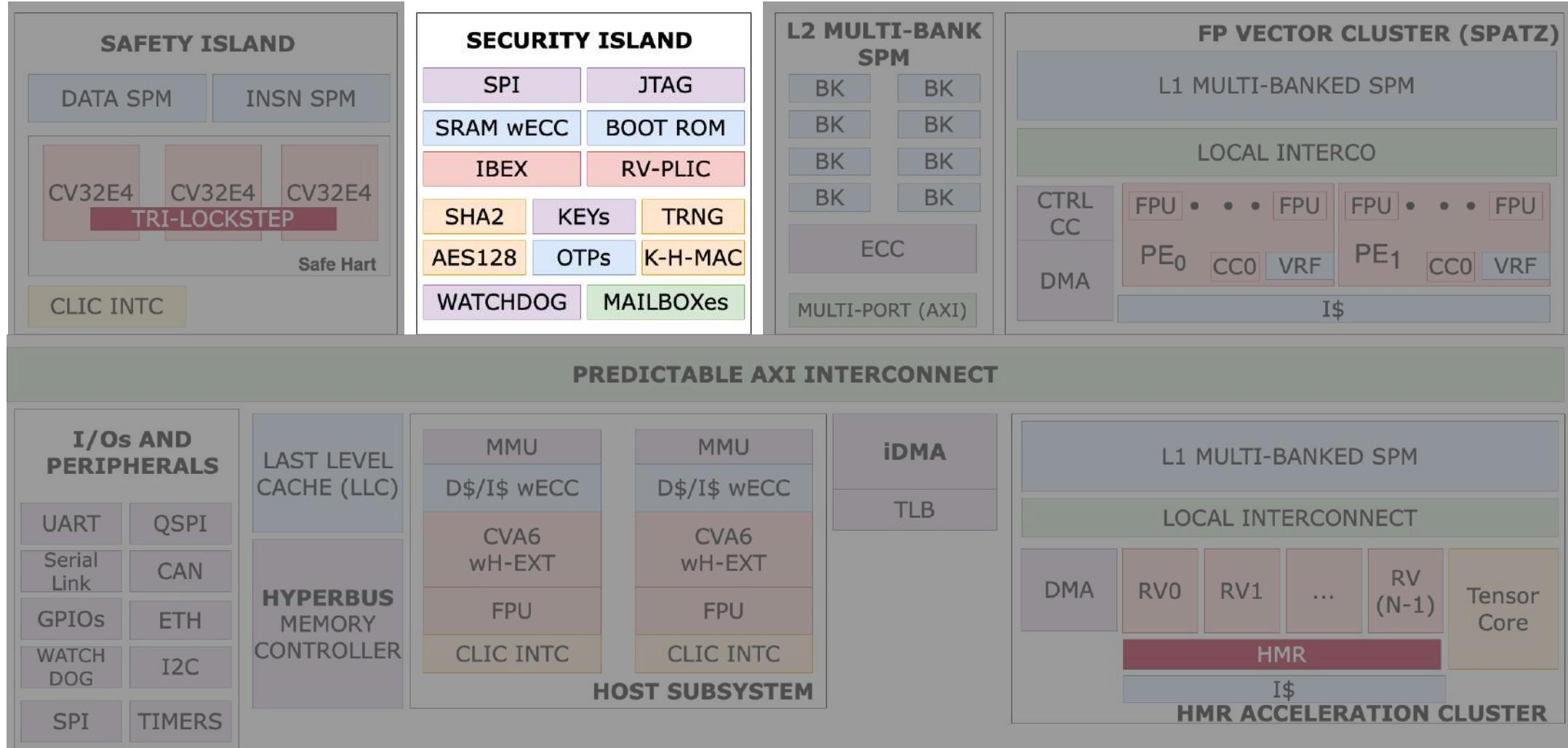
How Do We Handle Safety-Critical and Real-Time Tasks?



