

Sub-pJ-Operation Scalable Computing

The PULP Experience

Davide Rossi¹,

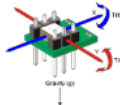
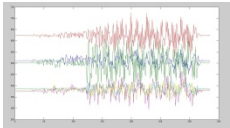
Antonio Pullini², Igor Loi¹, Francesco Conti¹, Andreas Traber², Michael Gautschi², Andreas Traber²,
Frank K. Gürkaynak², Davide Schiavone², Daniele Palossi², Alessandro Capotondi¹, Giuseppe Tagliavini¹,
Andrea Marongiu^{1,2}, Germain Haugou², Eric Flamand², Luca Benini^{1,2}

¹DEI-UNIBO, ²IIS-ETHZ



Sense

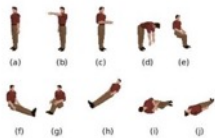
MEMS IMU



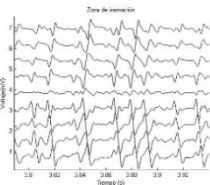
MEMS Microphone



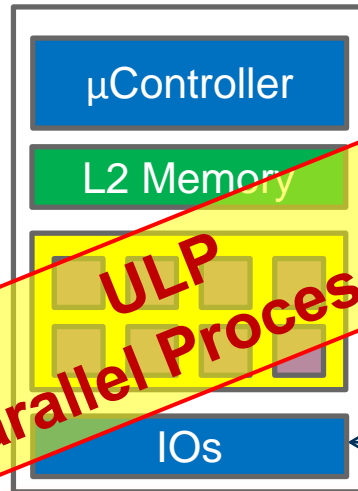
ULP Imager



EMG/ECG/EIT



Analyze and Classify



$1 \div 25 \text{ MOPS}$
 $1 \div 10 \text{ mW}$

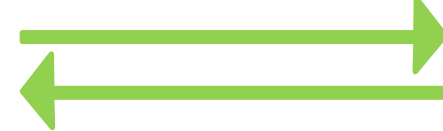


Transmit

Short range, medium BW



Low rate (periodic) data



SW update, commands

Long range, low BW

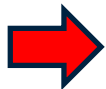


$100 \mu\text{W} \div 2 \text{ mW}$

*Battery + Harvesting powered
→ a few mW power envelope*

Idle: $\sim 1 \mu\text{W}$
Active: $\sim 50 \text{ mW}$

	INPUT BANDWIDTH	COMPUTATIONAL DEMAND	OUTPUT BANDWIDTH	COMPRESSION FACTOR
■ Image				
Tracking: [*Lagroce2014]	80 Kbps	1.34 GOPS	0.16 Kbps	500x
■ Voice/Sound				
Speech: [*VoiceControl]	256 Kbps	100 MOPS	0.02 Kbps	12800x
■ Inertial				
Kalman: [*Nilsson2014]	2.4 Kbps	7.7 MOPS	0.02 Kbps	120x
■ Biometrics				
SVM: [*Benatti2014]	16 Kbps	150 MOPS	0.08 Kbps	200x



Extremely compact output (single index, alarm, signature)



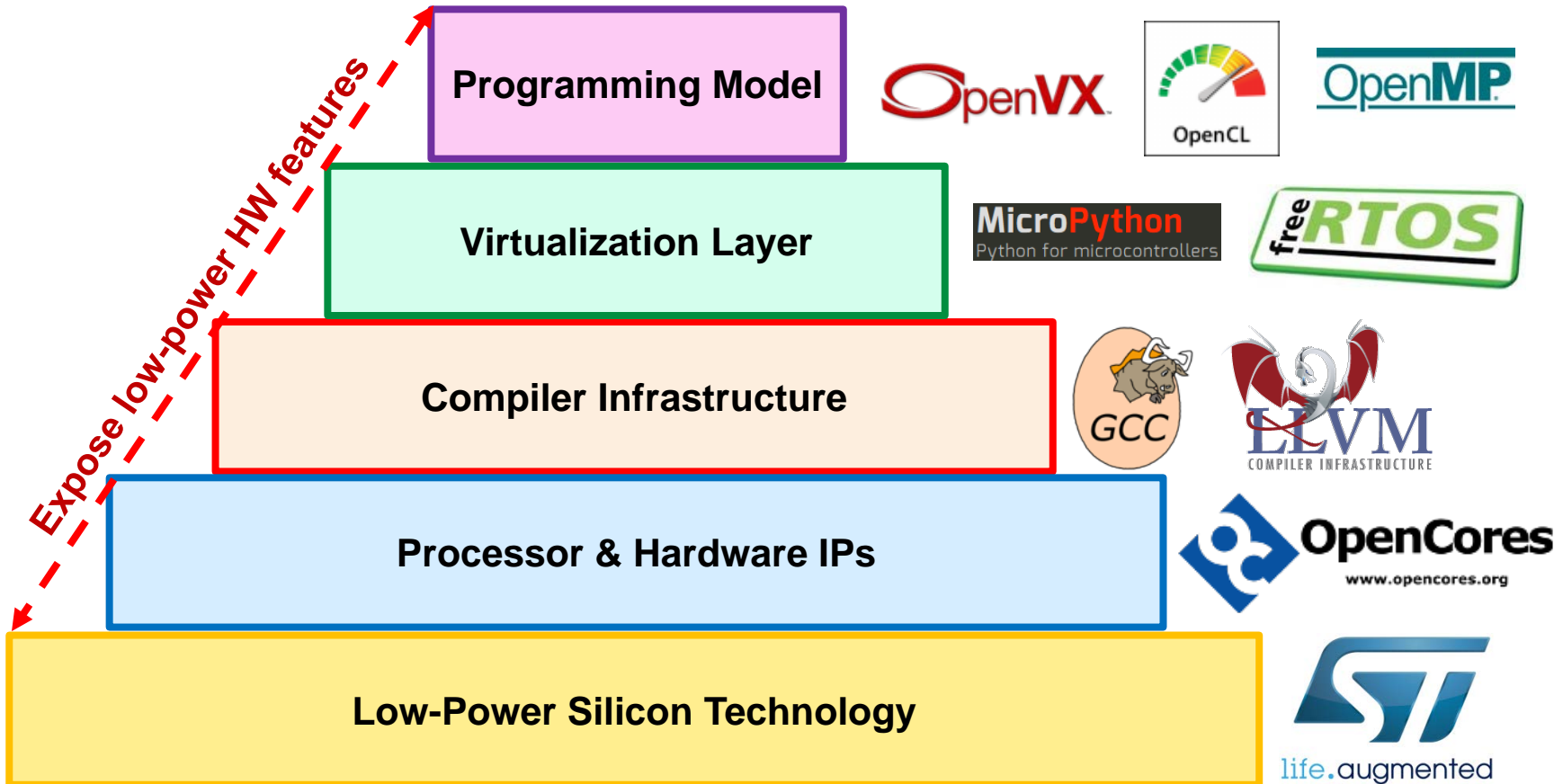
Computational power of ULP μ Controllers is not enough



Parallel workloads

PULP: pJ/op Parallel ULP computing

pJ/op is traditionally the target of ASIC + μ Controllers

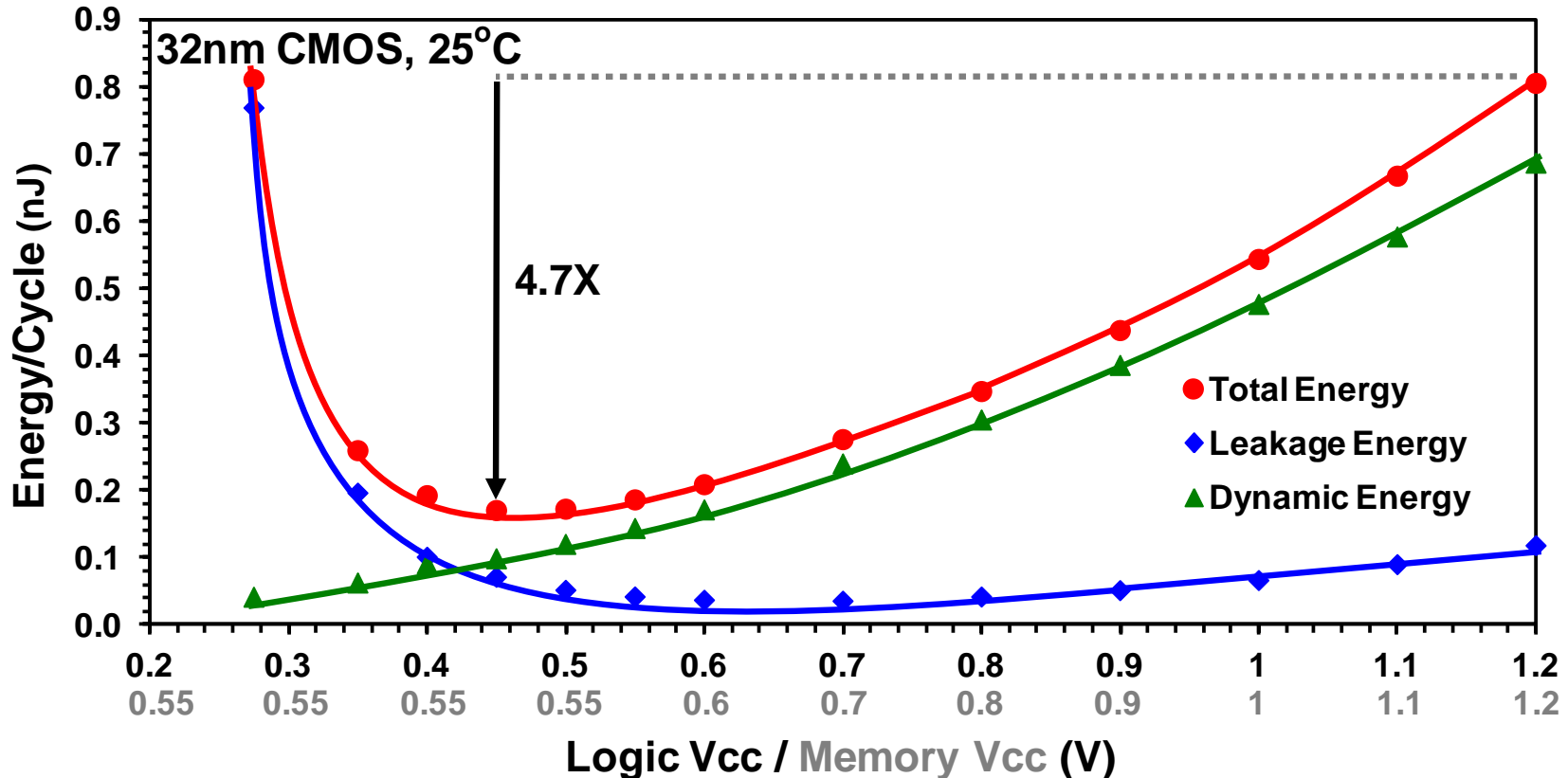


Parallel + programmable + heterogeneous ULP computing
1mW-10mW active power

Near-Threshold Multiprocessing



Source: Vivek De, INTEL – Date 2013



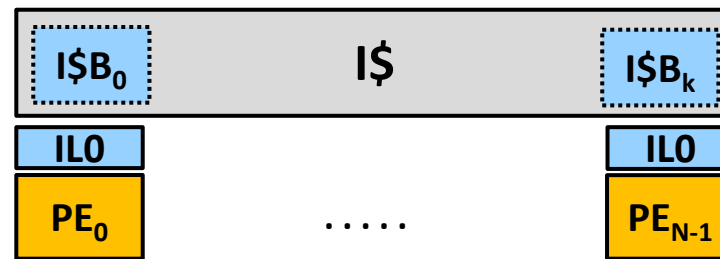
Near-Threshold Computing (NTC):

1. Don't waste energy pushing devices in strong inversion
2. Recover performance with parallel execution



Ultra-low power scalable computing

Shared L1 I\$ with Multi-instruction load



Private Loop/Prefetch Buffer

N Cores

4-stage, in-order RiscV
and OpenRISC ISAs

Micro-MMU (demux)

