

Sub-pJ-Operation Scalable Computing

Why, When, How?

Luca Benini

IIS-ETHZ & DEI-UNIBO

<http://www.pulp-platform.org>



IoT: Near-Sensor Processing

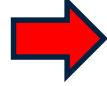
	INPUT BANDWIDTH	COMPUTATIONAL DEMAND	OUTPUT BANDWIDTH	COMPRESSION FACTOR
■ Image				
Tracking: [*Lagroce2014]	80 Kbps	1.34 GOPS 100 MOPS 7.7 MOPS 150 MOPS	0.16 Kbps 0.02 Kbps 0.02 Kbps 0.08 Kbps	500x 12800x 120x 200x
■ Voice/Sound				
Speech: [*VoiceControl]	256 Kbps			
■ Inertial				
Kalman: [*Nilsson2014]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



Extremely compact output (single index, alarm, signature)



Computational power of ULP µControllers is not enough

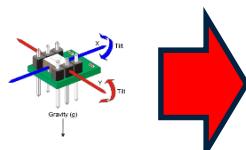
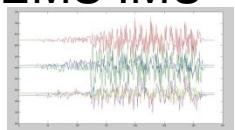


Parallel workloads

System View

Sense

MEMS IMU

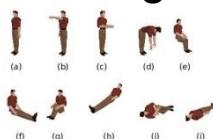


Analyze and Classify

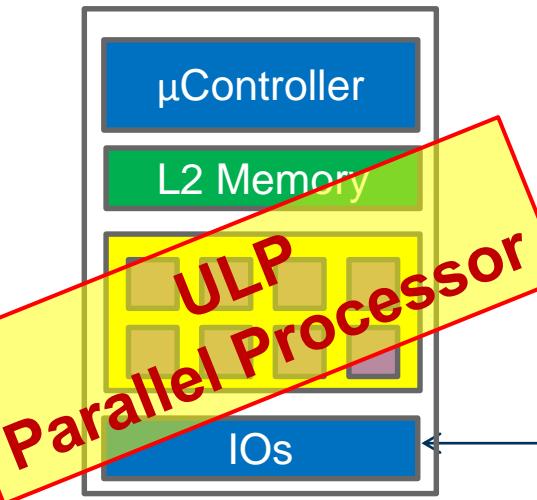
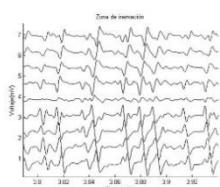
MEMS Microphone