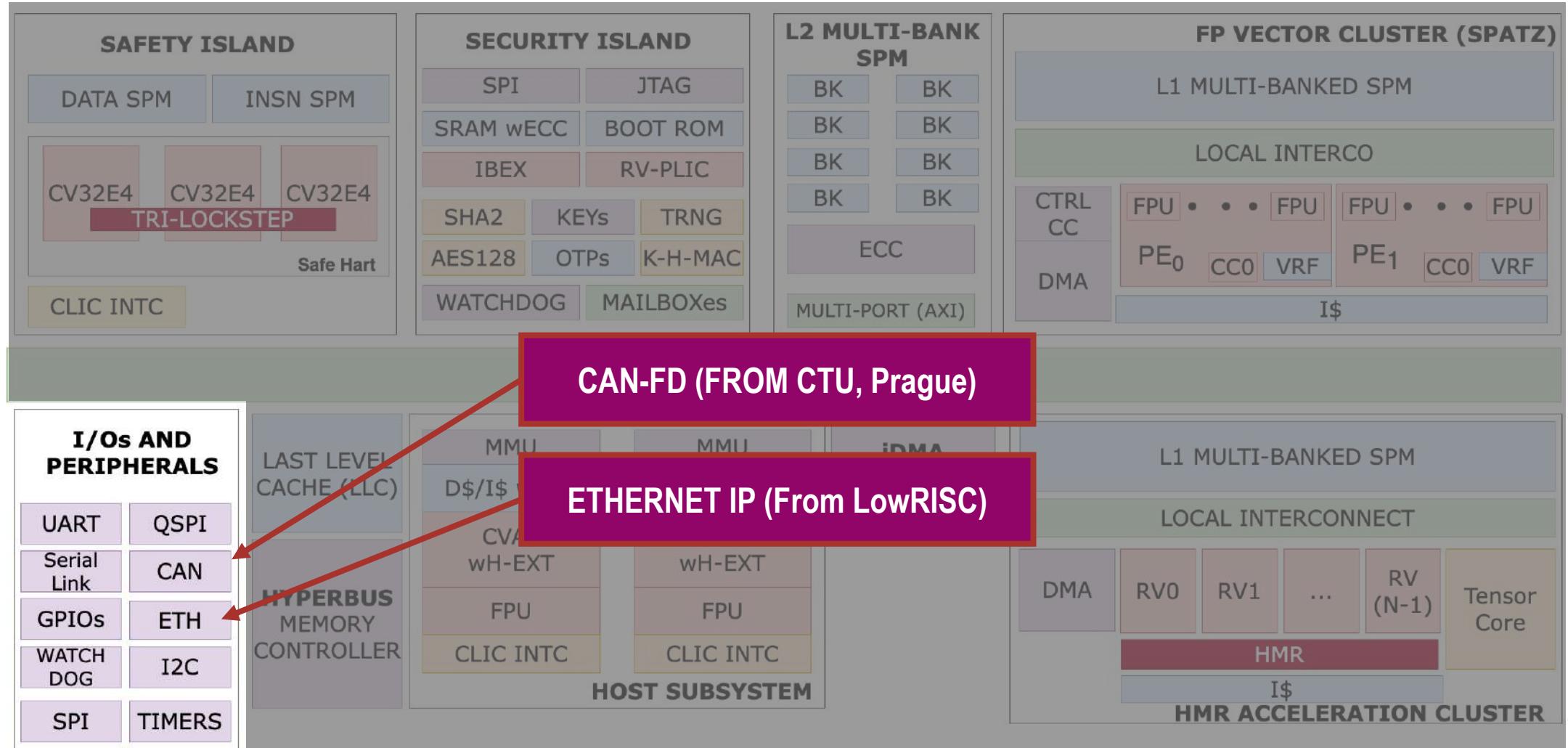
PREDICTABLE AXI INTERCONNECT