Periph
+ExtM



Tightly Coupled DMA

Shared L1 DataMem + Atomic Variables

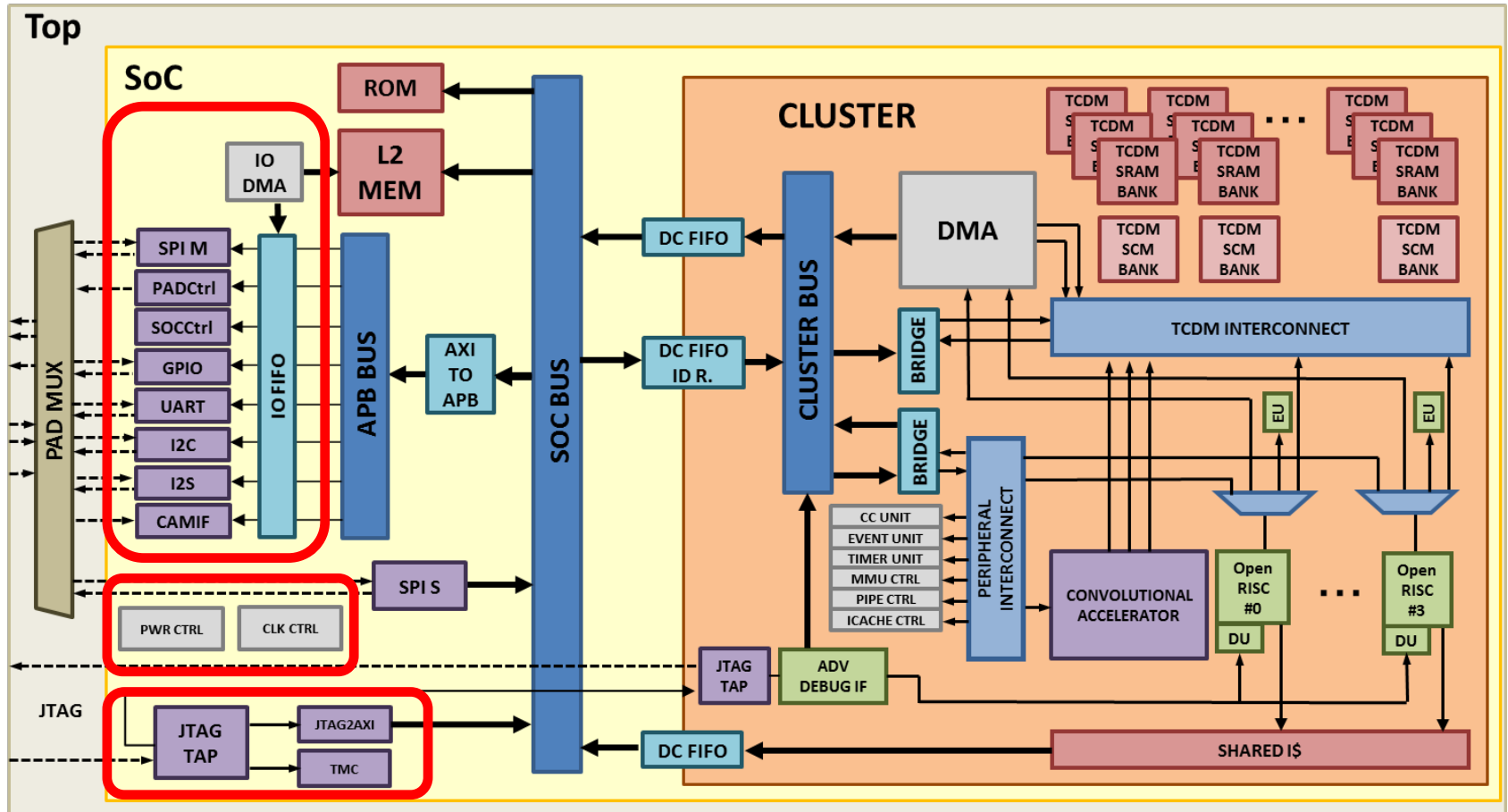
1.. 8 PE-per-cluster, 1...32 clusters

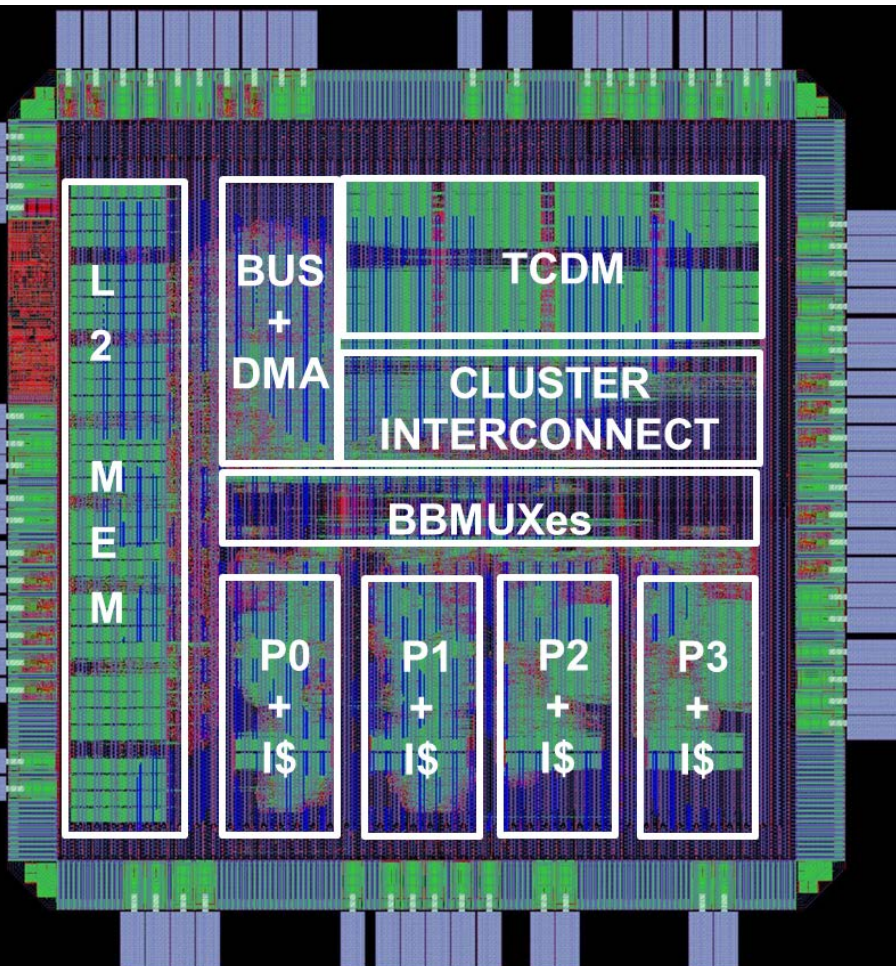
NT but parallel → Max. Energy efficiency when Active

+ strong PM for (partial) idleness



64-bit ——— 32-bit ———



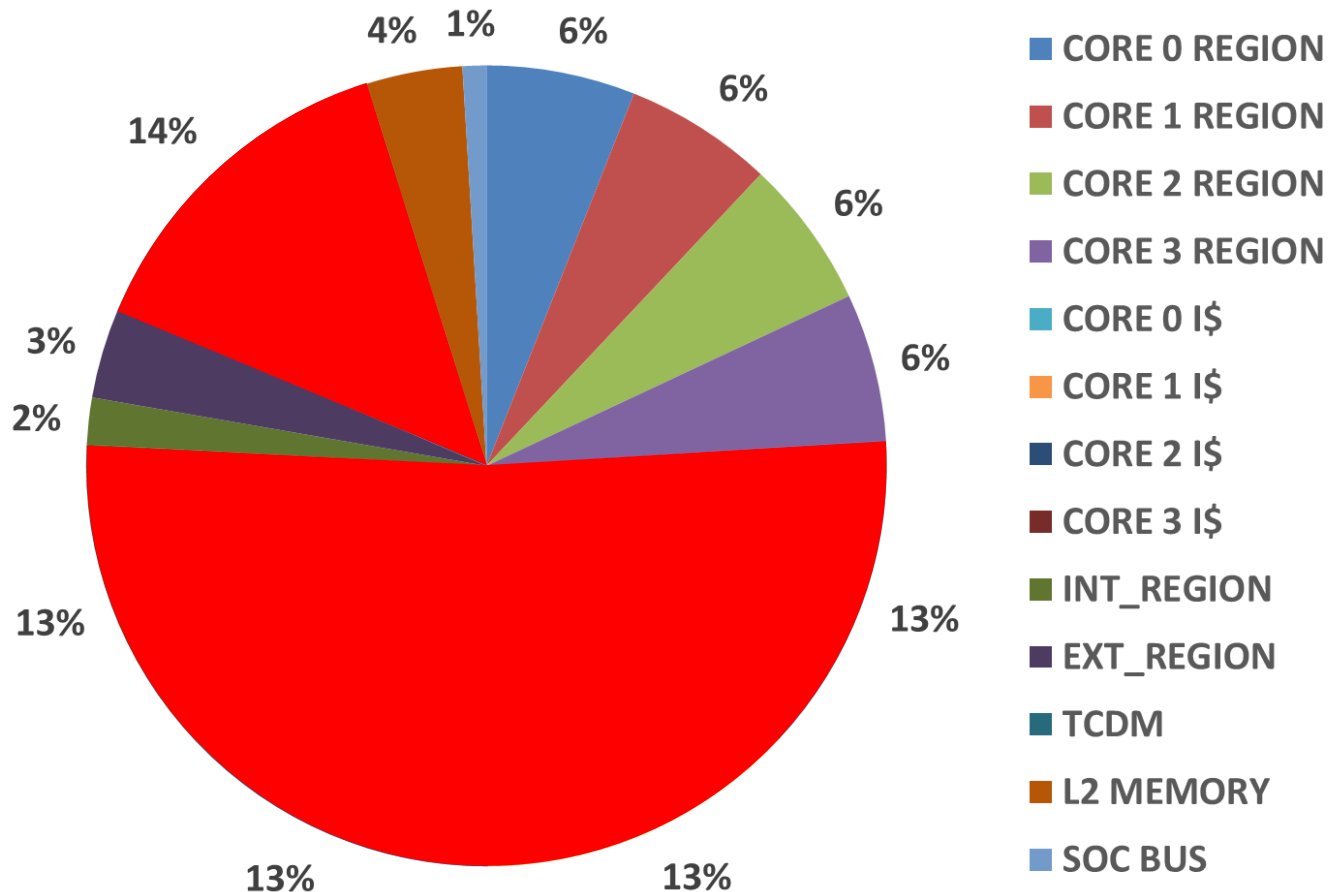


CHIP FEATURES

Technology	28nm FDSOI (RVT)
Chip Area	3mm ²
# Cores	4xOpenRISC
I\$	4x1kbyte (private)
TCDM	16 kbyte
L2	16 kbyte
BB regions	6
VDD range	0.45-1.2V
VBB range	-1.8V - +0.9V
Perf. Range	1 MOPS-1.9GOPS
Power Range	100 μ W - 127 mW
Peak Efficiency	60 GOPS/W@0.5V*

ISSCC15 (student presentations), Hot Chips 15, ISSCC16 (paper+student presentation)