ULP Imager



EMG/ECG/EIT



$1 \div 2000$ MOPS
 $1 \div 10$ mW



Transmit

Short range, 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**

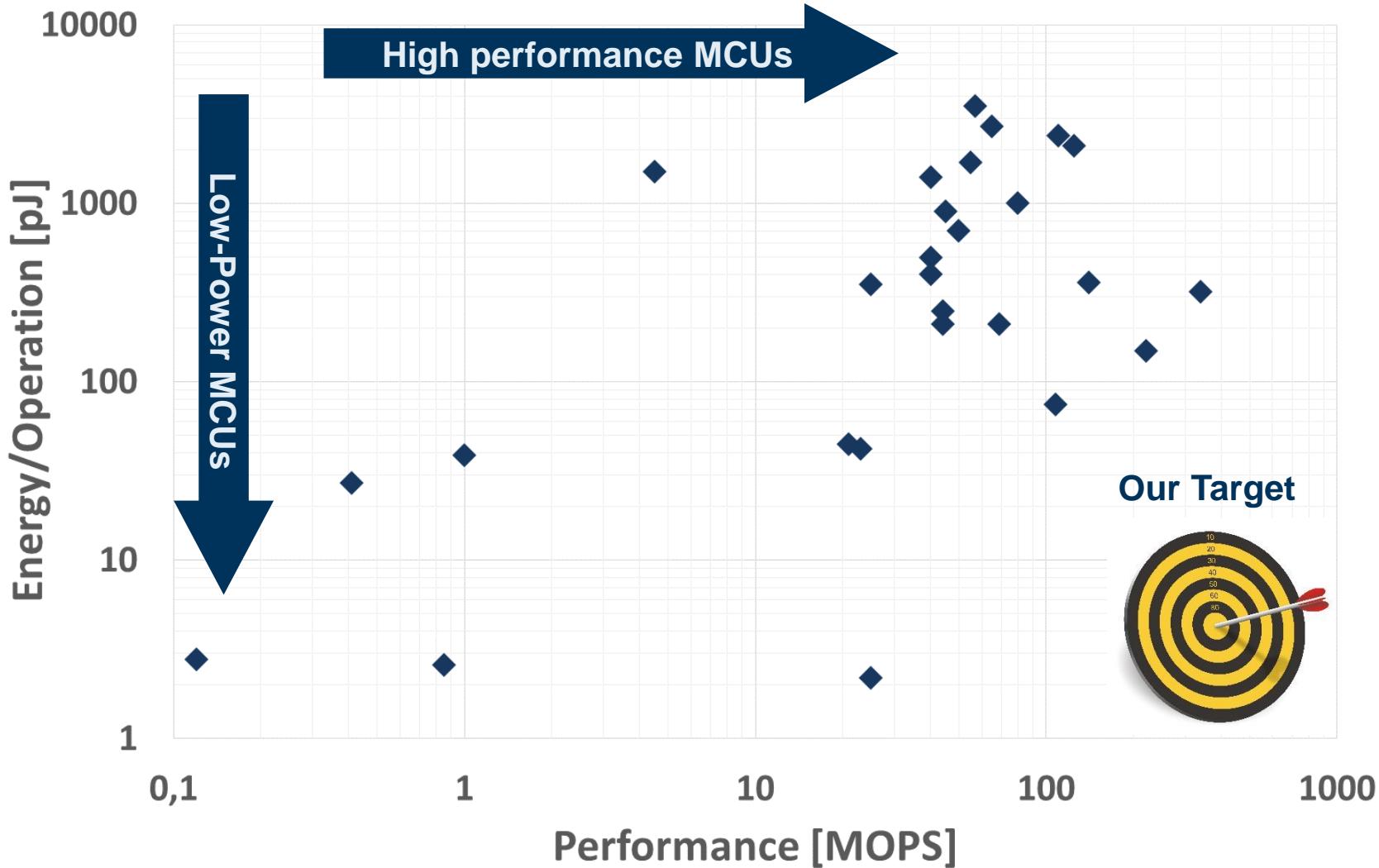
**Idle: ~1μW
Active: ~ 50mW**

Microcontrollers Landscape



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

*not exhaustive

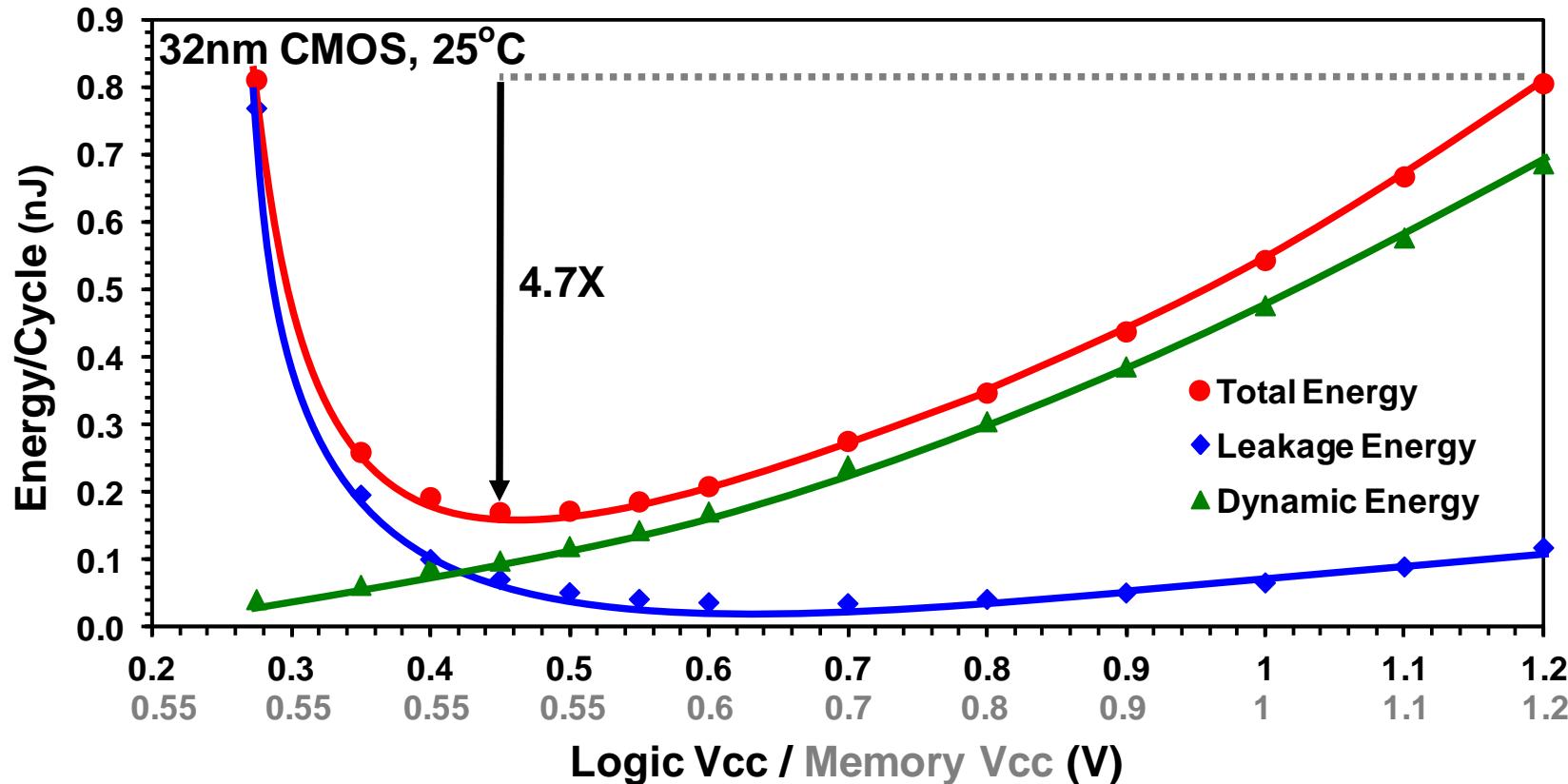


Near-Threshold Multiprocessing



Minimum energy operation

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

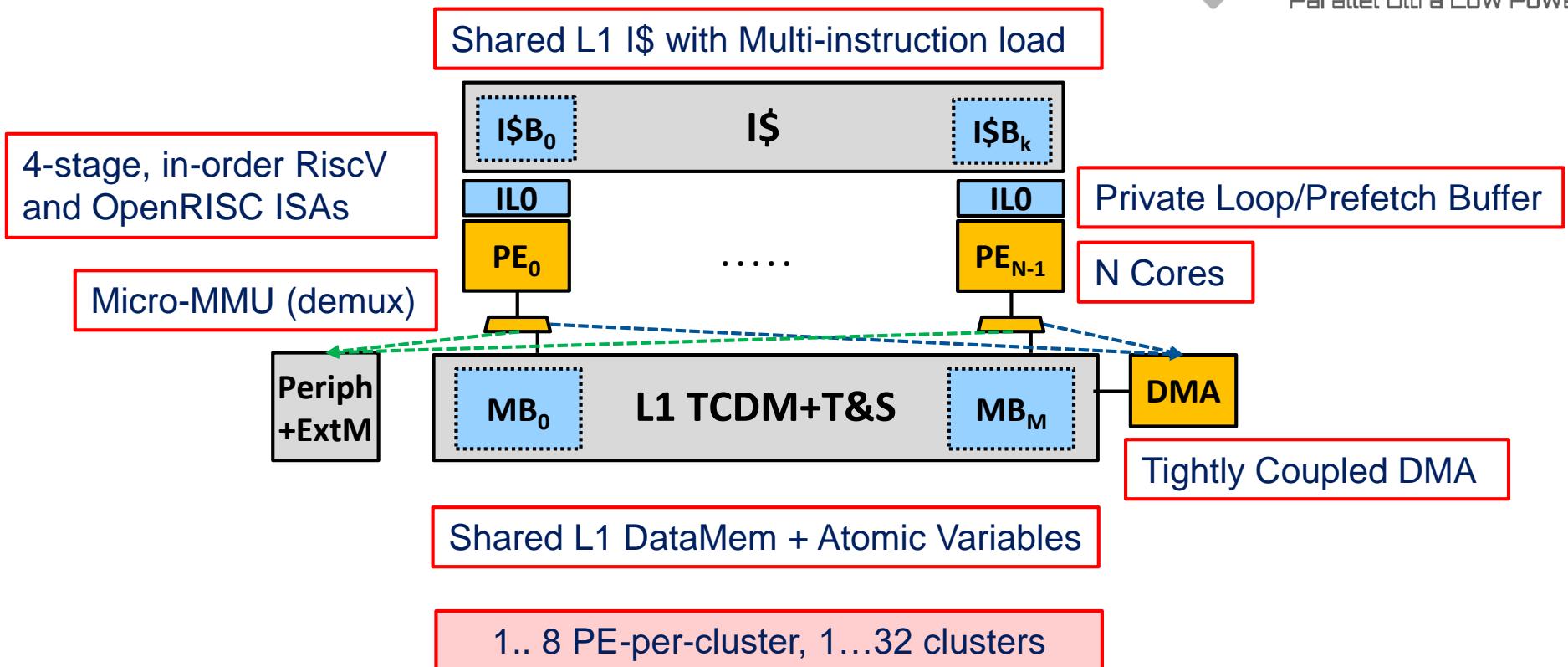
Near-Threshold Multiprocessing

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA



PULP
Parallel Ultra Low Power

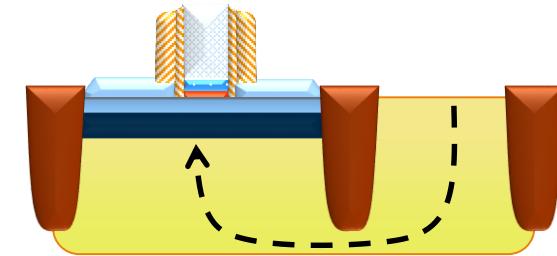
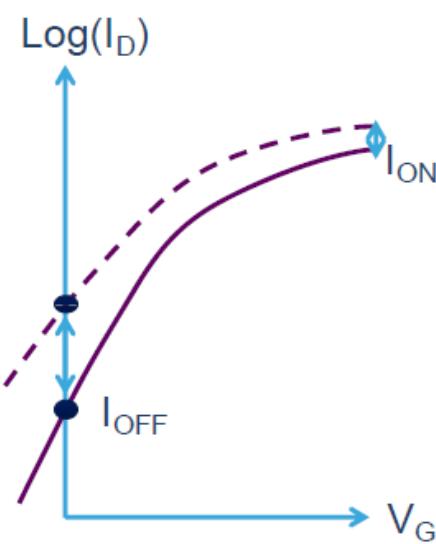
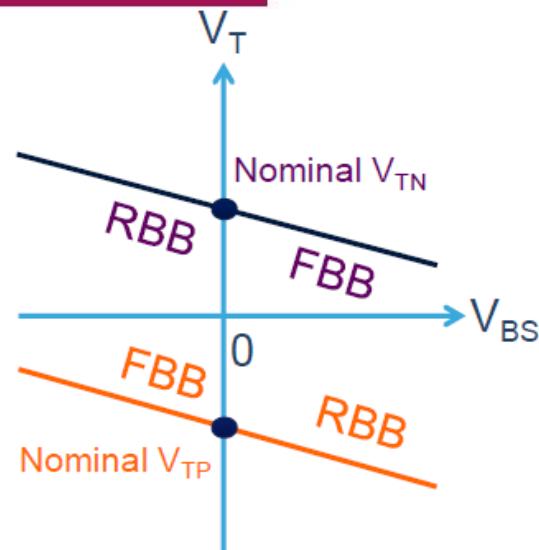
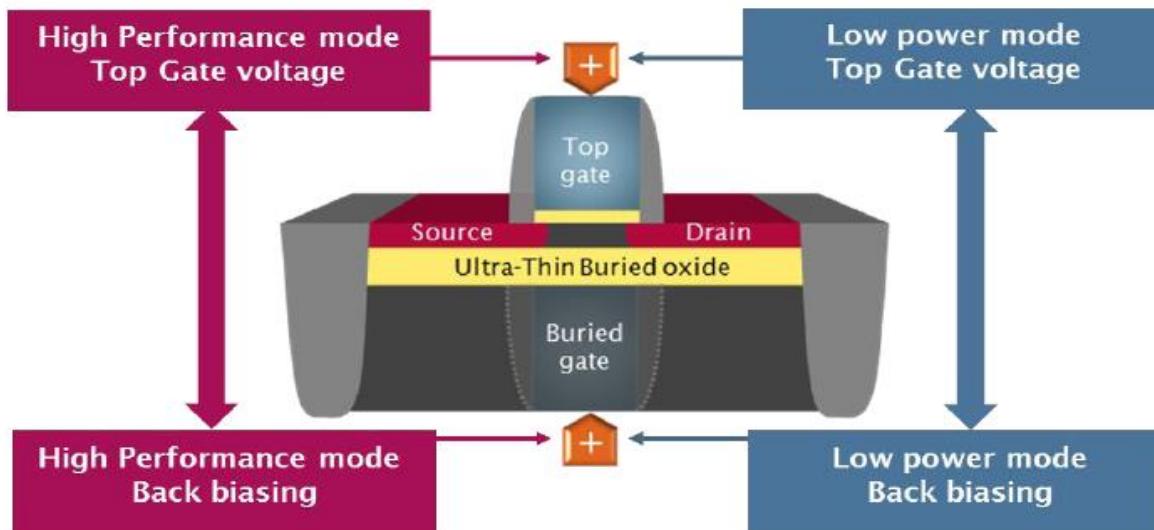
Ultra-low power scalable computing



NT but parallel → Max. Energy efficiency when Active

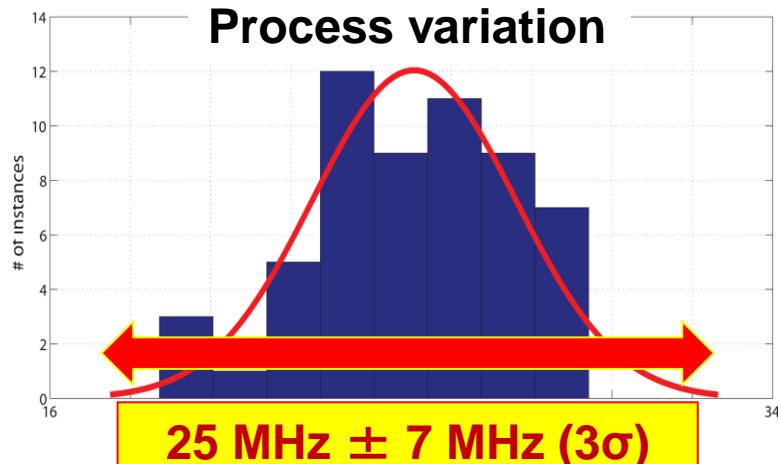
+ strong PM for (partial) idleness

Near threshold FDSOI technology



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

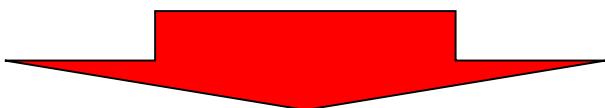
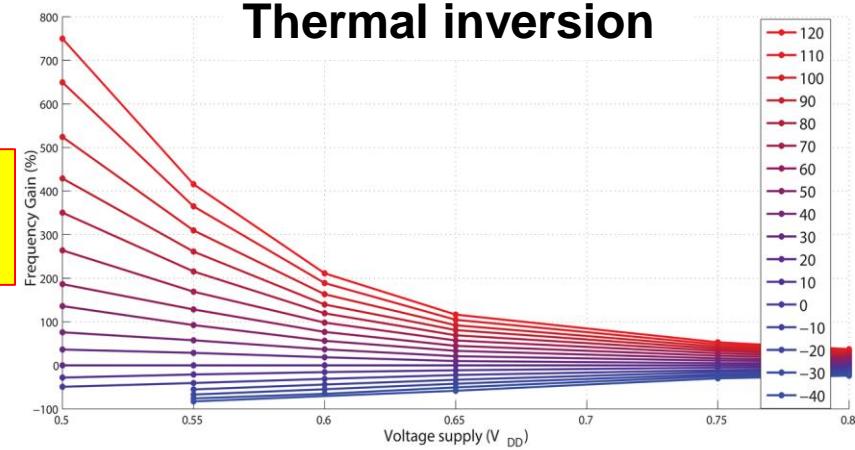
Body Biasing for Variability Management



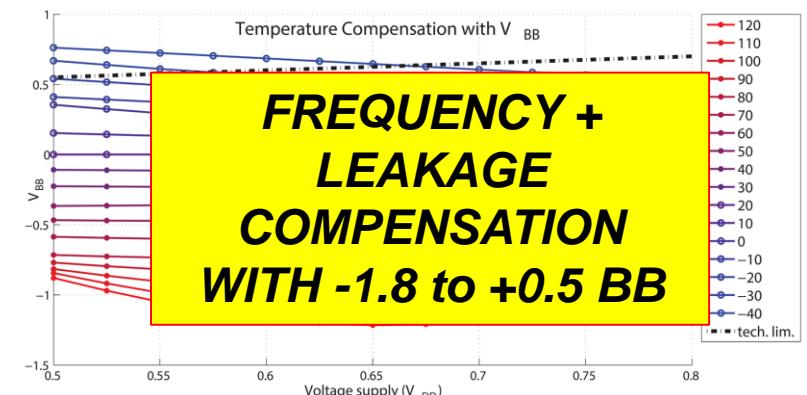
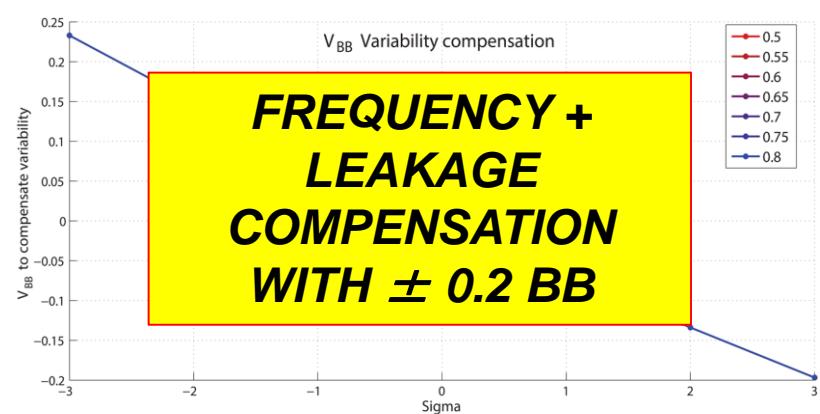
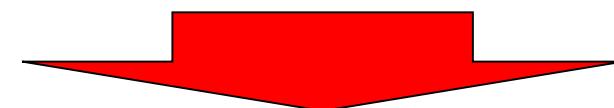
120° C



Thermal inversion



RVT transistors
FBB/RBB



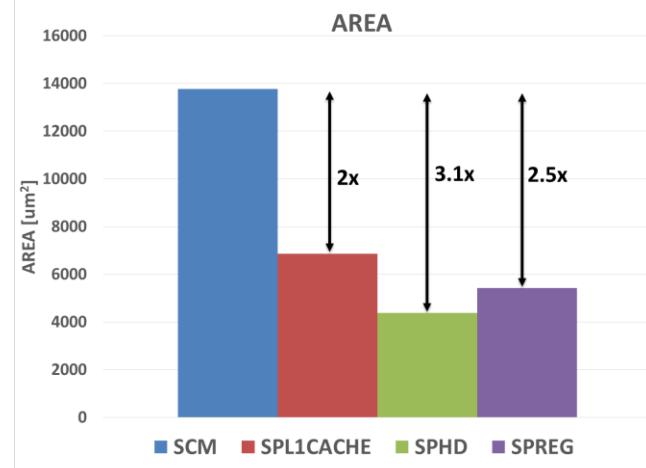
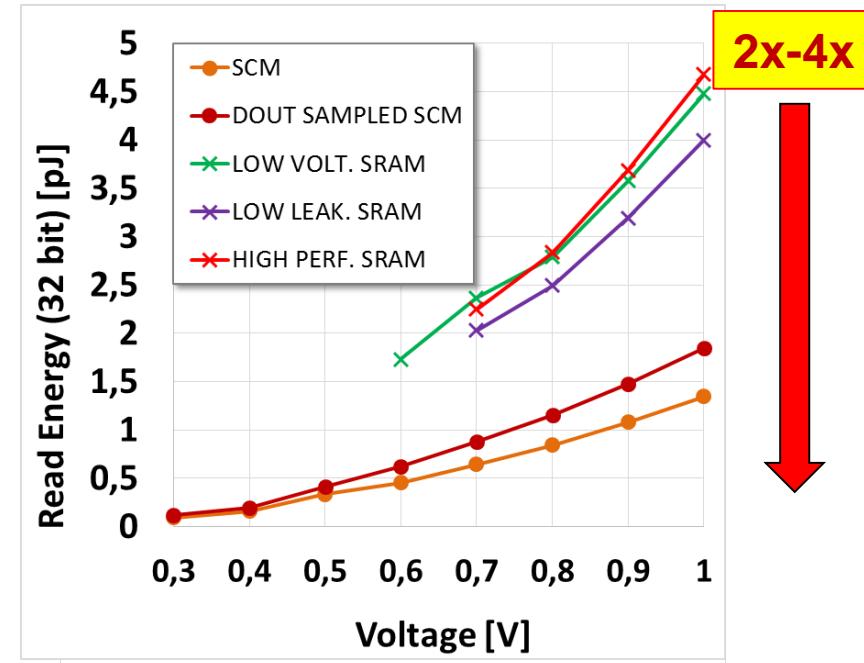
ULP (NT) Bottleneck: Memory