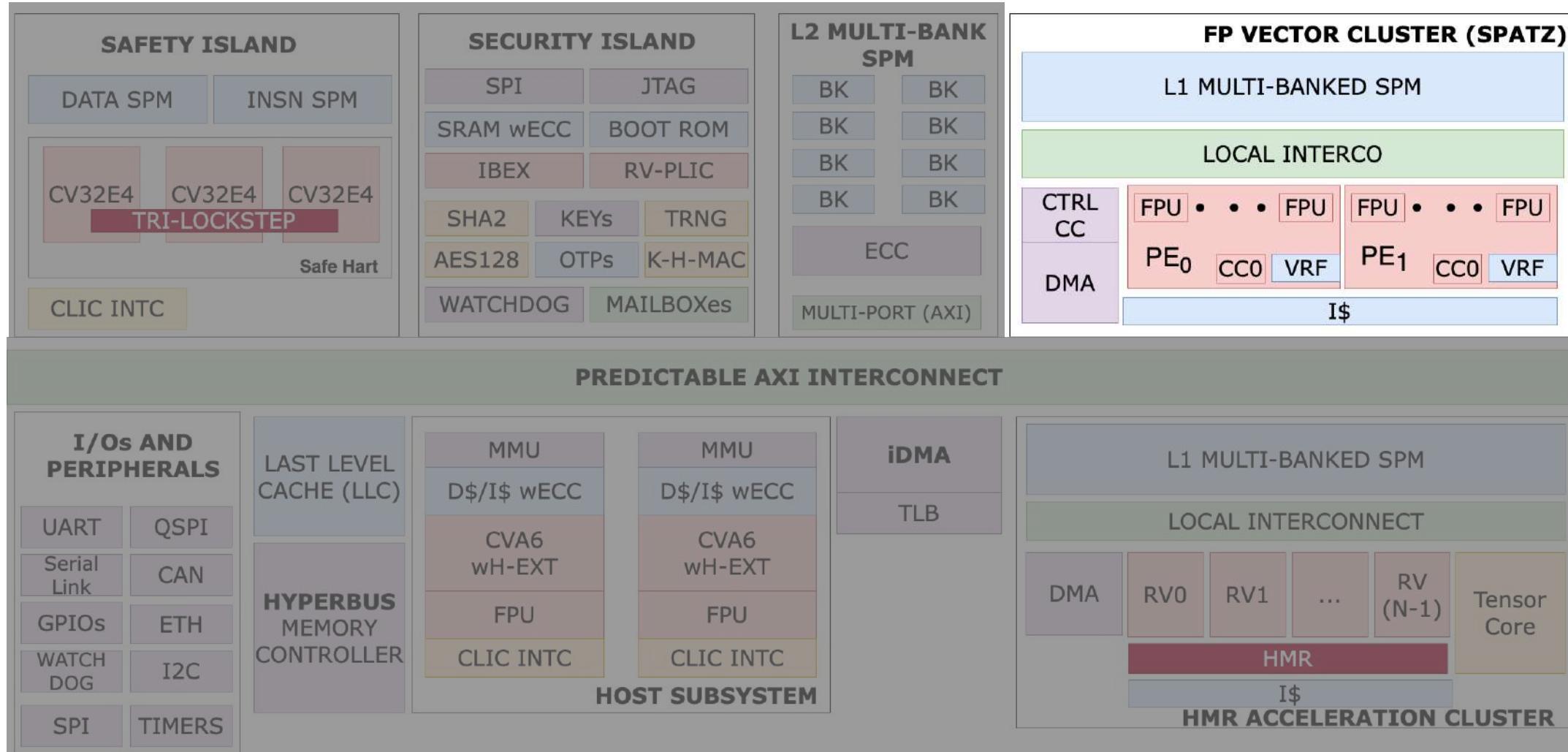
What About Security and Data Encryption/Decryption?