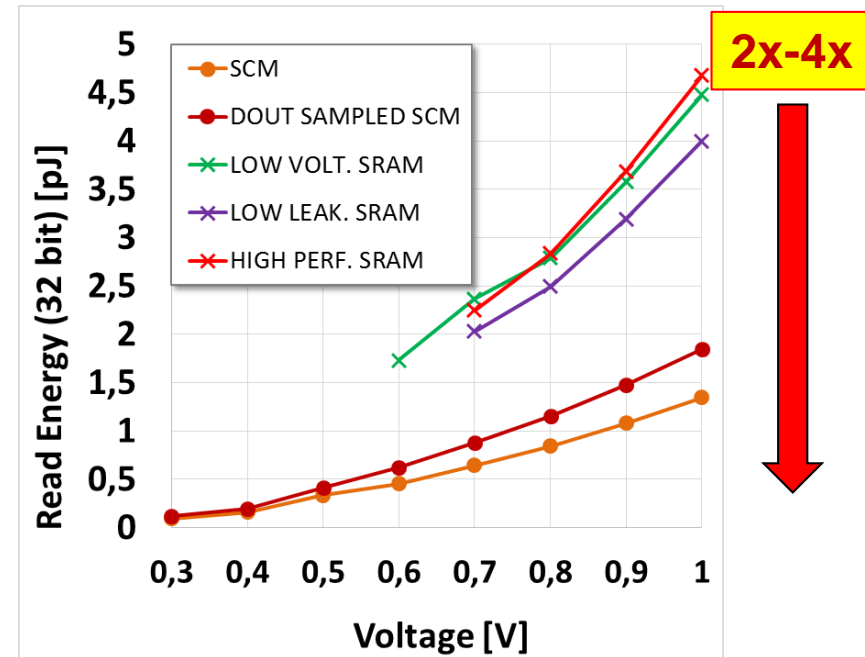
PULPv1 POWER BREAKDOWN @ BEST ENERGY POINT

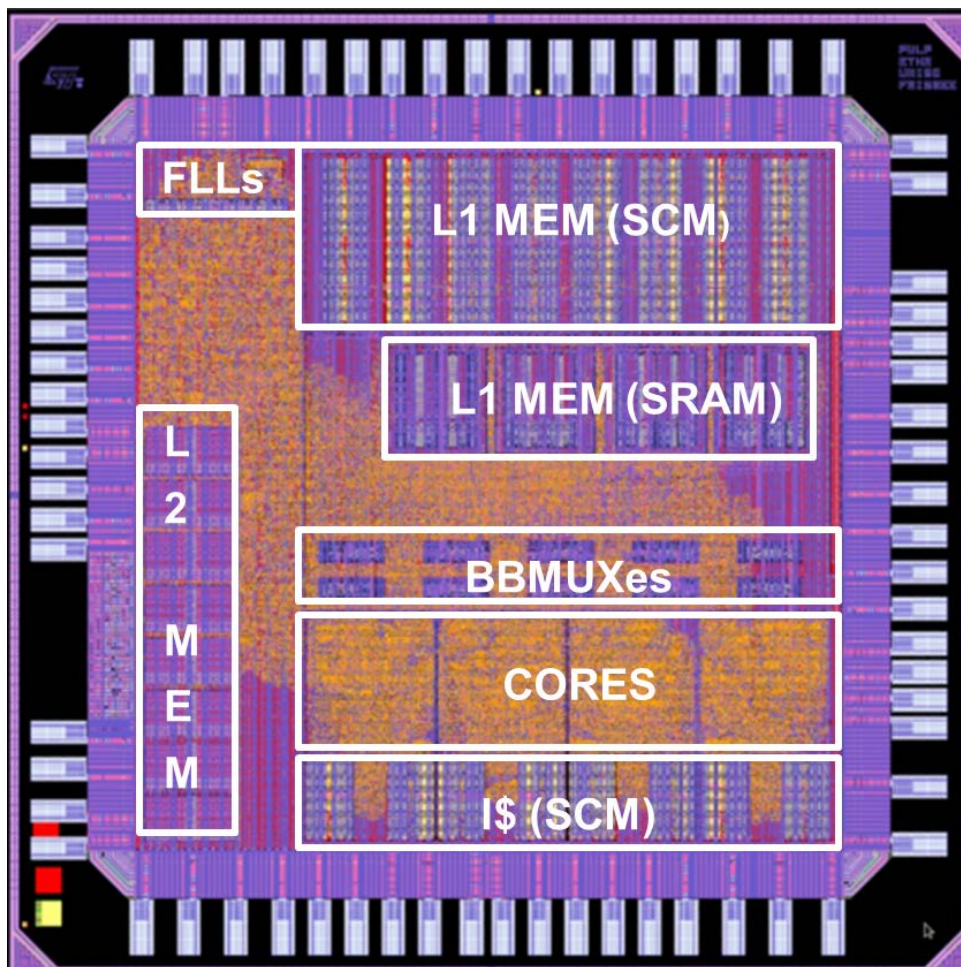


I\$s + TCDM consume > 60% of total power!!!

- “Standard” 6T SRAMs:
 - High VDDMIN
 - Bottleneck for energy efficiency
- Near-Threshold SRAMs (8T)
 - Lower VDDMIN
 - Area/timing overhead (25%-50%)
 - High active energy
 - Low technology portability
- Standard Cell Memories:
 - Wide supply voltage range
 - Lower read/write energy (2x - 4x)
 - Easy technology portability
 - Major area overhead (2x)

256x32 6T SRAMS vs. SCM





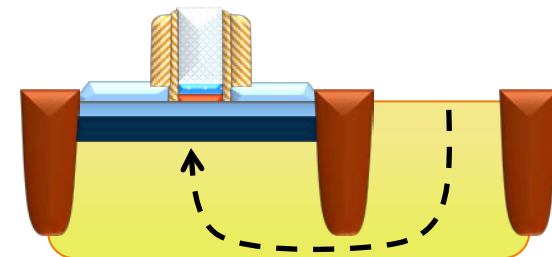
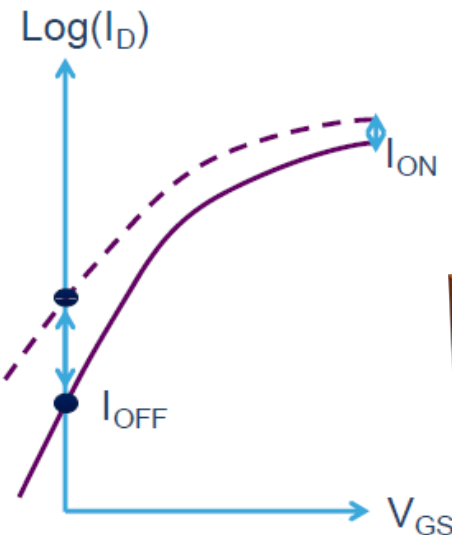
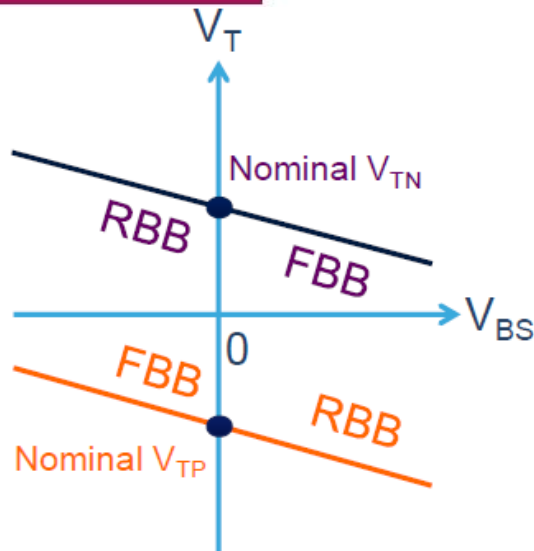
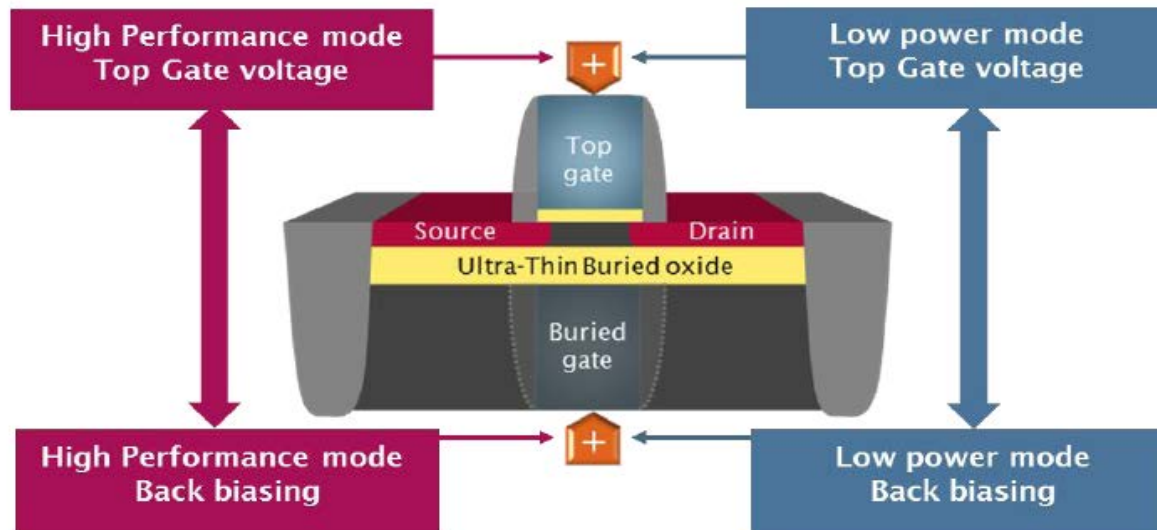
Taped out Nov. 2014!

CHIP FEATURES

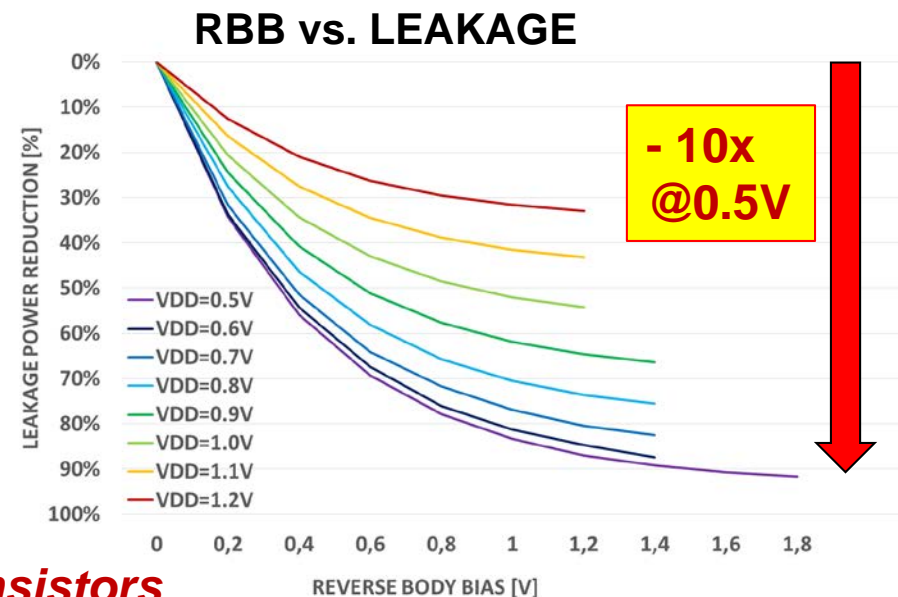
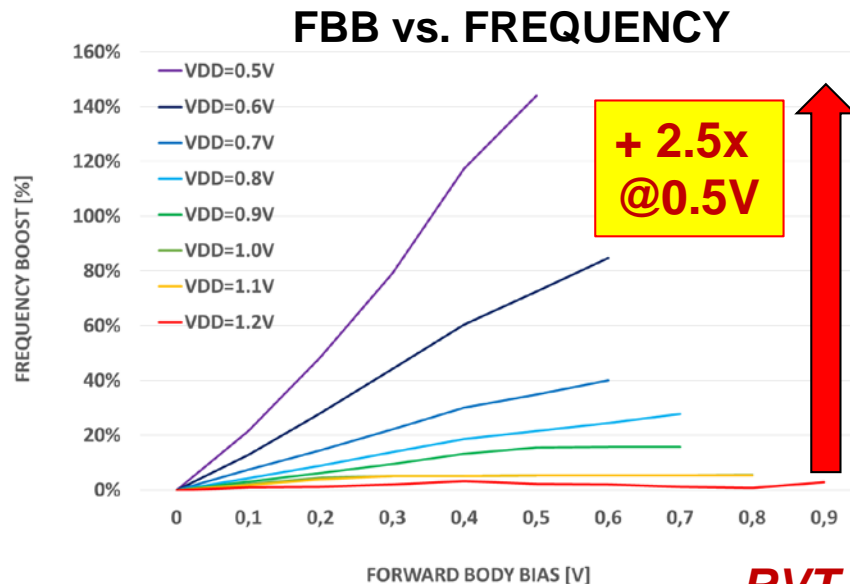
Technology	28nm FDSOI (LVT)
Chip Area	3mm ²
# Cores	4xOpenRISC
I\$ (SCM)	4x1kbyte (private)
TCDM	32 + 8 Kbyte
L2	64 kbyte
BB regions	10
VDD range	0.3-1.2V (0.45-1.2V)*
VBB range	0V-2V
Perf. Range	1 MOPS - 3.3 GOPS
Power Range	10μW - 480 mW
Peak Efficiency	193 GOPS/W@0.46V