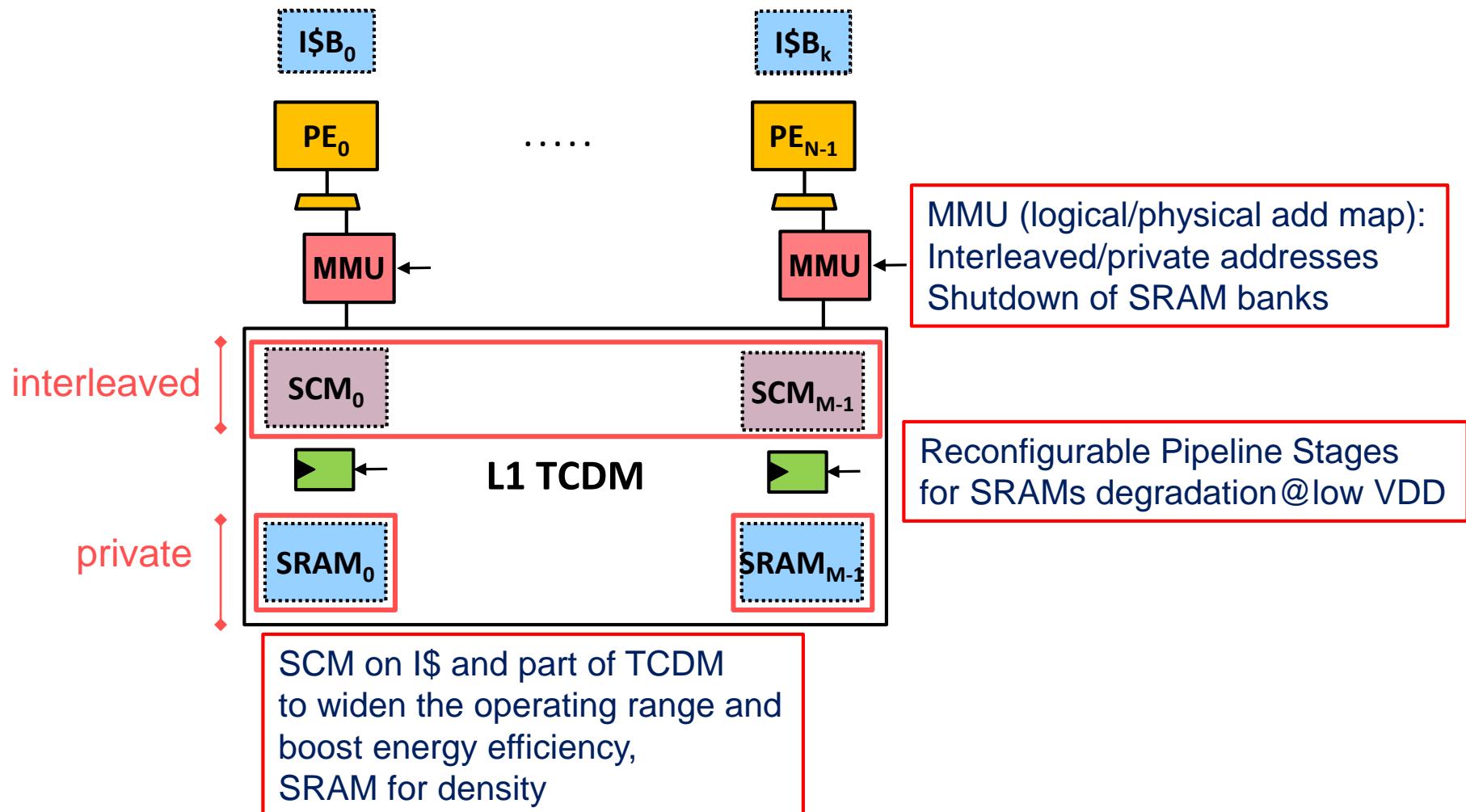
ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

256x32 6T SRAMS vs. SCM

- “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)



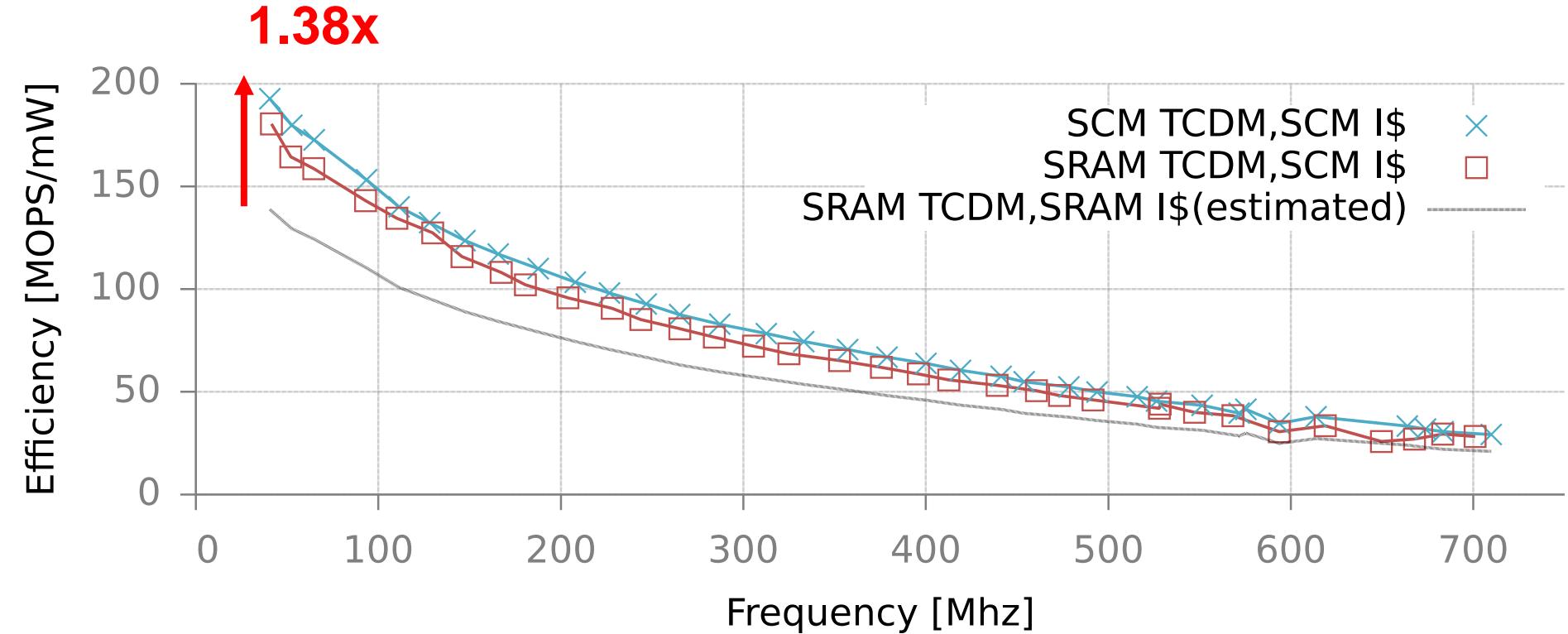
Heterogeneous + Reconfigurable Memory Architecture



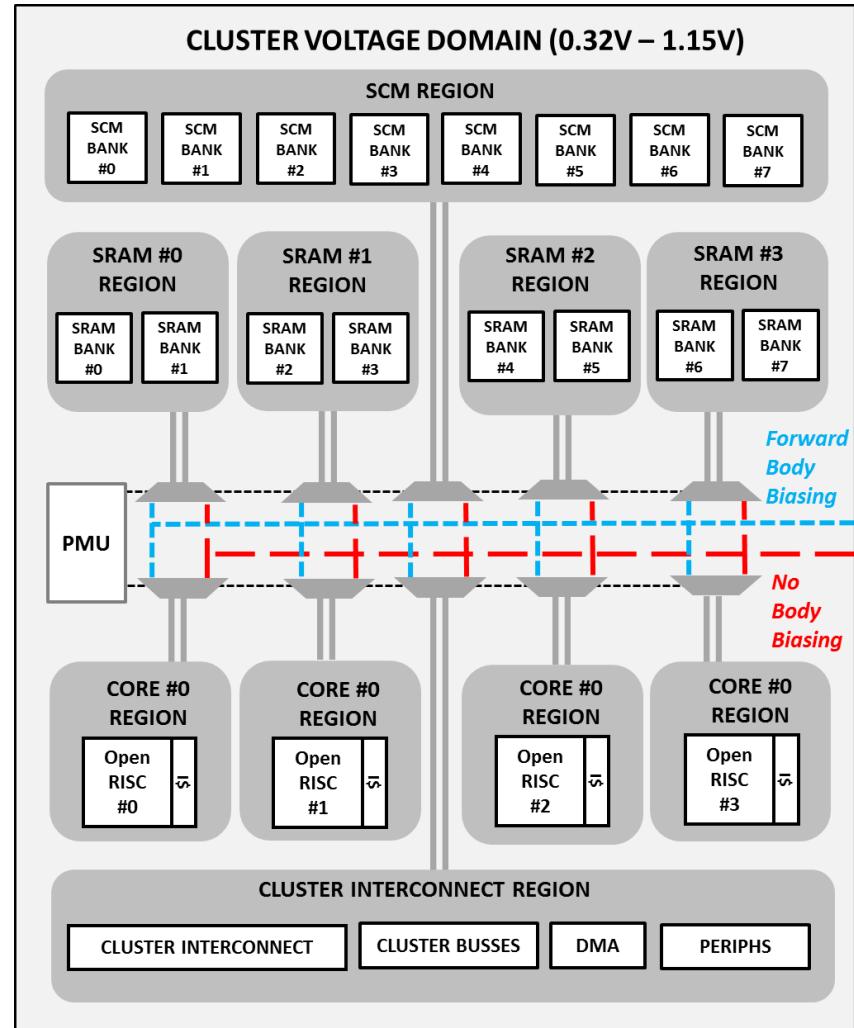
Heterogeneous Memory Energy Efficiency Boost



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

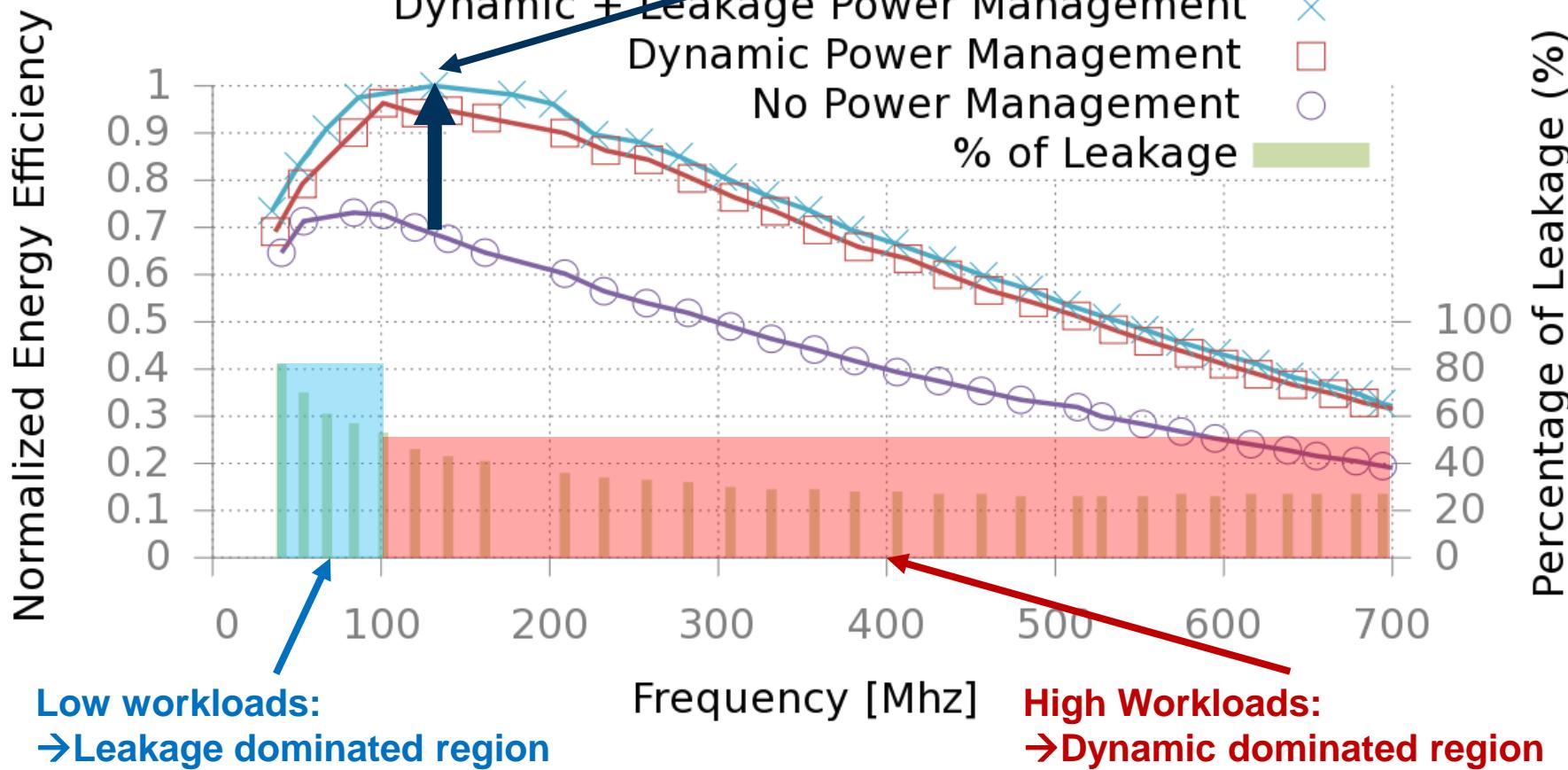


- The cluster is partitioned in separate **clock gating and body bias regions**
- Body bias multiplexers (BBMUXes)** control the well voltages of each region
- A **Power Management Unit (PMU)** automatically manages transitions between the operating modes
- Power modes of each region:
 - Boost mode:** active + FBB
 - Normal mode:** active + NO BB
 - Idle mode:** clock gated + NO BB