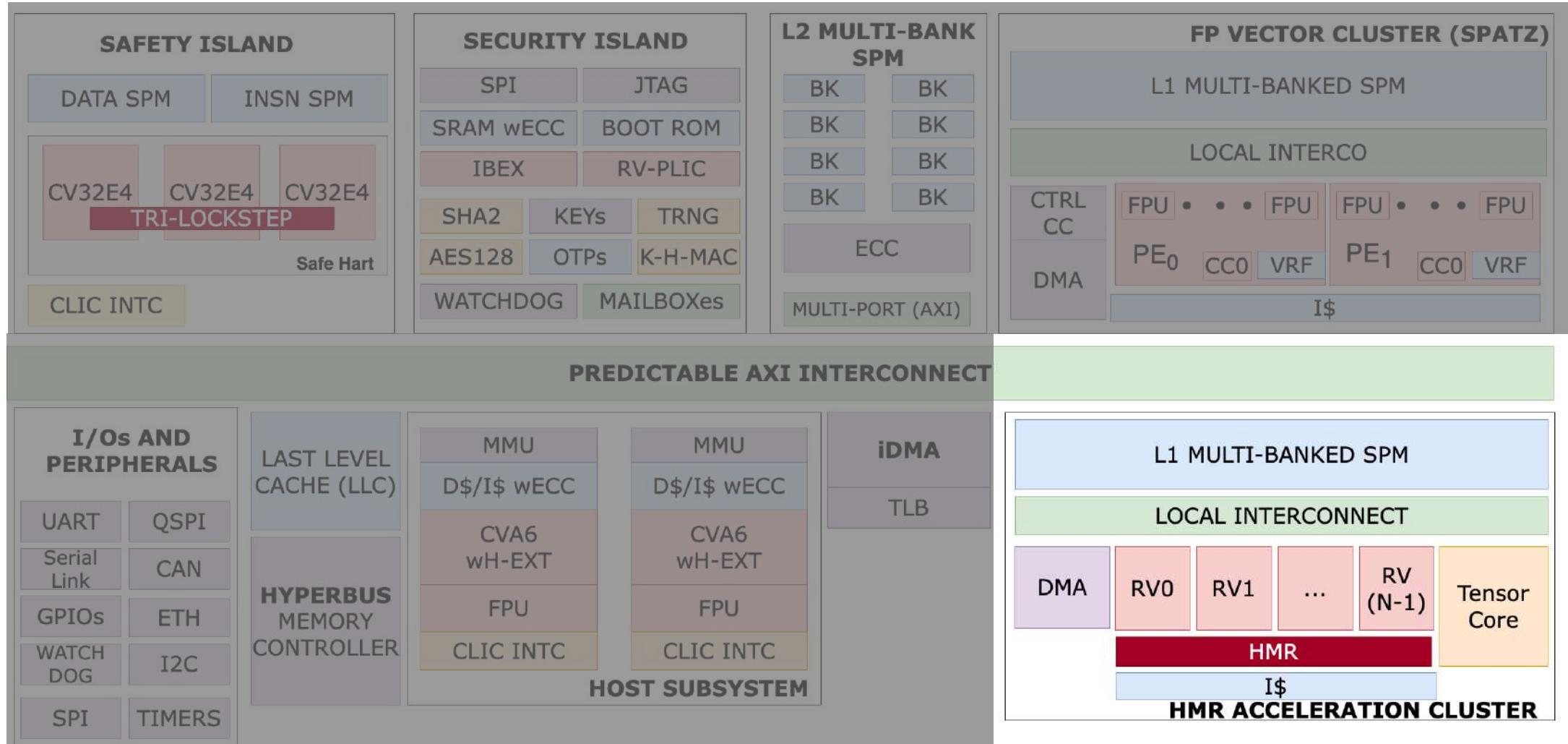
The I/O Communication



The Spatz Acceleration Cluster



The Inference Acceleration Cluster





Questions?



<http://pulp-platform.org>



@pulp_platform