	[2]	[3]	[4]	[5]	This Work
Technology	CMOS 32nm	CMOS 28nm LP	FD-SOI 28nm flip-well	FD-SOI 28nm conventional-well	FD-SOI 28nm flip-well
Data format	2x 32-bit superscalar	4x 32-bit VLIW	32-bit	32-bit	32-bit
# of cores	1	1	1	4	4
I\$/D\$/L2	8K/8K/n.a.	16K/32K/256K	4K/4K/n.a.	1Kx4/16K/16K	1Kx4/48K/64K
Voltage range (SRAMs)	0.28V – 1.0V (0.5V – 1.0V)	0.6V - 1.05V	0.4V - 1.3V	0.44V – 1.2V (0.54V – 1.2V)	0.32V – 1.15V (0.45V – 1.15V)
Max frequency	915 MHz	1.2 GHz	2.6 GHz	475 MHz	825 MHz
Best power density	170 μ W/MHz	58 μ W/MHz	62 μ W/MHz	65 μ W/MHz	20.7 μW/MHz
Best performance	1.8 GOPS	3 GOPS	2.6 GOPS	1.8 GOPS	3.3 GOPS
Peak energy efficiency (MAX)	11.7 MOPS/mW @ 50 MOPS	43.1 MOPS/mW @ 230 MOPS	16.1 MOPS/mW @ 460 MOPS	60 MOPS/mW @ 25.6 MOPS	193 MOPS/mW @ 162 MOPS





Body bias: Highly effective knob for power & variability management!

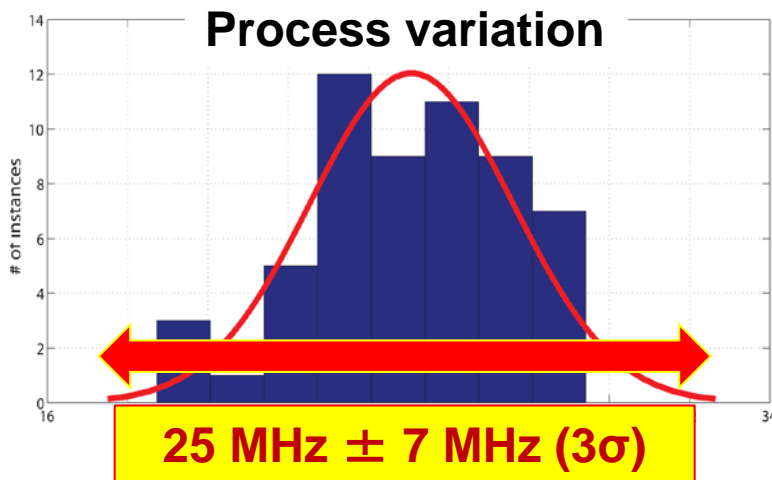


RVT transistors

- ➡ State retentive (no state retentive registers and memories)
- ➡ Ultra-fast transitions (tens of ns depending on n-well area to bias)
- ➡ Low area overhead for isolation (3μm spacing for deep n-well isolation)
- ➡ Thin grids for voltage distribution (small transient current for wells polarization)
- ➡ Simple circuits for on-chip VBB generation (e.g. charge pump)

But even with aggressive RBB leakage is not zero!

Body Biasing for Variability Management



120° C

100x
@0.5V

-40° C

Thermal inversion



Frequency Gain (%)

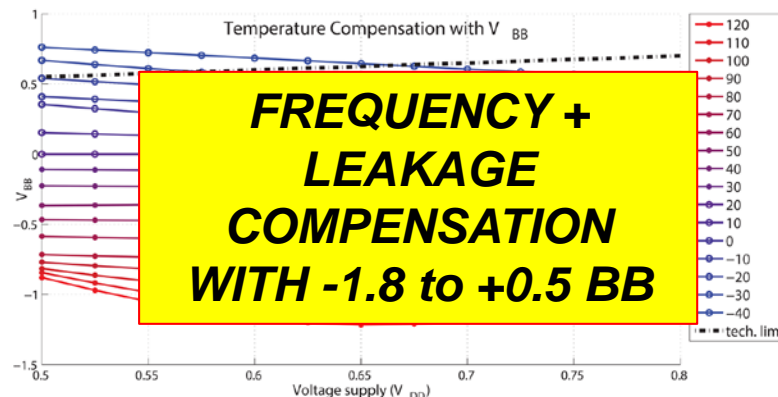
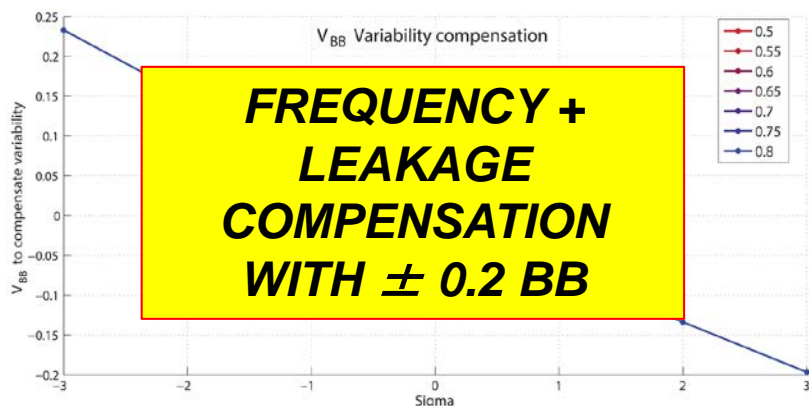
800 700 600 500 400 300 200 100 0 -100

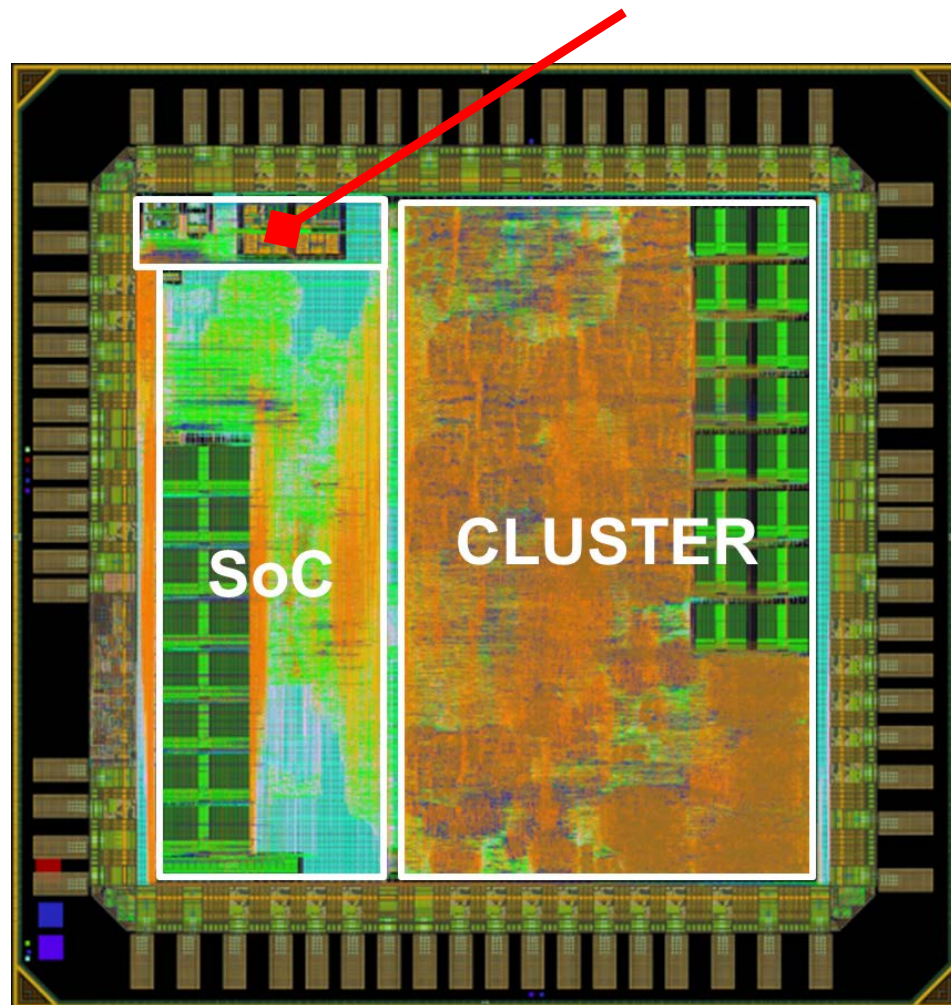
0.5 0.55 0.6 0.65 0.7 0.75 0.8

Voltage supply (V_{DD})

120 110 100 90 80 70 60 50 40 30 20 10 0 -10 -20 -30 -40

RVT transistors
FBB/RBB

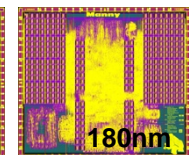
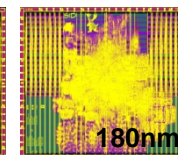
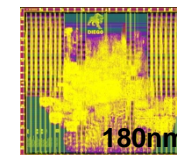
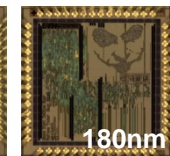
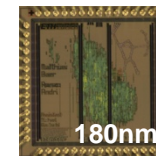
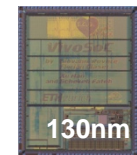
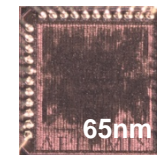
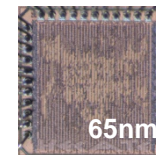
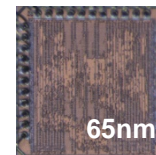
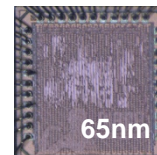
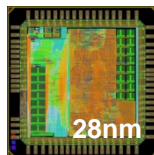
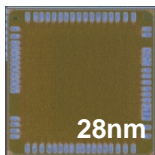
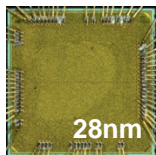


FLLs + BBGEN + PMBs + Multiprobes***under test NOW!*****CHIP FEATURES**

Technology	28nm FDSOI (RVT)
Chip Area	3mm ²
# Cores	4xOR1ON
I\$	4 kbyte (shared)
TCDM	64 + 8 kbyte
L2	128 kbyte
BB regions	2 (SoC, Cluster)
VDD range	0.4-0.7V
VBB range	-1.8V - +0.9V
Perf. Range	1 MOPS - 1.8 GOPS*
Power Range	20μW - 5.6 mW*
Peak Efficiency	387 GOPS/W@0.5V*