Sequential processing with selective FBB

Energy efficiency improvement in best energy point: 1.6x

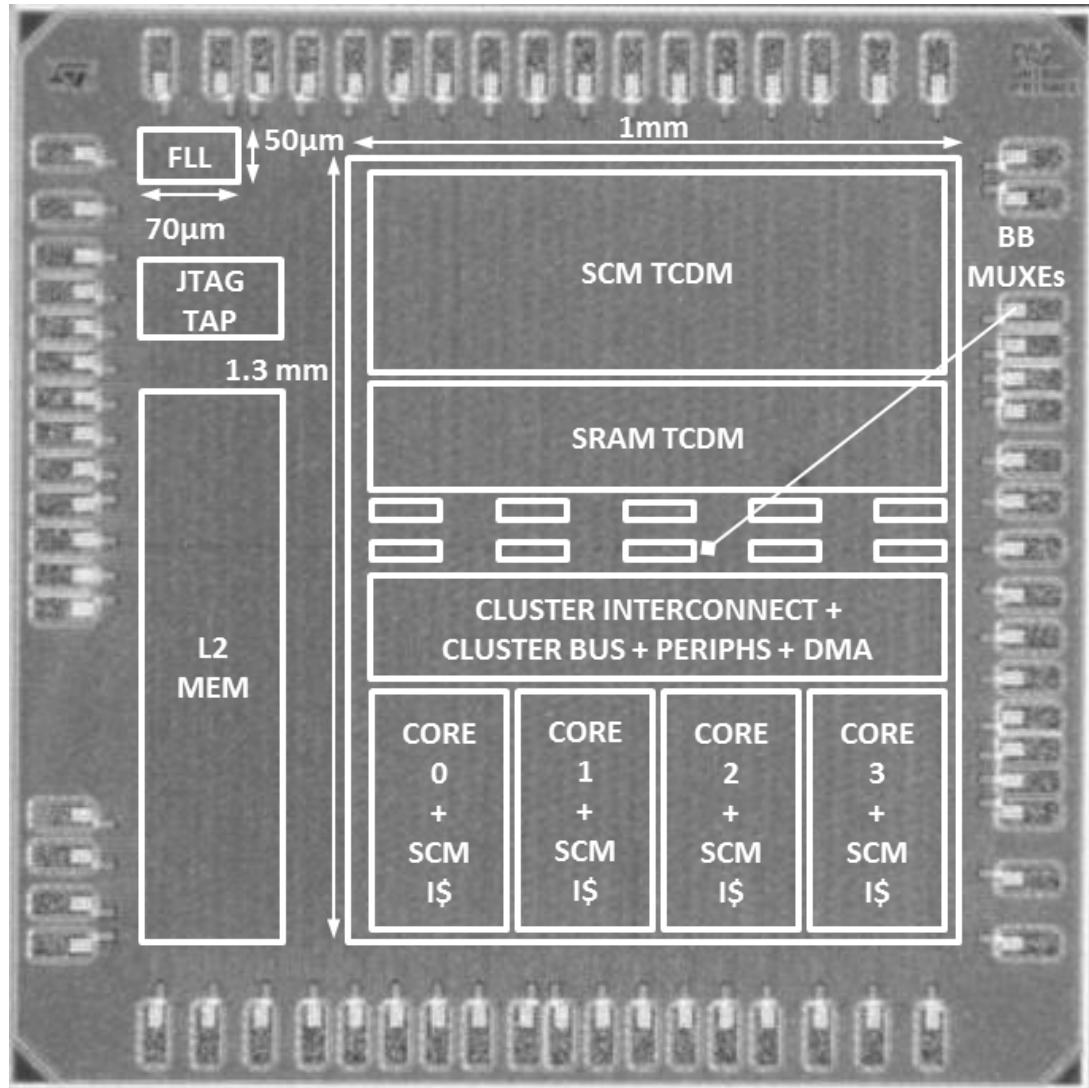


Sequential Processing, Temperature=25° C, FBB=1.75V

Silicon Results



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

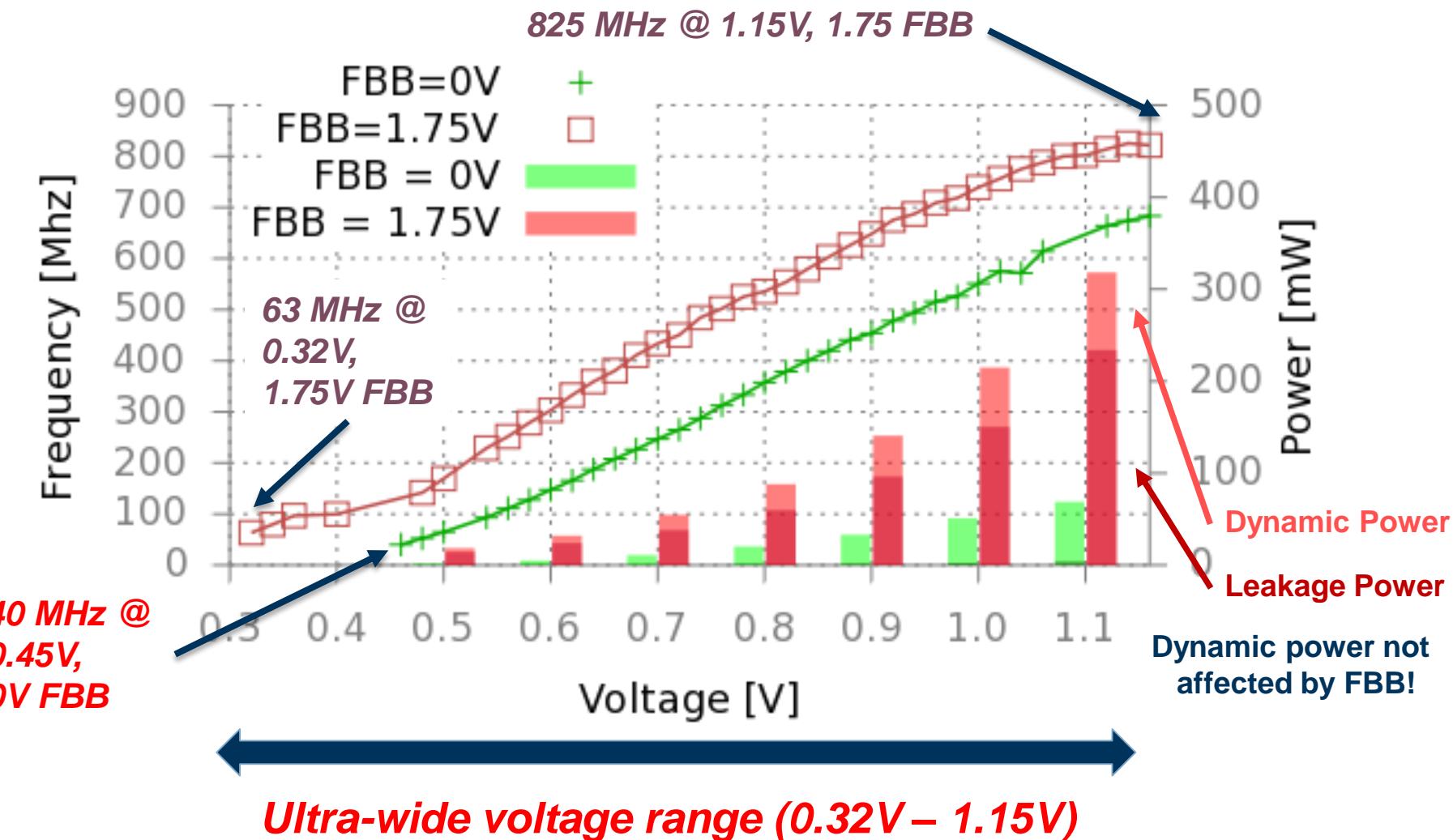


Technology	UTBB FD-SOI 28nm
Transistors	Flip well $L = 24\text{ nm}$
Cluster area	1.3 mm2
VDD range (memories)	0.32V - 1.15V (0.45 – 1.15V)
BB range	0V - 1.75V
SRAM macros	8 x 32 kbit (TCDM)
SCM macros	16x4 kbit (TCDM) 4x 2x4 kbit (I\$)
Gates	200K
Frequency range	NO BB: 40.5-710 MHz MAX FBB: 63.5 - 825 MHz
Power range	NO FBB: 0.56 - 85 mW MAX FBB: 6.9 - 480 mW

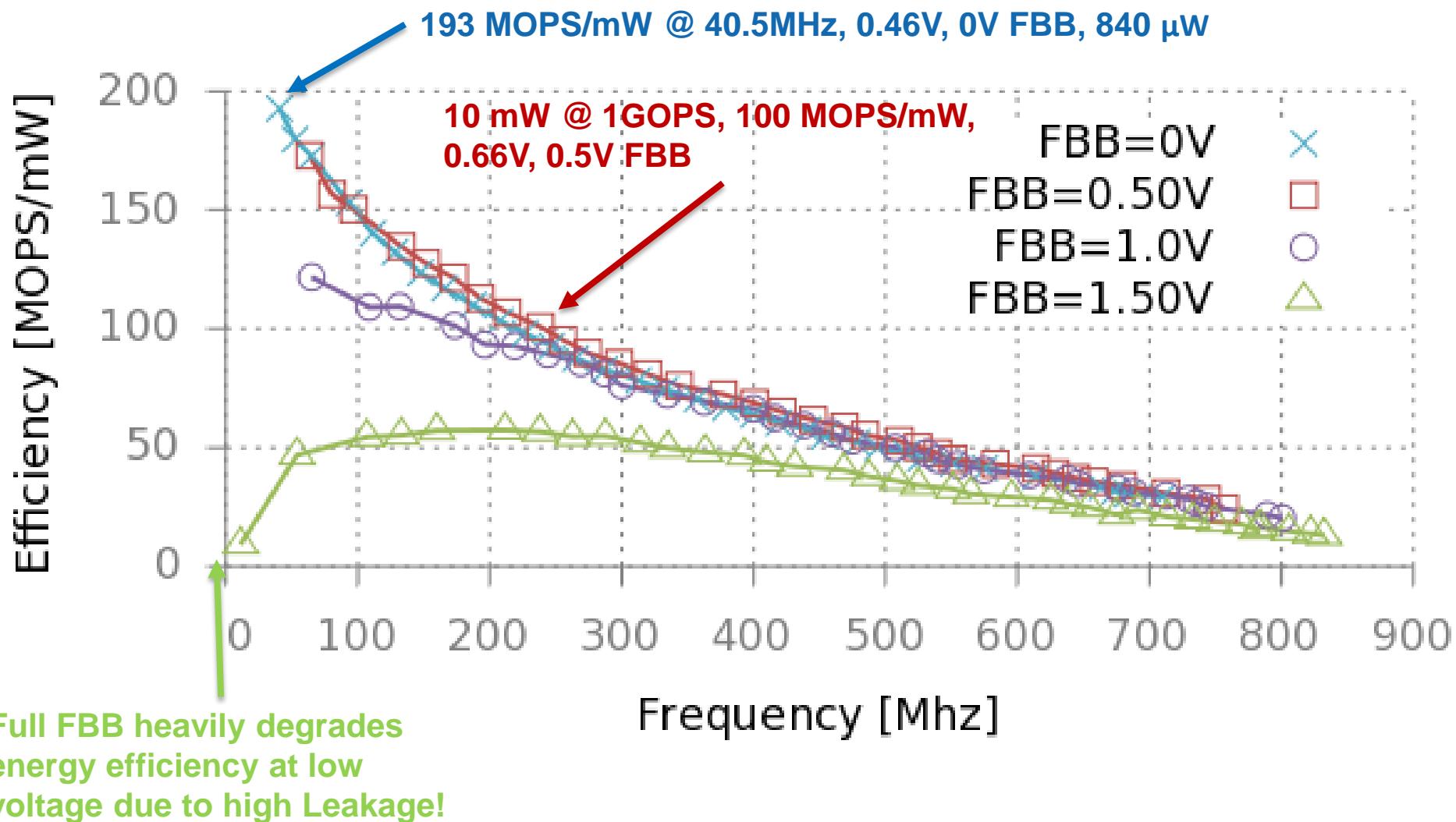
Hot Chips 15 V1 → Cool Chips 16 V2

193MOPS/mW @162MOPS → ~5pJ/OP

Cluster Performance



Cluster Energy Efficiency



Approximate Computing to the Rescue



Approximate?



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Less-than-perfect results perceived as correct by the users
e.g. image processing (filtering)



RGB to GRayscale



RGB to GRayscale (+ 10% error)

Approximation is not always acceptable
→ Application and program phase dependent!

Approximate Storage?

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

- **Retention voltage**

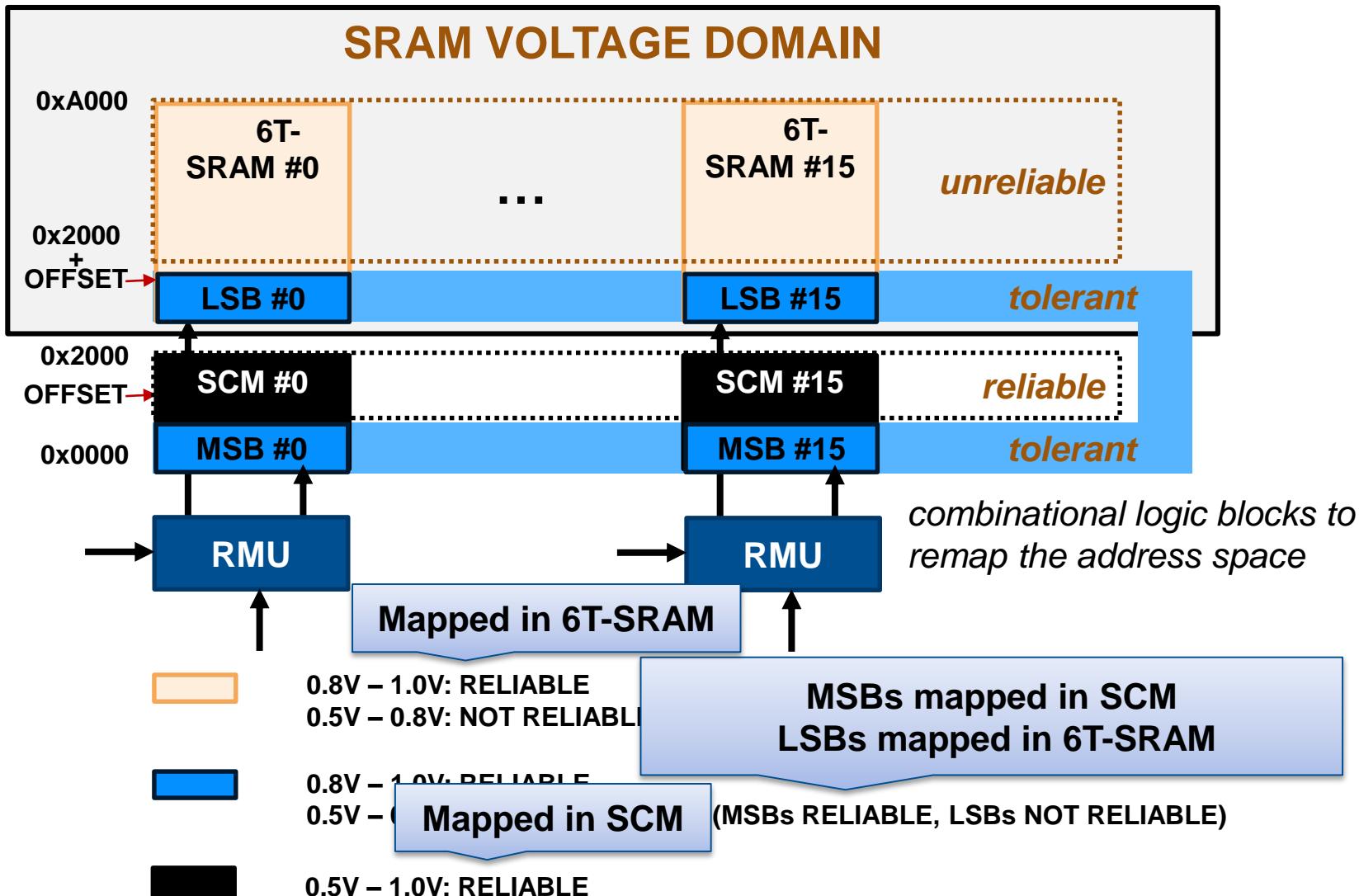
	Retention
SCM	0.25V
6T-SRAM	0.29V

- Probability of **flip-bit error** on a single bit during read/write operations

Voltage (V)	0.50	0.55	0.60	0.65	0.70	0.75	0.80
P(flip-bit) SCM	0.0	0.0	0.0	0.0	0.0	0.0	0.0
P(flip-bit) 6T	0.0037	0.0012	0.0003	5.24e-5	4.35e-6	4.16e-8	0.0

**6T-SRAM
32K**

**SCM
8K**



Programming model



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

- To conveniently control energy-efficient data mapping with an approximation model, we propose an ***extension to OpenMP***

- Two additional directives:
 - **#pragma tolerant var** → coupled to a *variable declaration* to specify it tolerates approximation
 - **#pragma tolerant computation** → coupled to a *program statement* to specify the memory addressed therein can be accessed in tolerant (low-voltage) mode

Example

```
int main()
{
    #pragma omp tolerant var
    int sparse_A[];
    int i=0, index;

    while(func(i))
    {
        ...
        index = compute_index(i);
        #pragma omp tolerant computation
        sparse_A[index] = compute_element(sparse_A, index);

        ...
        update(i);
    }
    use_value(sparse_A);
}
```

Marking index as tolerant would probably lead to fatal errors!

Tolerant computation here would probably lead to fatal errors!

Savings and Accuracy



- System-level Energy savings

	1 domain	2 domains		4 domains	
		Average	Best	Average	Best
Area (μm^2)	2943648	2968768		3062400	
Normalized energy	1.00	0.91	0.80	0.88	0.72
Normalized area	1.00	1.01		1.04	
Norm en. \times area	1.00	0.92	0.81	0.92	0.81

10-30%

- Accuracy «Max MSE» → for max end-to-end loss <1%

Benchmark	Zero-ing	Flip-bit	Max MSE
Color Tracking	676	64	225
HOG	2.12E+13	58564	2814663
CNN	17114769	26244	36864
Health	6867734	378	25000
Navi	1681	36	100

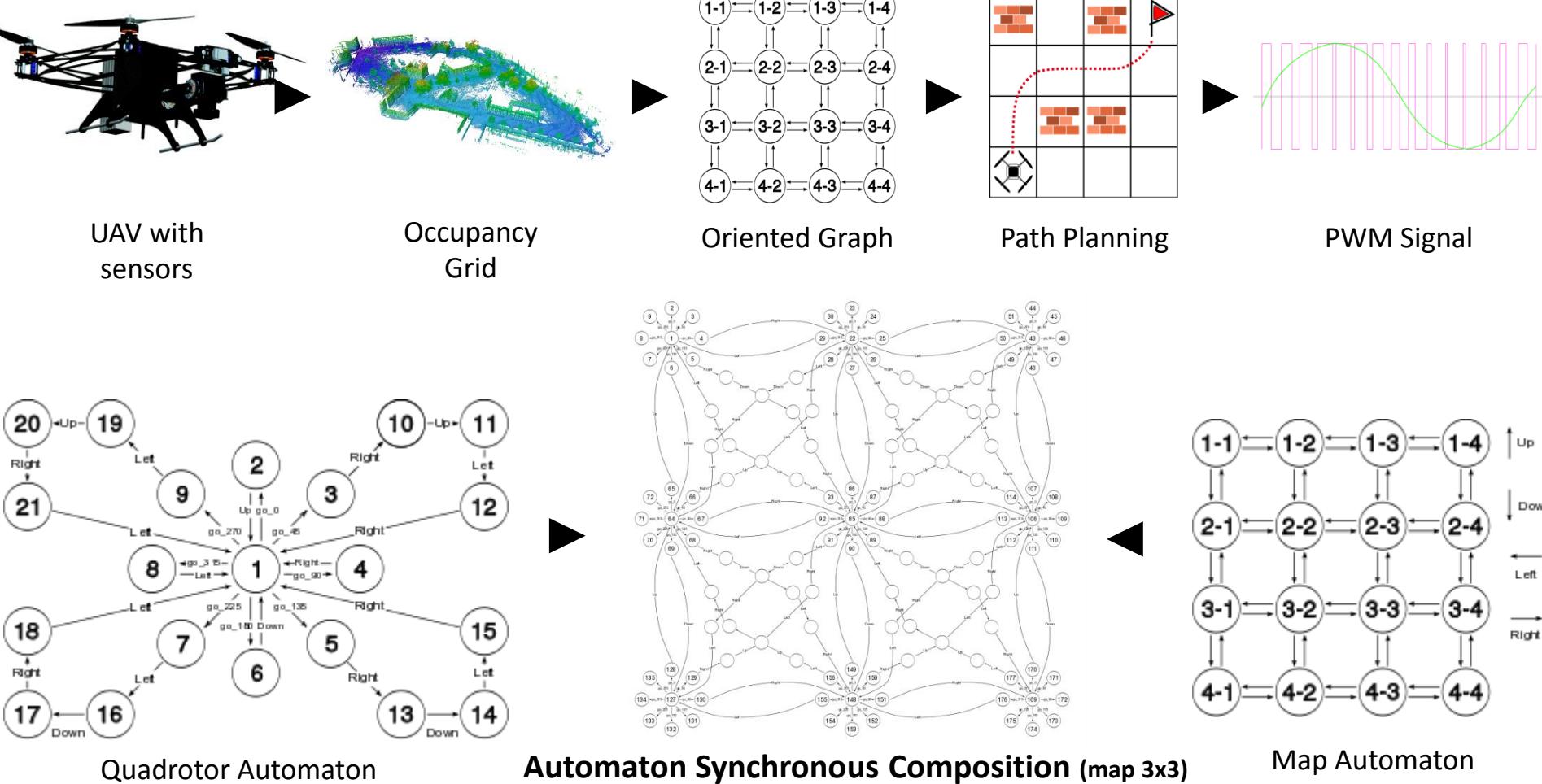
650x

Approximate ≠ Low Quality!!



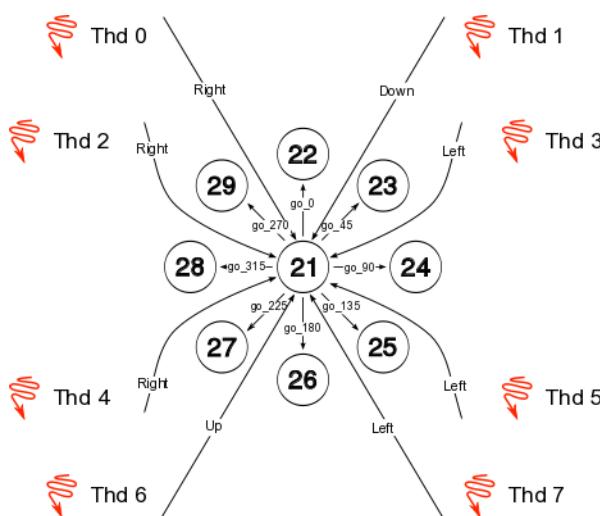
A new viewpoint on approximation

- A practical example: real-time UAV path planner



- Parallel Dijkstra algorithm
 - Each thread (in parallel)
 - handles several nodes “to be explored” (reference nodes)
 - for each reference node, explores all the neighbors
 - for each neighbor, eventually, updates both cost and predecessor id