****Cluster, Estimated***

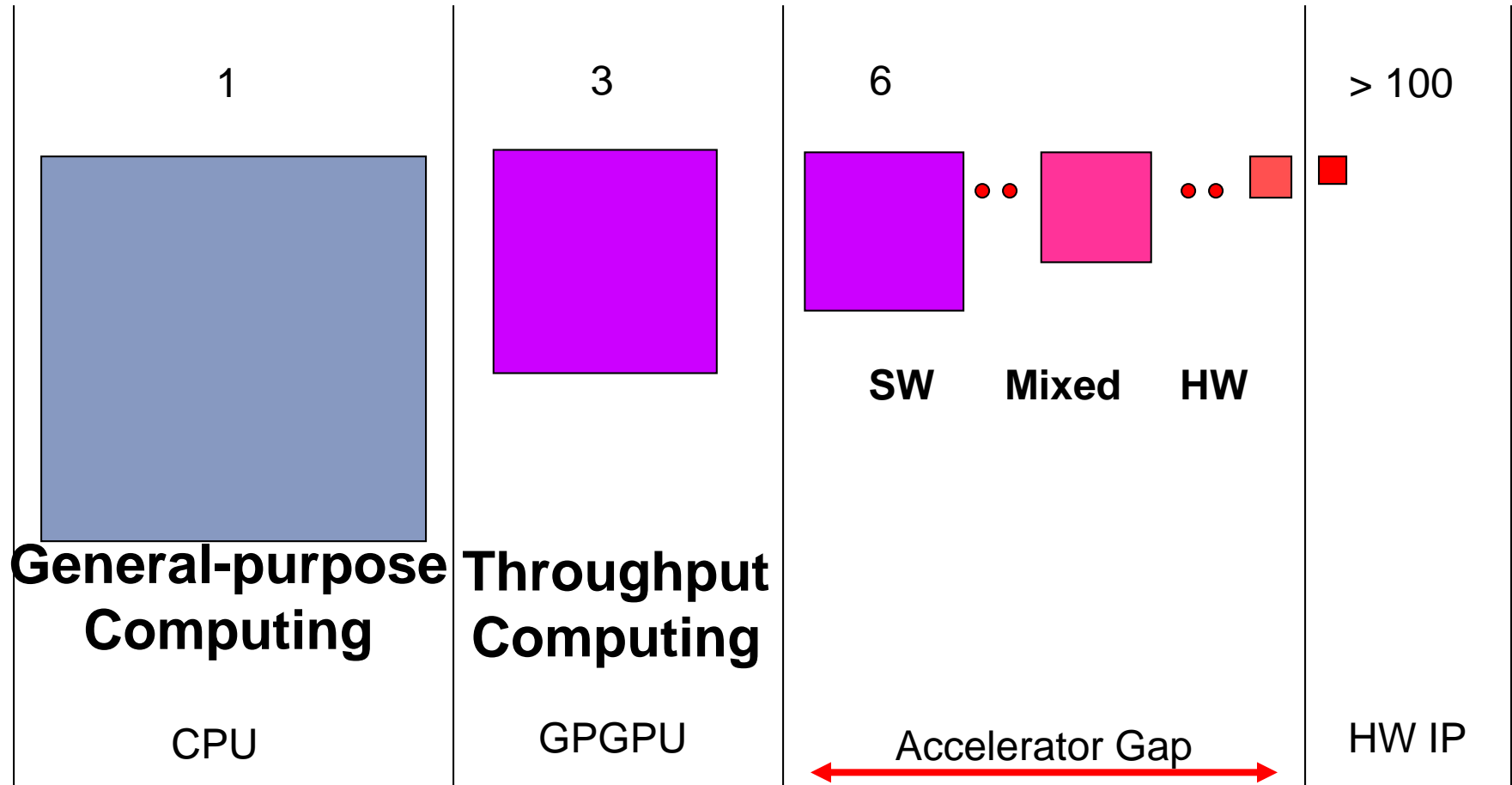
- **Main PULP chips** (ST 28nm FDSOI)
 - PULPv1
 - PULPv2
 - PULPv3 (under test)
 - PULPv4 (in progress)
- **PULP development** (UMC 65nm)
 - Artemis - IEEE 754 FPU
 - Hecate - Shared FPU
 - Selene - Logarithmic Number System FPU
 - Diana - Approximate FPU
 - Mia Wallace – full system
 - *Imperio* - PULPino chip
 - *Fulmine* – Secure cluster (Jan 2016)
- **RISC-V based systems** (GF 28nm)
 - Honey Bunny (Nov 2015)
- **Early building blocks** (UMC180)
 - Sir10us
 - Or10n
- **Mixed-signal systems** (SMIC 130nm)
 - VivoSoC
 - EdgeSoC (in planning)
- **IcySoC chips approx. computing platforms** (ALP 180nm)
 - Diego
 - Manny
 - Sid



Pushing Beyond pJ/OP



GOPS/W



Closing The Accelerator Efficiency Gap with Agile Customization

- Brain-inspired (deep convolutional networks) systems are high performers in many tasks over *many domains*



leopard

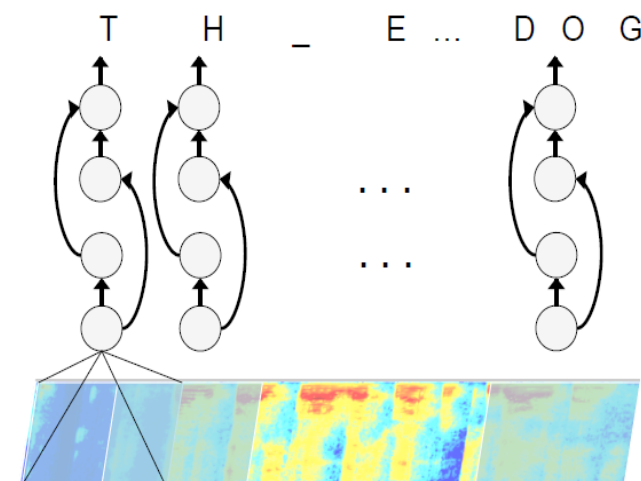


Image recognition

[Russakovsky et al., 2014]

CNN:
93.4% accuracy
(Imagenet 2014)
Human:
85% (untrained),
94.9% (trained)

[Karpahy15]

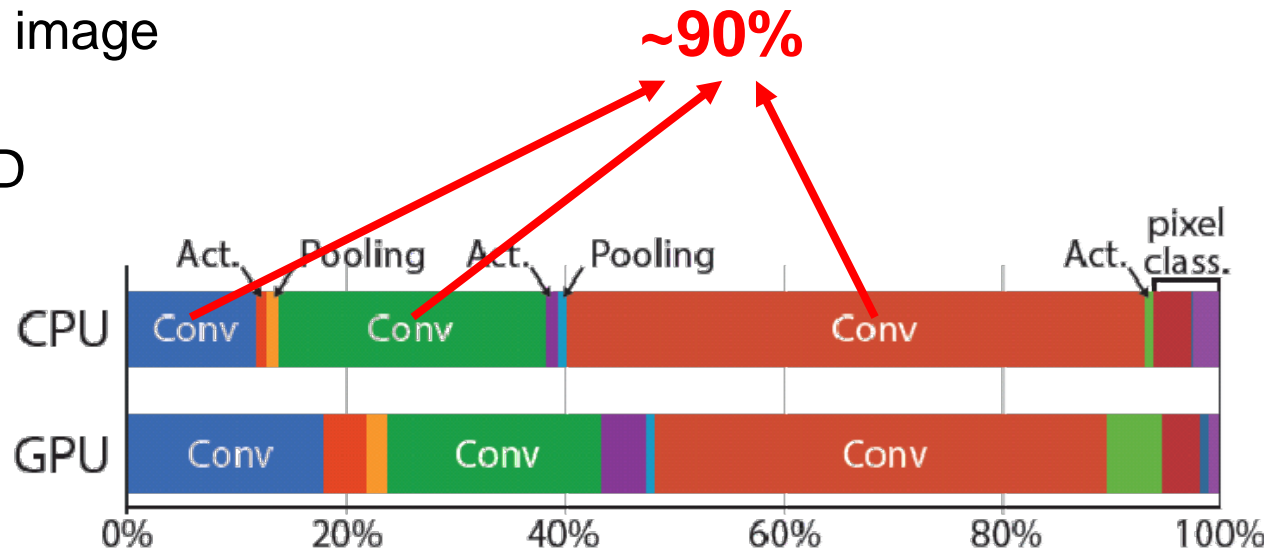


Speech recognition

[Hannun et al., 2014]

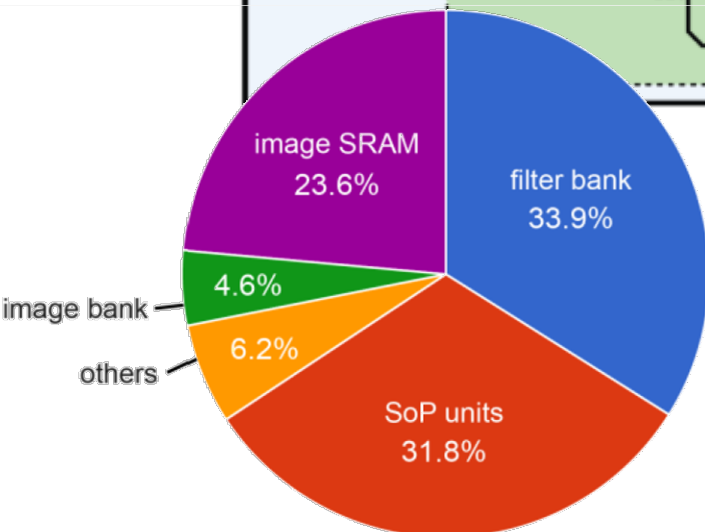
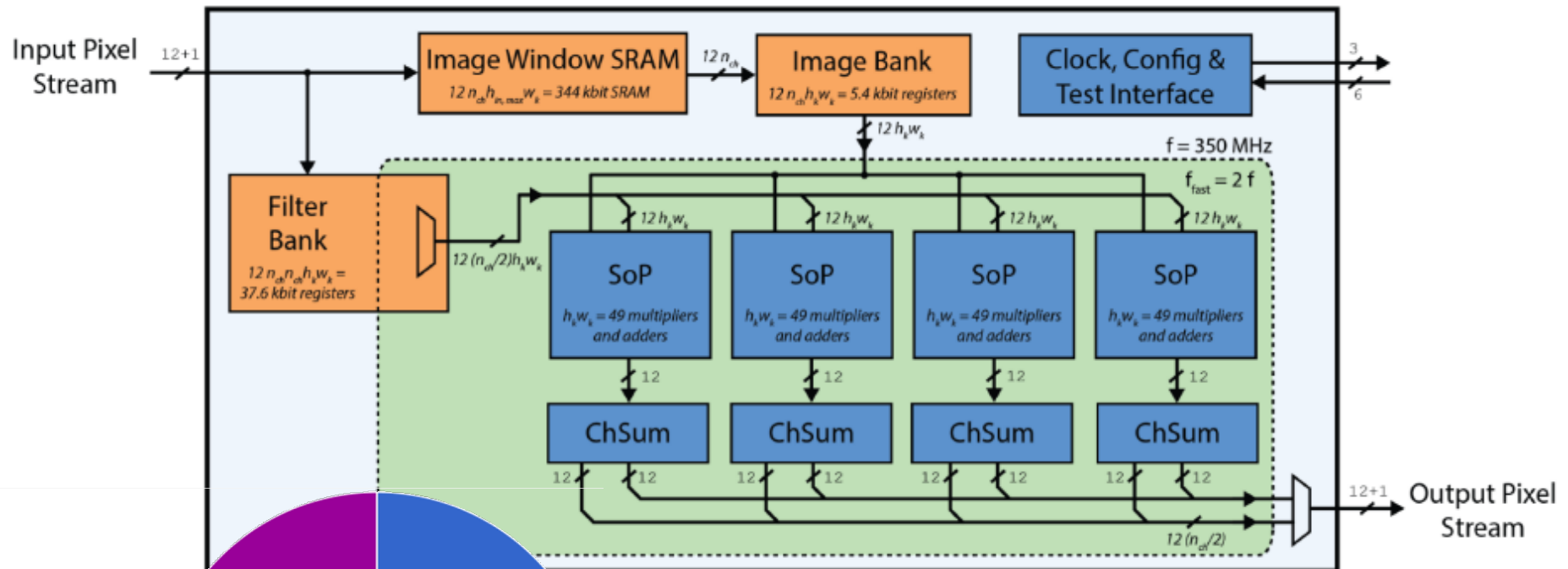
- Flexible** acceleration: learned CNN weights are “the program”