Potential **race condition** → non-optimal and **non-deterministic** path



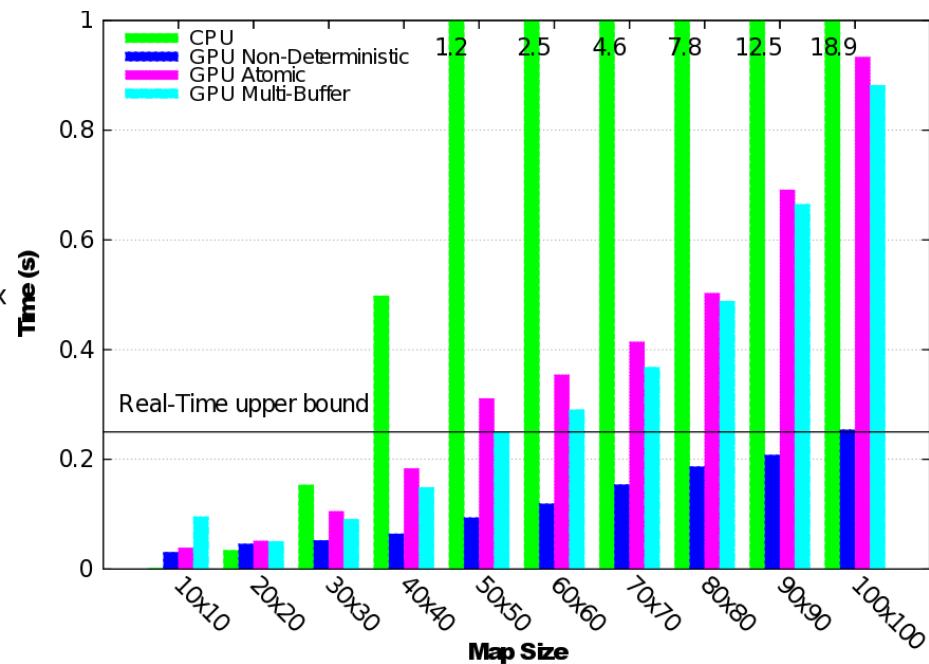
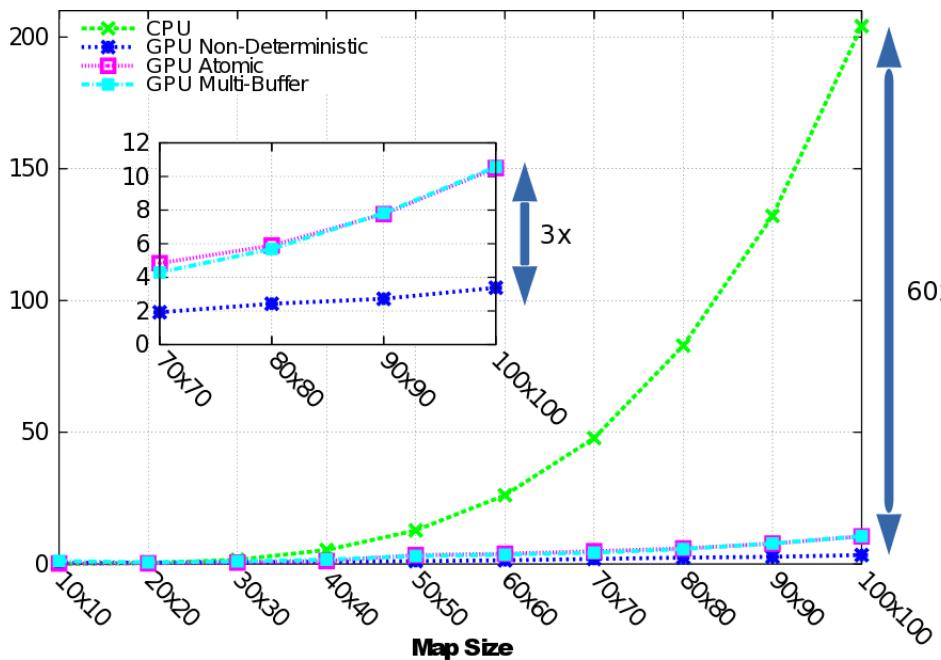
3 different approaches:

1. Data packing + keep non-determinism: *Non-Deterministic version*
2. Data packing + OpenCL atomic operation: *Atomic version*
3. Data packing + extra-size cost buffer, selecting minimal value in a second stage: *Multi-Buffer version*

Approximate is better!

- Sequential << Parallel << Parallel+Approximate

Of course energy is better, but also...

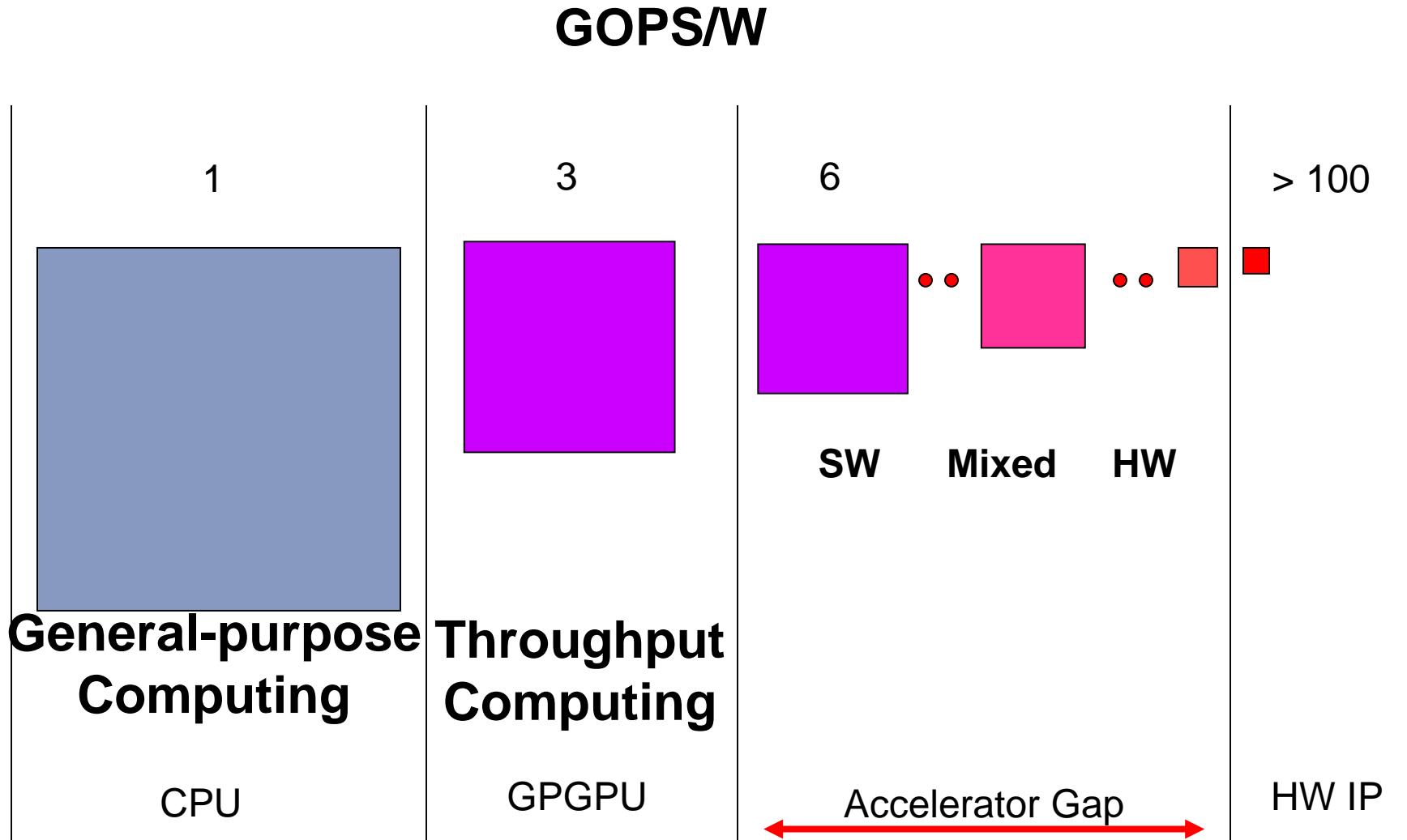


Plan quality is much better if the surveyed area (within the RT bound) is larger!

Pushing Beyond pJ/OP



Recovering more silicon efficiency



Closing The Accelerator Efficiency Gap with Agile Customization

Learn to Accelerate

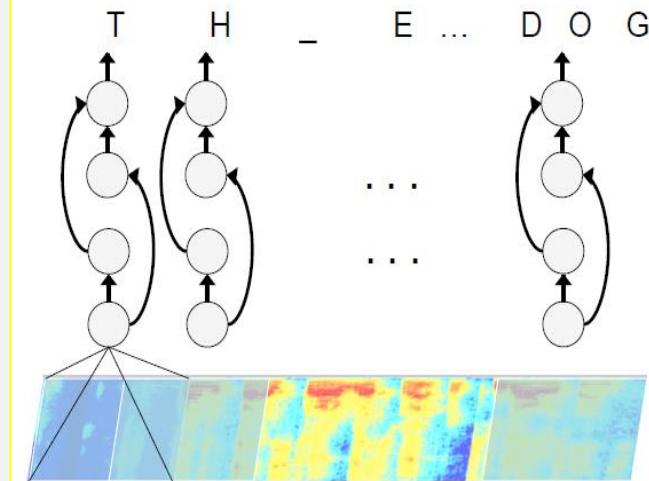
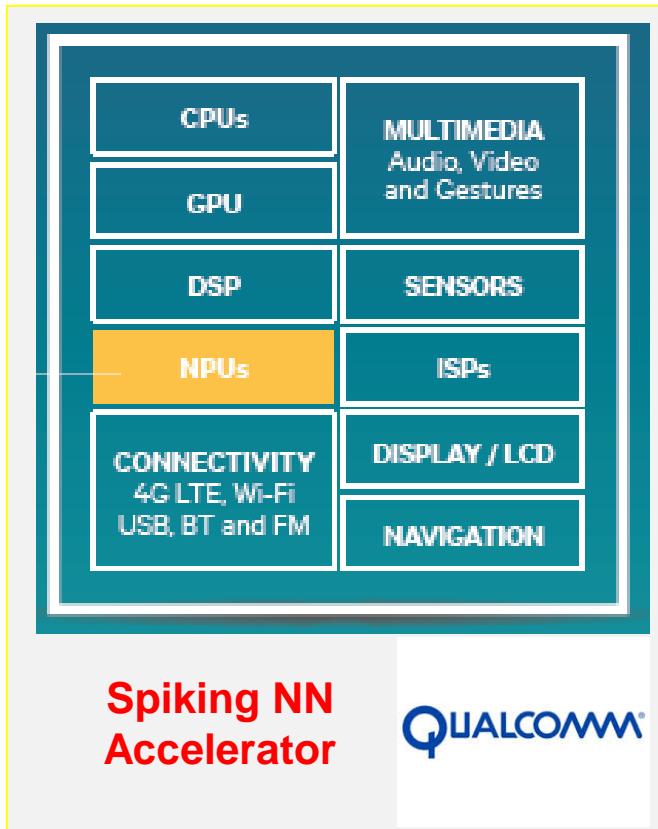


ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

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



Image recognition
[Russakovsky et al., 2014]



Speech recognition
[Hannun et al., 2014]

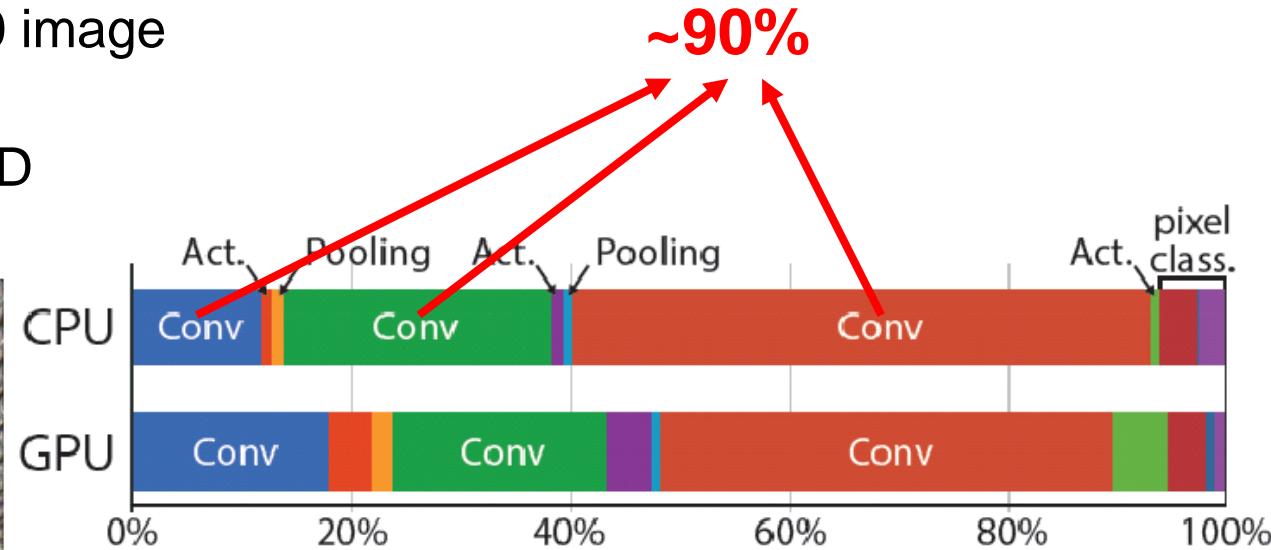
- **Flexible** acceleration: learned CNN weights are “the program”

Computational Effort

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

■ Computational effort

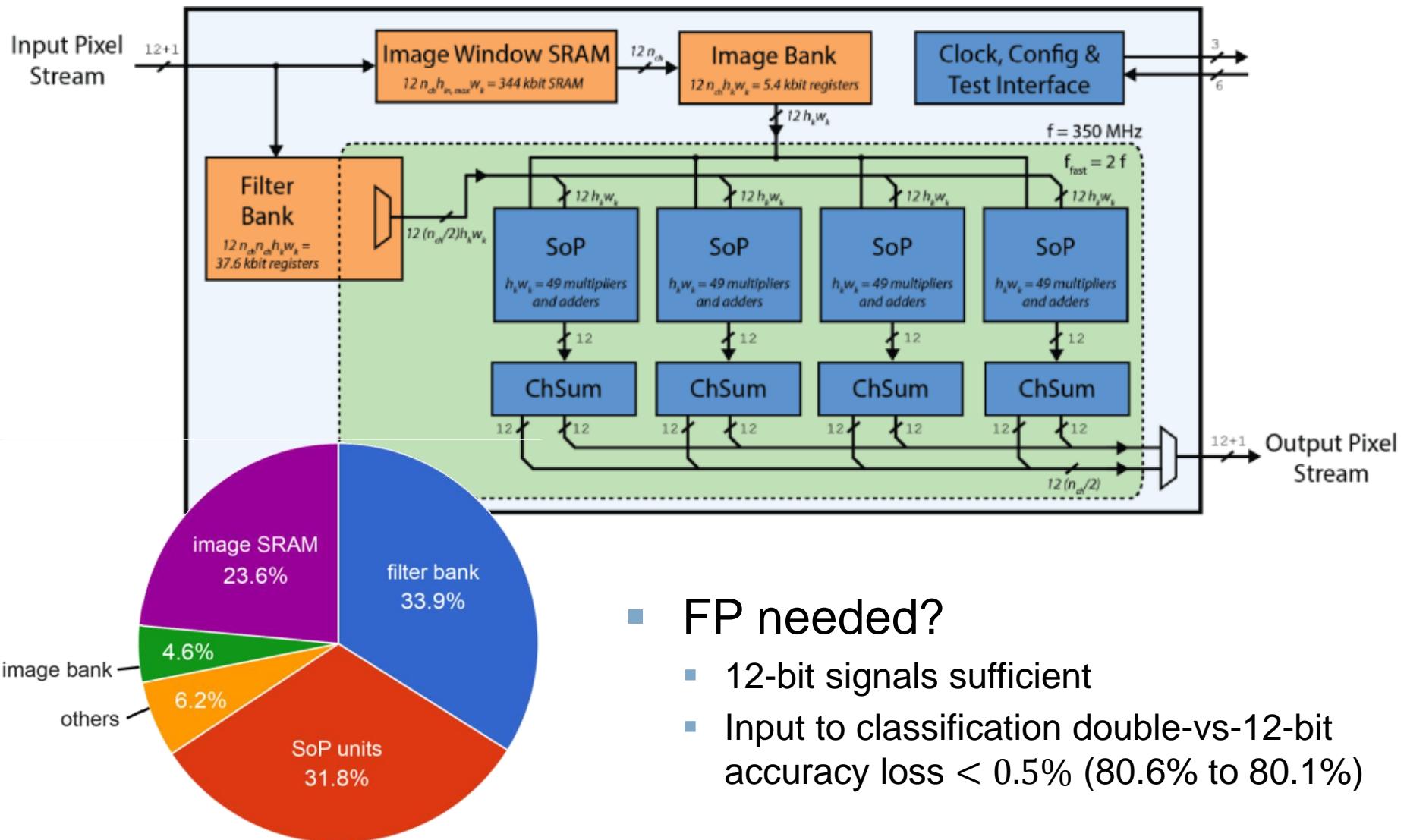
- 7.5 GOp for 320x240 image
- 260 GOp for FHD
- 1050 GOp for 4k UHD



Origami chip

Origami: A CNN Accelerator

ALMA MATER STUDIORUM
 UNIVERSITÀ DI BOLOGNA



Sub pJ/OP?

437 GOPS/W @1.2V

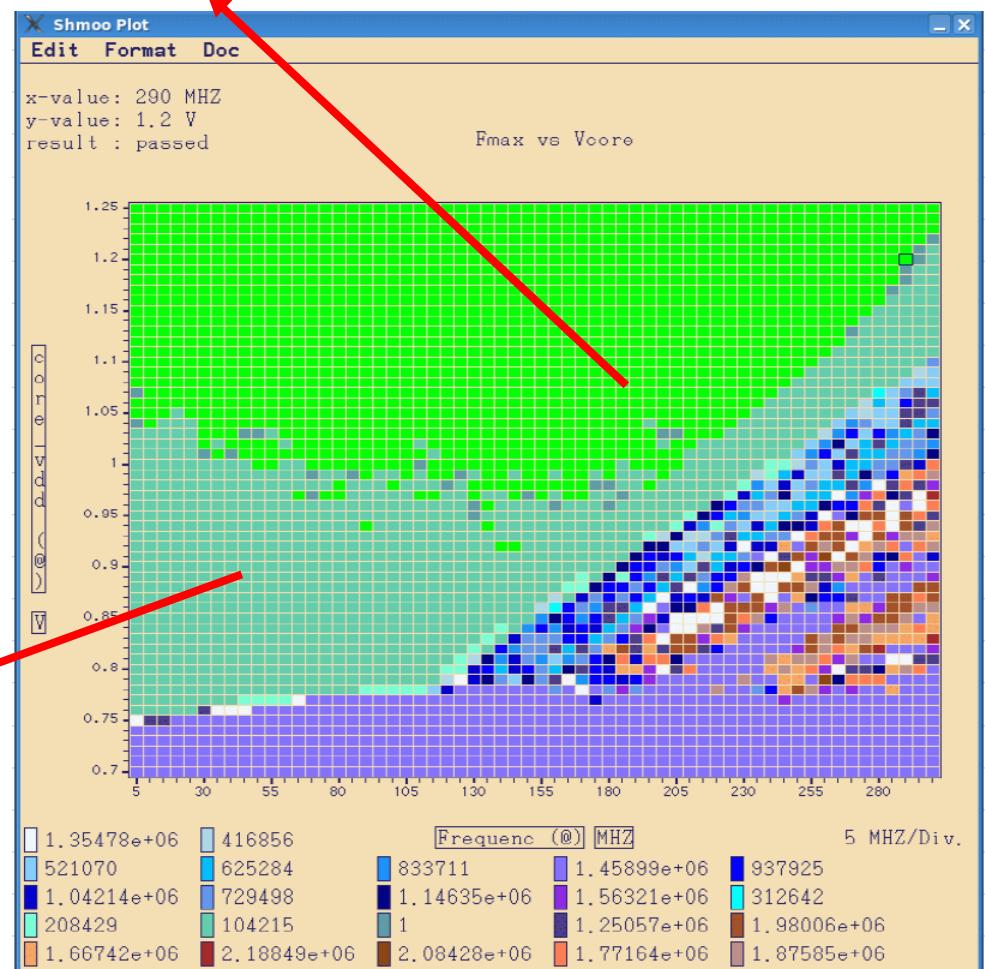
803 GOPS/W @0.8V

0% bit flips

1% bit flips
1.84x energy improvement



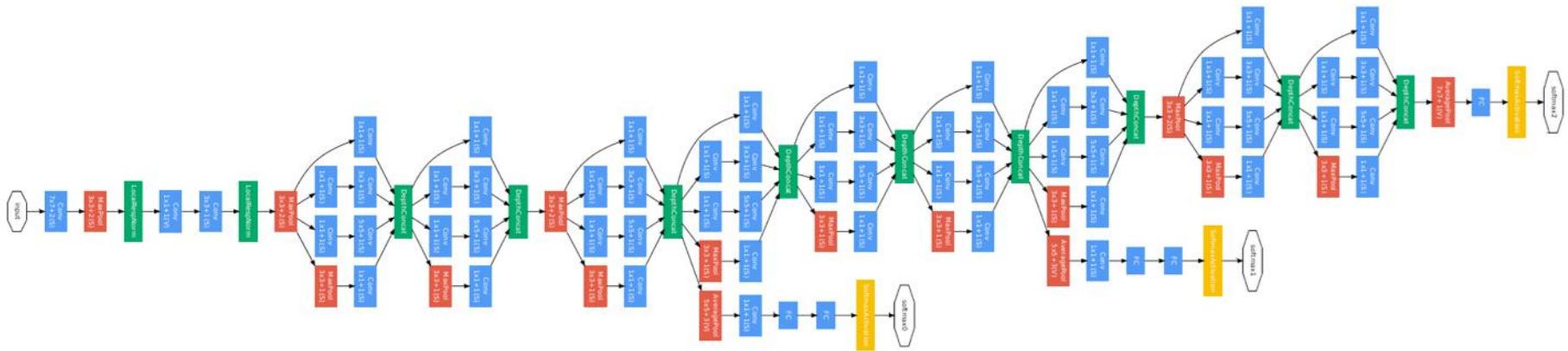
1.2pJ/OP



CNNs: typical workload



ALMA MATER STUDIORUM
 UNIVERSITÀ DI BOLOGNA



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

$\sim 7 \times 10^6$ parameters

$\sim 2.3 \times 10^9$ 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

$$w_{b,stoch} = \begin{cases} -1 & p_{-1} = \sigma(w) \\ 1 & p_1 = 1 - p_{-1} \end{cases} \quad w_{b,det} = \begin{cases} -1 & w < 0 \\ 1 & w > 0 \end{cases}$$

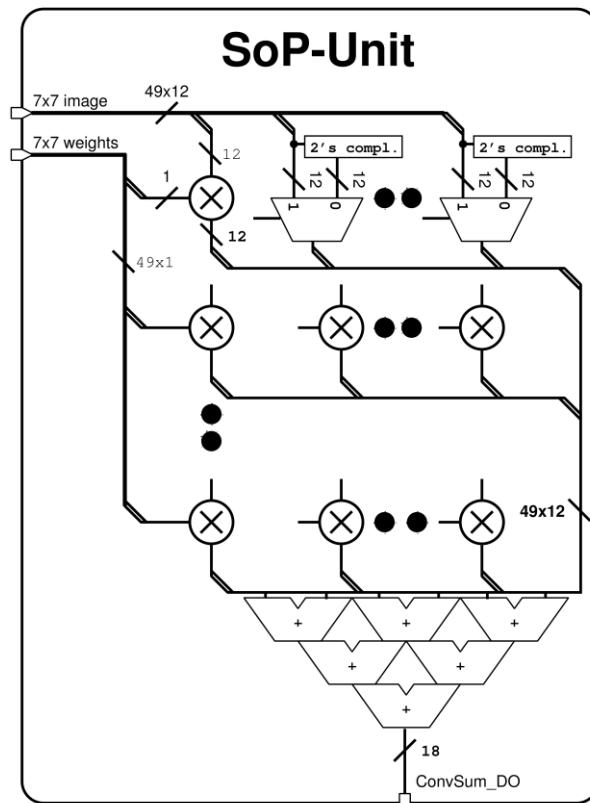
- Learning large networks is still an issue with binary connect...
- Ultra-optimized HW is possible!
 - Power reduction because of arithmetic simplification
 - Major arithmetic density improvements
 - Area can be used for more energy-efficient weight storage
 - SCM memories for lower voltage → E goes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