- Computational effort
 - 7.5 GOp for 320x240 image
 - 260 GOp for FHD
 - 1050 GOp for 4k UHD



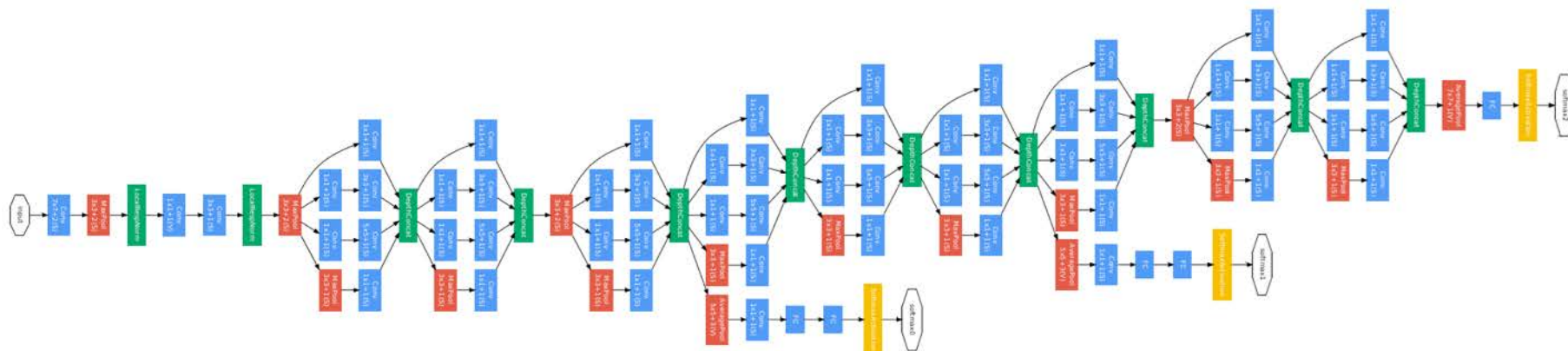
Origami chip





■ FP needed?

- 12-bit signals sufficient
- Input to classification double-vs-12-bit accuracy loss $< 0.5\%$ (80.6% to 80.1%)



Example: **GoogLeNet** [ILSVRC 2014 winner]

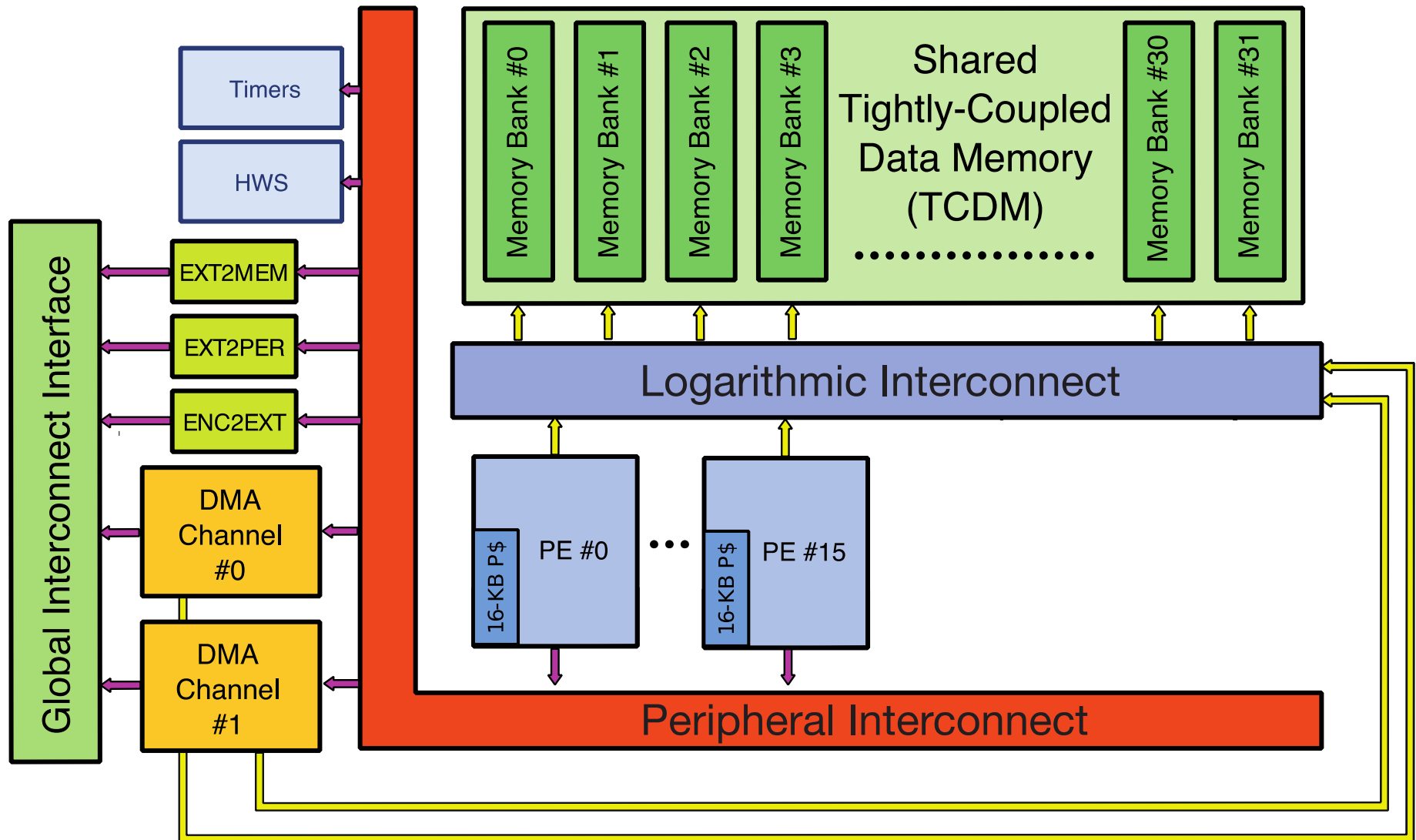
~ 7×10^6 parameters

~ 2.3×10^9 MAC operations on a **320x240** RGB image

Realtime (10 fps): ~**23 GMAC/s** performance

Realtime & Low-Power (10 fps @ 10mW): ~**2300 GMAC/s/W** efficiency

Origami core in 28nm FDSOI → GoogLeNet with ~10mW



- Approximation at the algorithmic side → Binary weights
- BinaryConnect [Courbariaux, NIPS15]
 - Reduce weights to a binary value -1/+1
 - Stochastic Gradient Descent with Binarization in the Forward Path

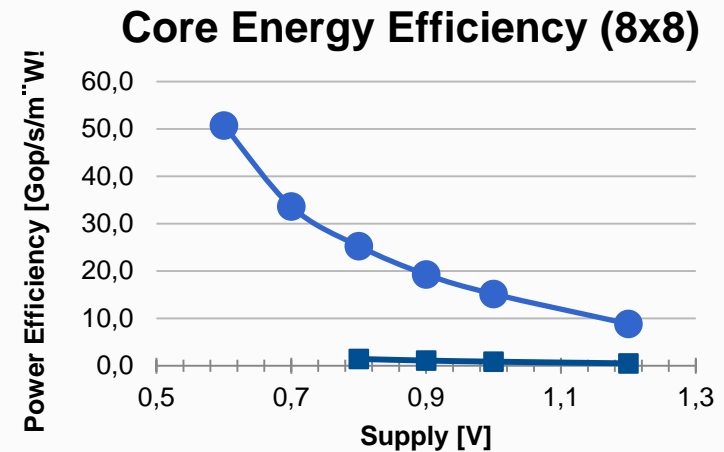
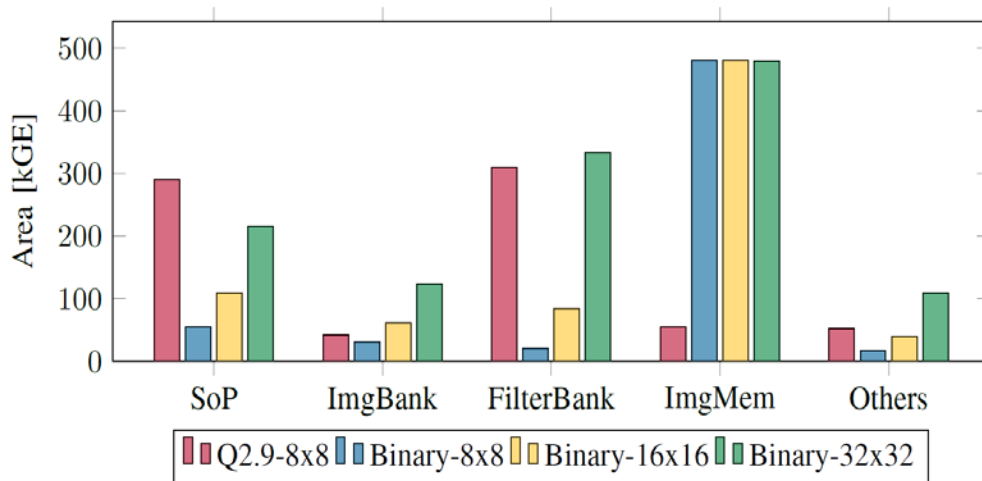
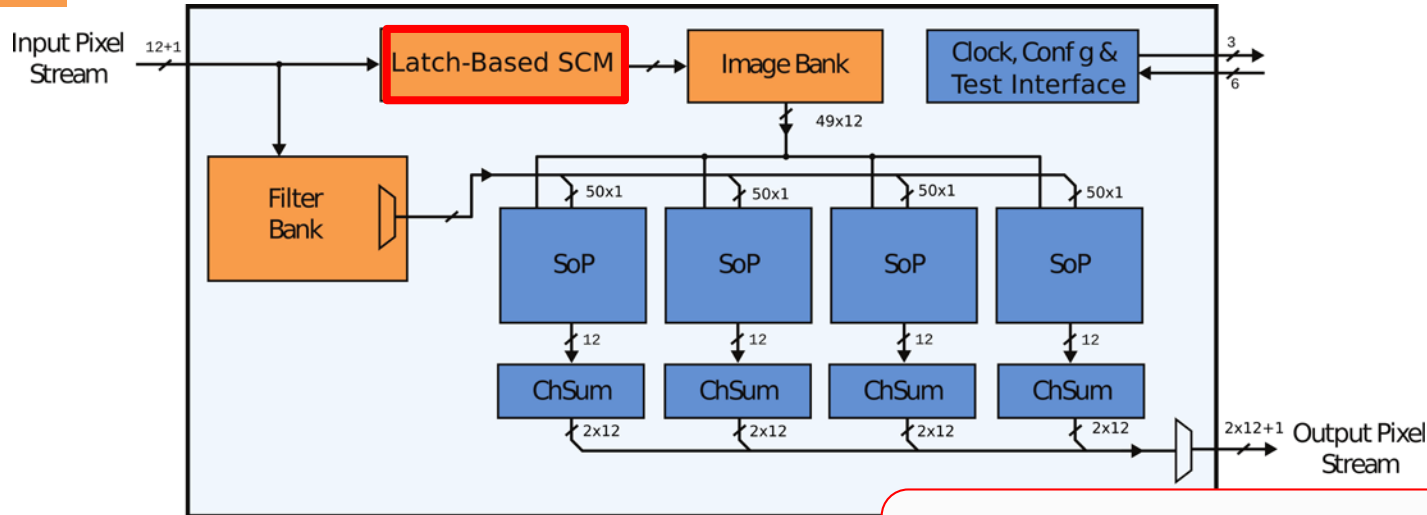
$$w_{b,stoch} = \begin{cases} -1 & p_{-1} = \sigma(w) \\ 1 & p_1 = 1 - p_{-1} \end{cases}$$

$$w_{b,det} = \begin{cases} -1 & w < 0 \\ 1 & w > 0 \end{cases}$$

- Learning large networks is still an issue with binary connect...
- Ultra-optimized HW is possible!
 - Power reduction because of arithmetic simplification
 - Major arithmetic density improvements
 - Area can be used for more energy-efficient weight storage
 - SCM memories for lower voltage → E goes with $1/V^2$

¹After the Yedi Master from Star Wars - “Small in size but wise and powerful” cit. www.starwars.com

SCM (down to 0.6V in 65nm)



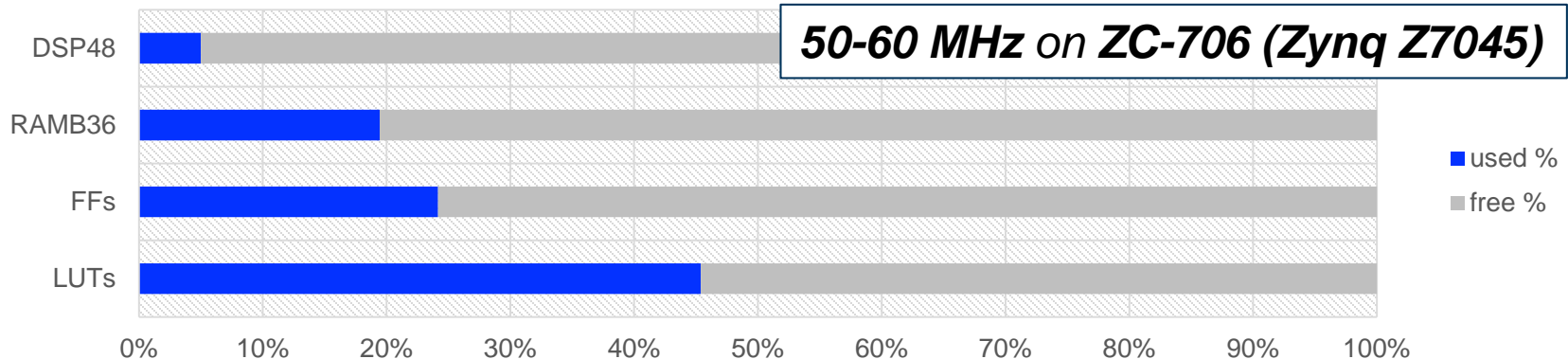
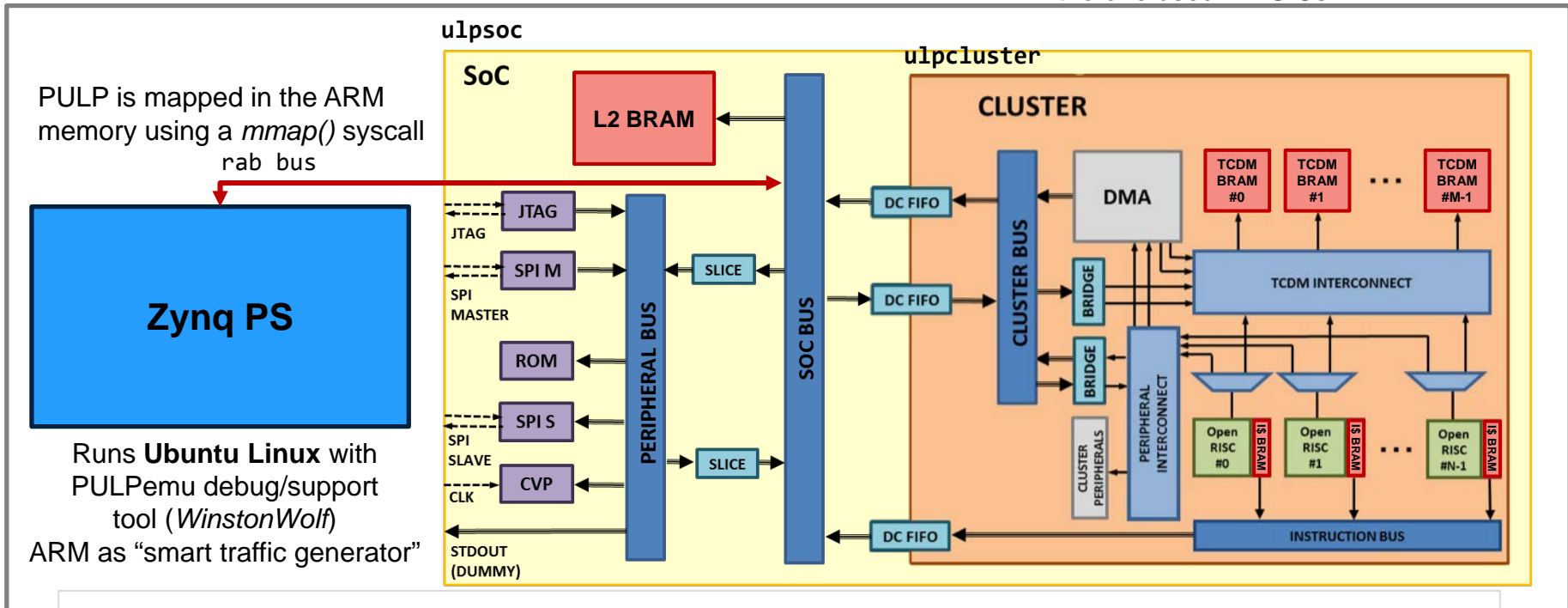
16x Energy efficiency improvement: 0.5pJ/OP → 50GOPS/mW

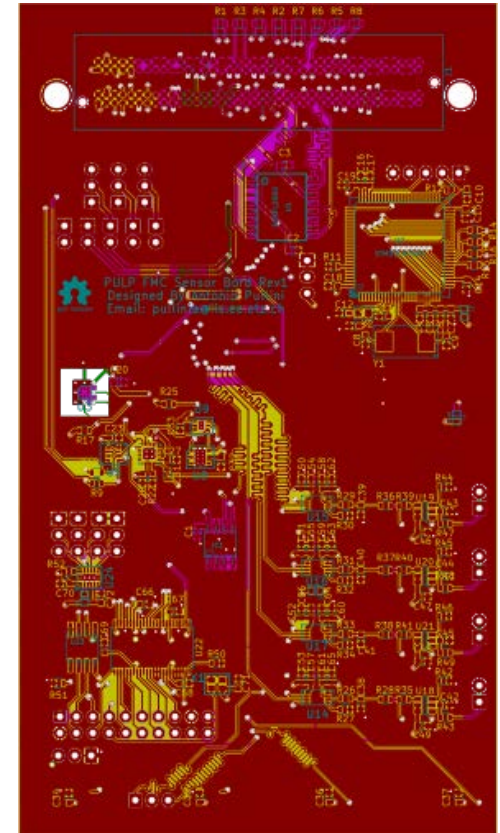
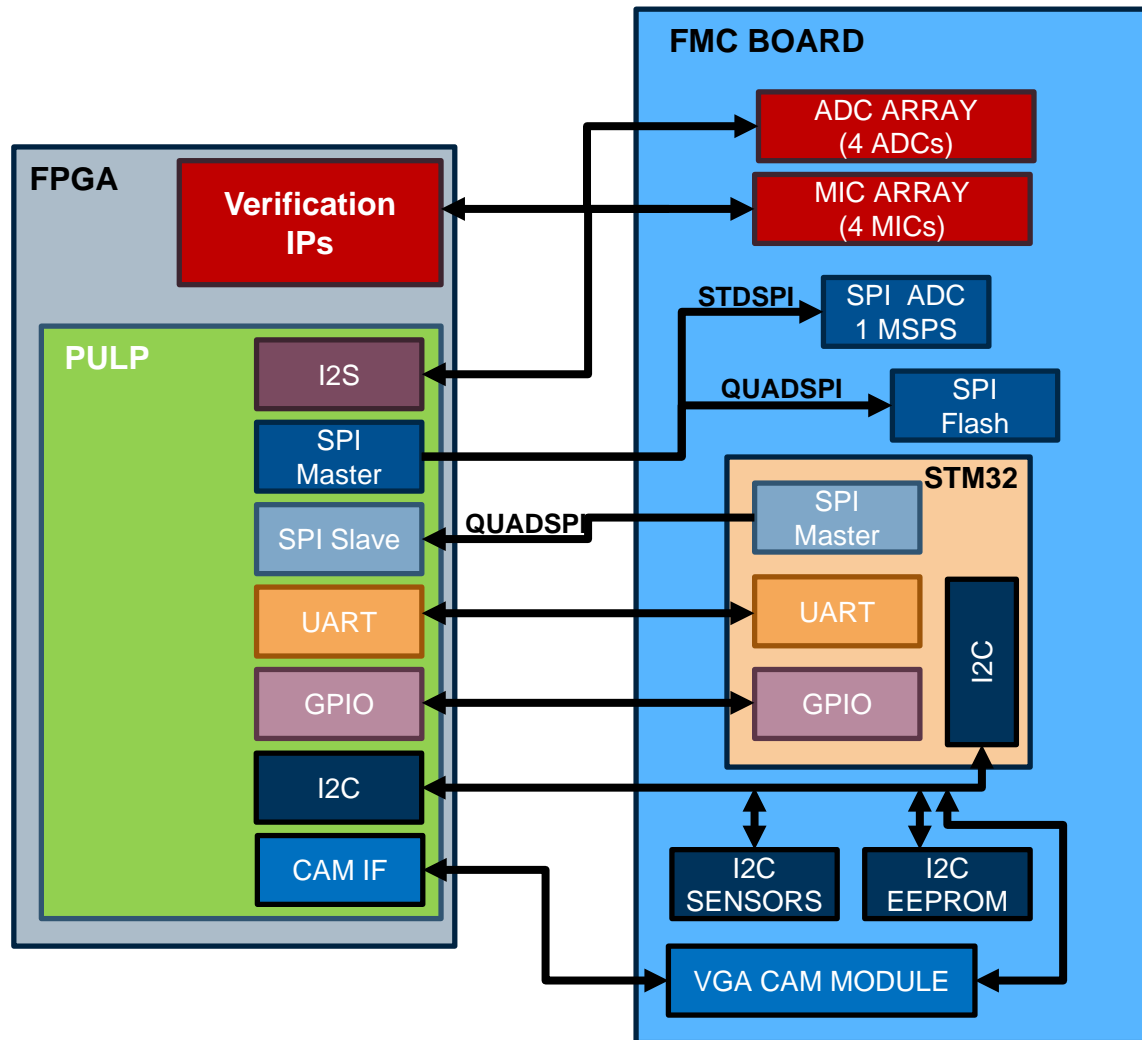
PULP infrastructure: not only hardware...



pulpemu

The PULP cluster is RTL-equivalent to the one used in ASICs





Board design is **open source**
(part of **PULPino** release)

Two toolchains:

- GCC 4.9
- LLVM 3.7



Transparent compiler support for most HW features:

- Hardware loops, Post-Increment, Register offset, MAC, vectorization...
- Bit-counting operations (1.ff1, 1.f11, 1.cnt) supported through intrinsics
- All extensions can be disabled via compiler flags
- *libc* support (I/O based on semi-hosting)

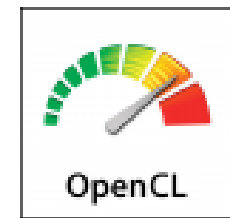
Performance in industry standard benchmark:

- **2.37** CoreMark/MHz with **GCC**
- **2.16** CoreMark/MHz with **LLVM**

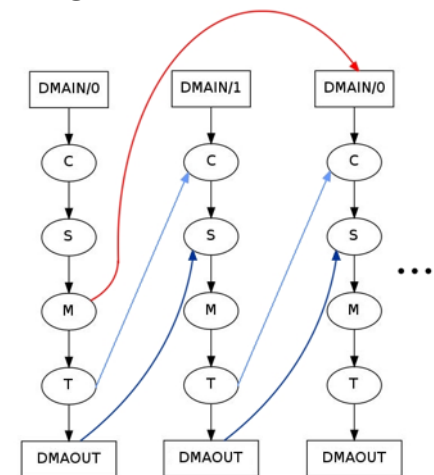
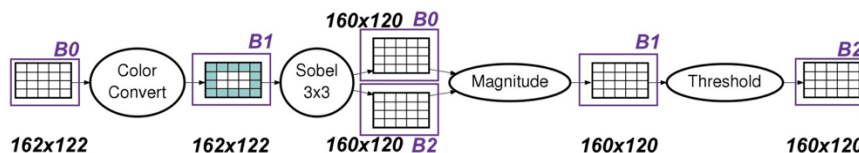


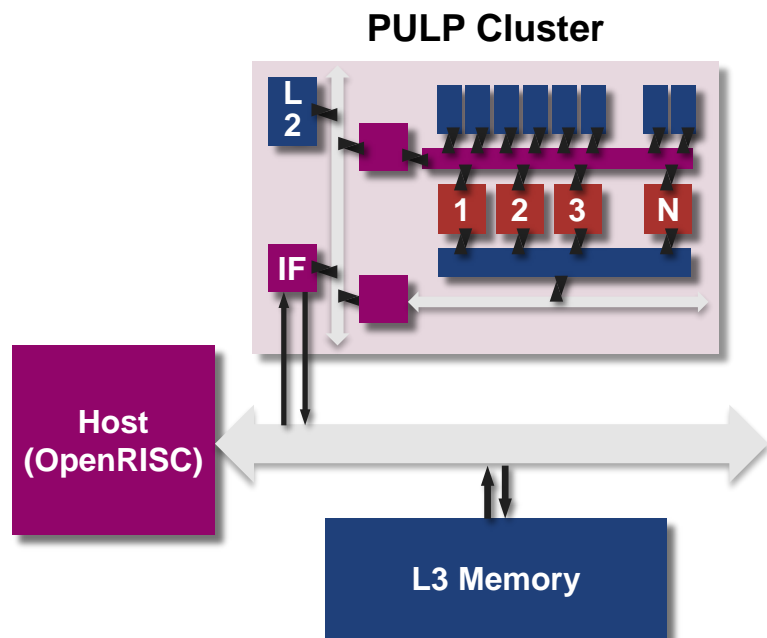
Full **OpenMP v3.0** support:

- Optimized for performance and power with HW support (event unit + HW test-and-set)
- 4-core parallel region: **70 cycles**
- barrier: **25 cycles**
- critical section: **60 cycles**

**OpenVX** support on PULP virtual platform:

- Vision acceleration layer for mobile/embedded, C-based
- Application = Directed Acyclic Graph with data as linkage
- *Khronos* runtime running on host OR10N
- Kernels are parallized with OpenMP on PULP
- Automatic tiling and insertion of DMA transfers





Virtual platform implementation:

- C++ for fast native simulation
- *Python* for instantiation + configuration
- SWIG is used to generate C++ stubs from Python
- Any ISS can be integrated (*or1ksim*, *gdb* simulator, *riscv* on-going)

Timing model:

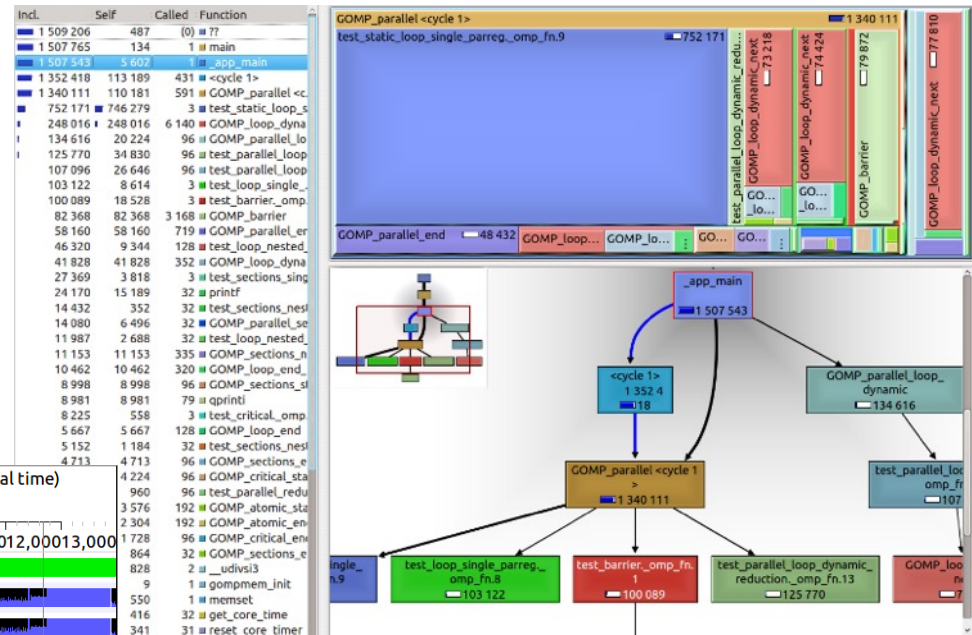
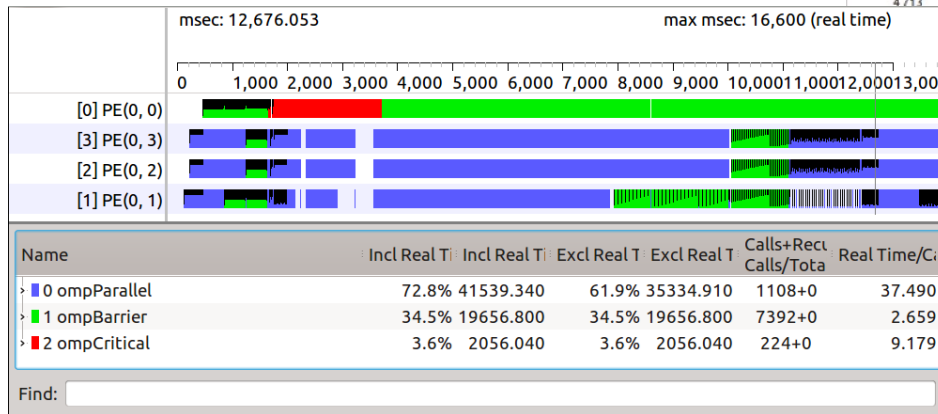
- Fully-event based, instances can generate events at specific time
- Includes timing models for interconnects, DMACs, memories...

Simulation performance:

- Around 1MIPS simulation speed
- Functionally aligned with HW
- Timing accuracy is within **10-20%** of target HW

Profiling based on **KCachegrind**:

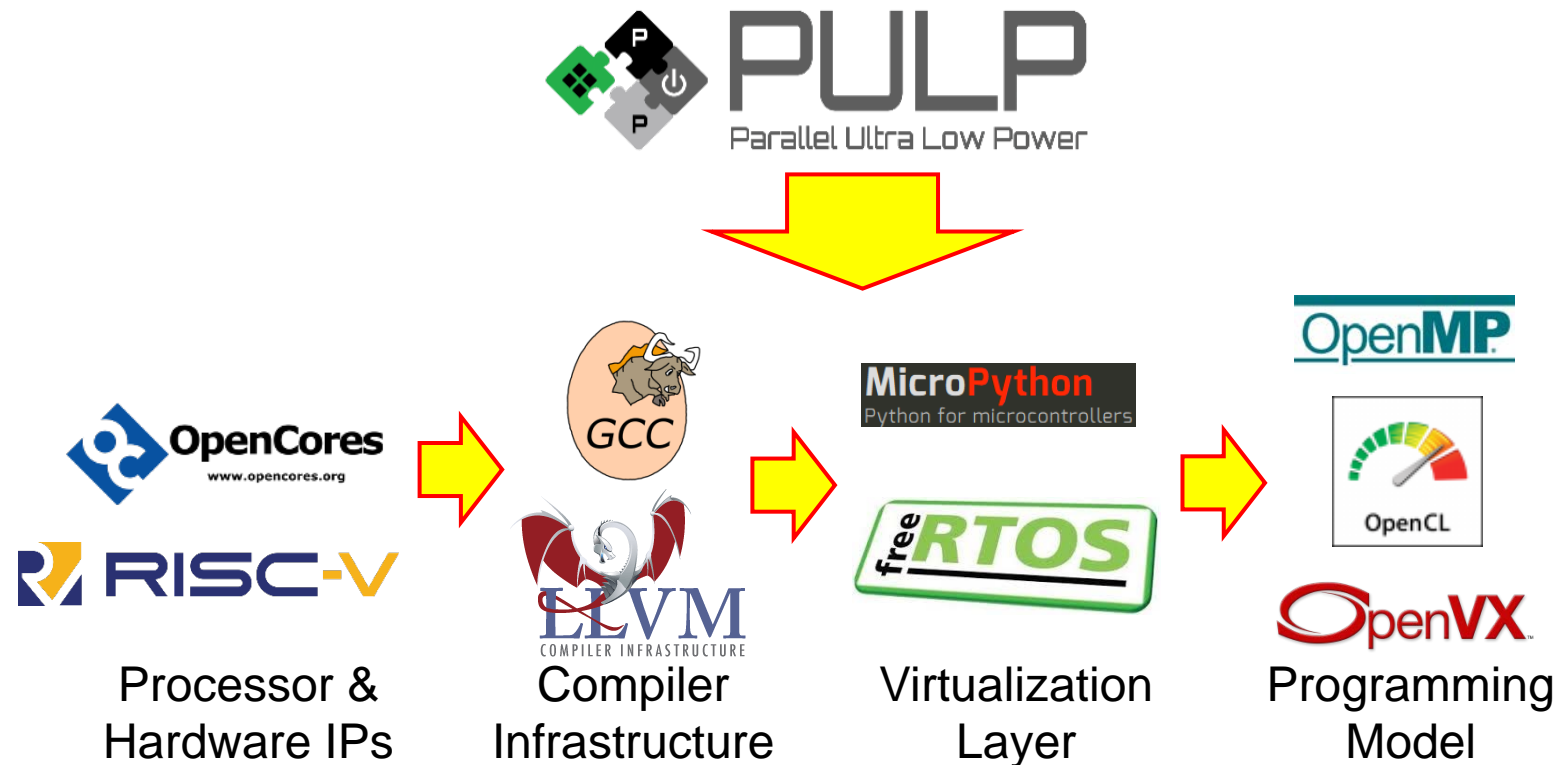
- Supports PC traces from RTL and virtual platform (FPGA emulator on-going)
- Several events can be caught (PC, cycles, I\$ miss, stalls...)

Debug with **GDB**:

- Supports RTL and virtual platform
- Uses a bridge to inject JTAG requests
- Main GDB features work (step-by-step, breakpoints, introspection)
- Forked from *minsoc* project

Open Source Parallel ULP computing for the IoT

(sub)-pJ/op computing platform - **let's make it Open!**



[About](#)[Release Plan](#)[Resources](#)[Download](#)<https://github.com>

git

We are happy to share our FREE and OP

You can download the entire source code, test programs,
completely for free under the [Solderpad license](#).**Hundreds of git forks in a few weeks!****EE Times**Connecting the Global
Electronics Community[Home](#) [News](#) [Opinion](#) [Messages](#) [Authors](#) [Video](#) [Slideshows](#) [Teardown](#) [Education](#) [EEL](#)[designlines](#)[Android](#)[Automotive](#)[Embedded](#)[Industrial Control](#)[Internet of Things](#)**BREAKING NEWS****NEWS & ANALYSIS: Severance Clash in Microchip/Atmel Merger**[designlines](#) [INTERNET OF THINGS](#)**News & Analysis**

Open-Source Processor Core Ready For IoT

Peter Clarke

3/31/2016 10:48 AM EDT

4 comments

NO RATINGS

2 saves

[LOGIN TO RATE](#)[f Like](#)

105

[t Tweet](#)

in

[Share](#)

165

[G+1](#)

20

Researchers at ETH Zurich (Swiss Federal Institute of Technology in Zurich) and the University of Bologna have developed PULPino, an open-source processor optimized for low power consumption and application in wearables and the Internet of Things (IoT).

Thanks for your attention!!!



www.pulp-platform.org

www-micrel.deis.unibo.it/pulp-project

iis-projects.ee.ethz.ch/index.php/PULP