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

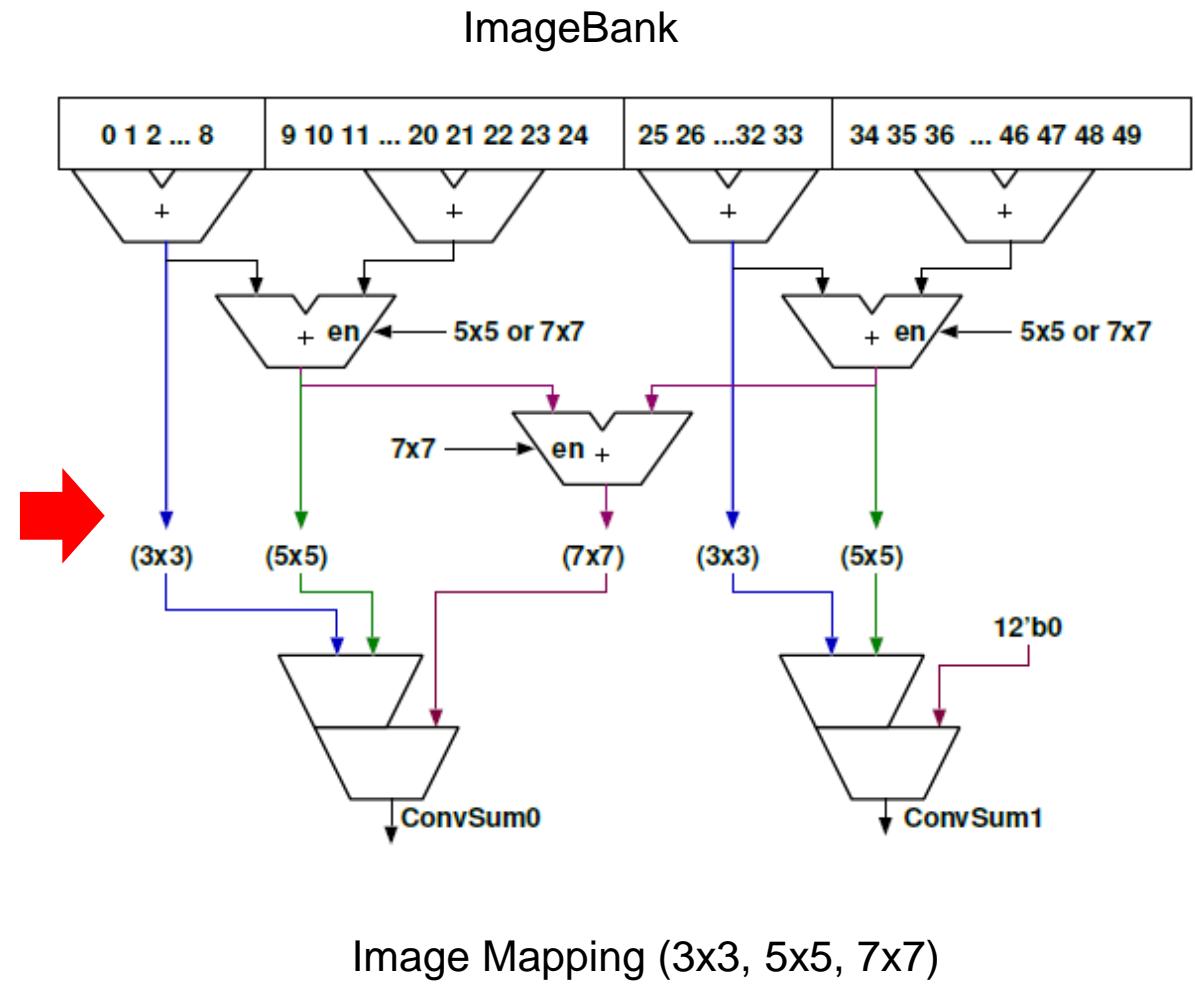
SoP-Unit Optimization



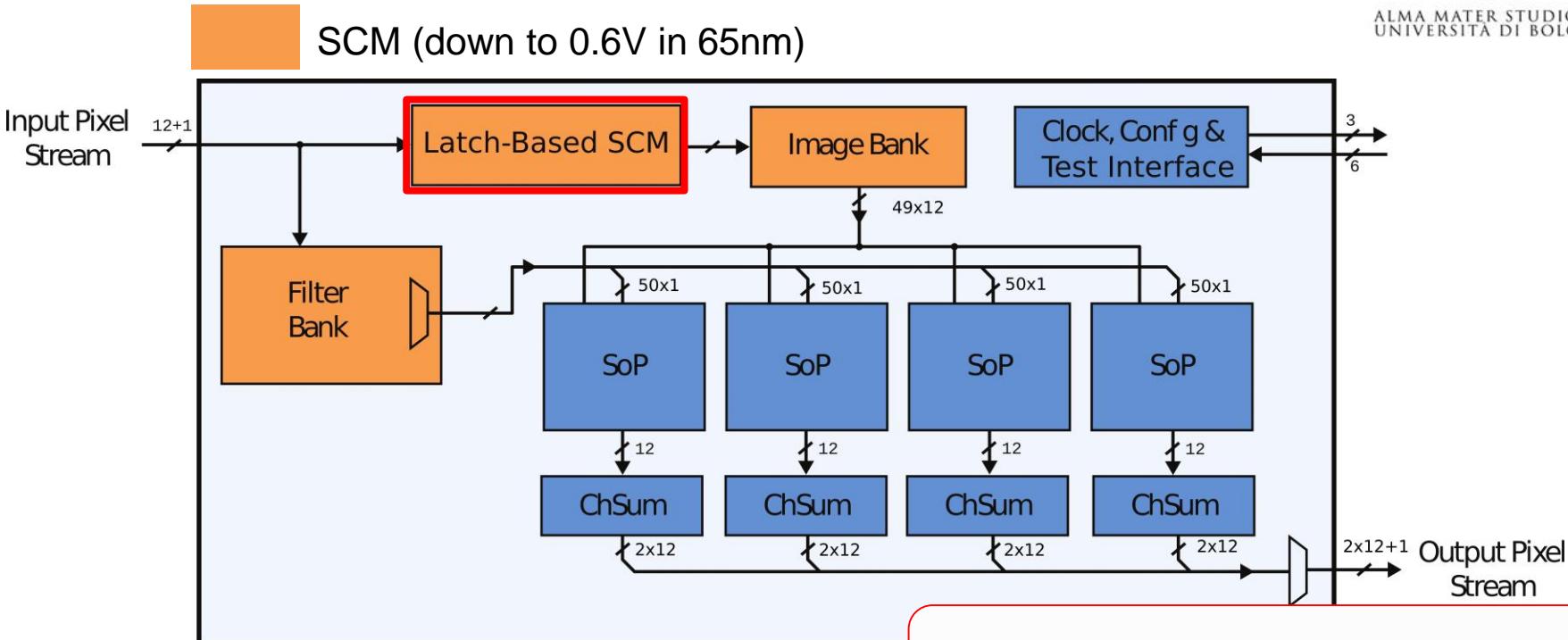
ALMA MATER STUDIORUM
 UNIVERSITÀ DI BOLOGNA



Equivalent for 7x7 SoP

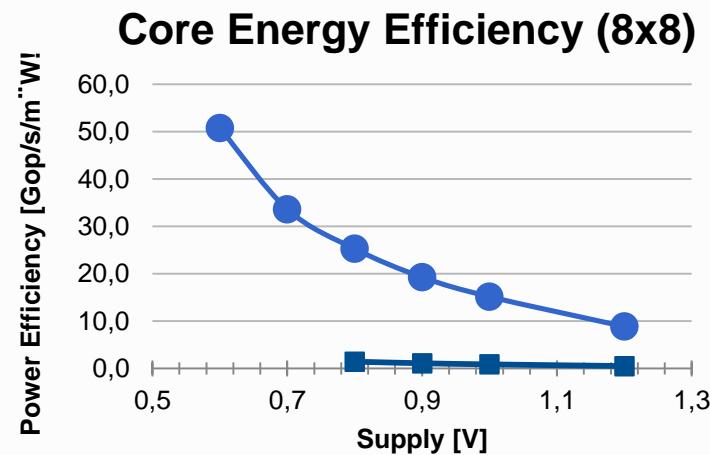


SCM for Energy efficiency



Same area 8 → 32 SoP units + all SCM

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

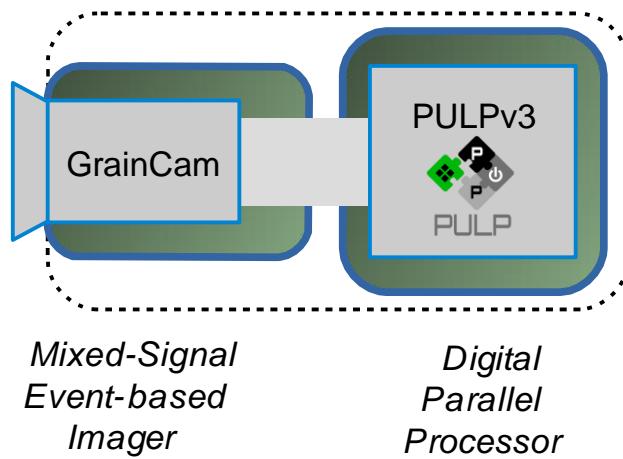


Back to System-Level



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Smart Visual Sensor → idle most of the time (nothing interesting to see)



- **Event-Driven Computation**, which occurs only when relevant events are detected by the sensor
- **Event-based sensor interface** to minimize IO energy (vs. Frame-based interface)
- **Mixed-signal event triggering** with an ULP imager with internal processing AMS capability

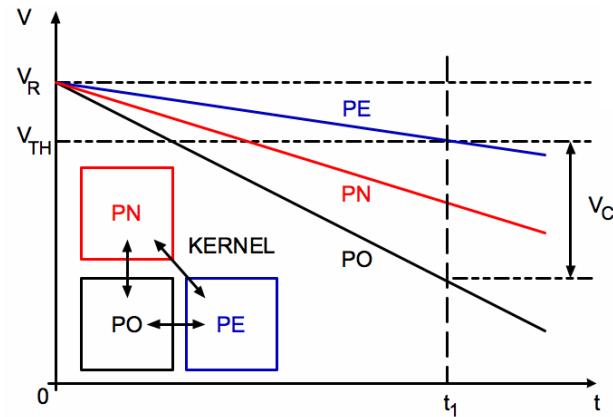
Doing *nothing* well is also very important!

GrainCam Imager (FBK)

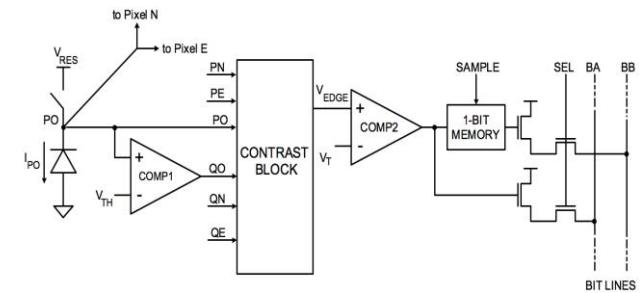
ALMA MATER STUDIORUM
 UNIVERSITÀ DI BOLOGNA



Pixel-level spatial-contrast extraction



$$V_C = V_{PE}(t_1) - V_{PO}(t_1) = (V_R - V_{TH}) \left(\frac{I_{PO} - I_{PE}}{I_{PE}} \right)$$



Analog internal image processing

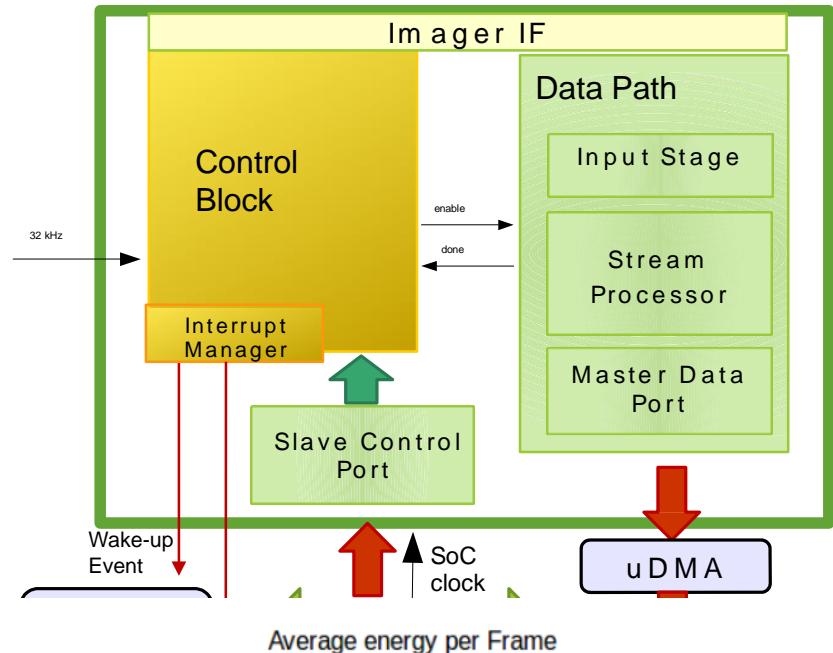
- Contrast Extraction
- Motion Extraction, differencing two successive frames
- Background Subtraction, differencing the reference image, stored in the memory, with the current frame



Event-based + Prefiltering + Buffering

Control Block:

- Always-on 32KHz controller for continuously handling sensor timing signals (e.g. integration, readout periods) and control signals (e.g. readout mode, reference frame sampling)
- Activates data-path block during imager readout phase
- Triggers events

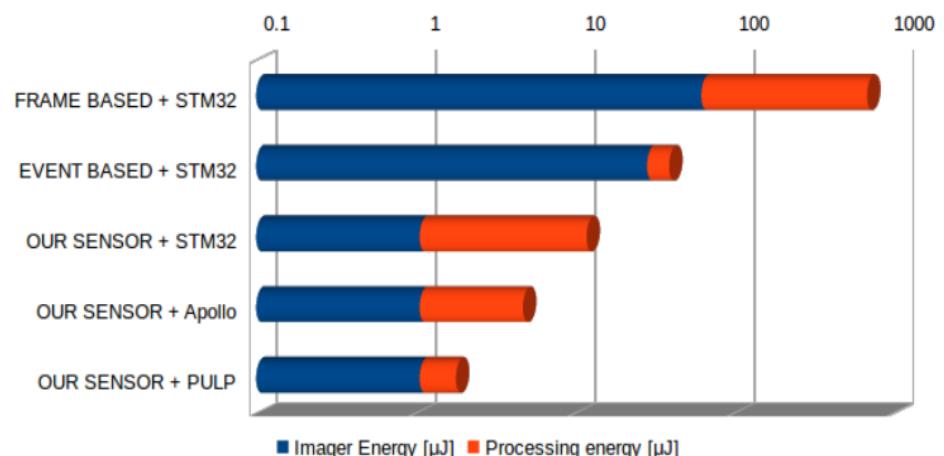


Data Path: stream processing on sensed post-triggering data

u-DMA: autonomously transfers pixel-events to L2 memory during imager readout, without requiring cluster action

Sensor IF Peripheral as transfer initiator.

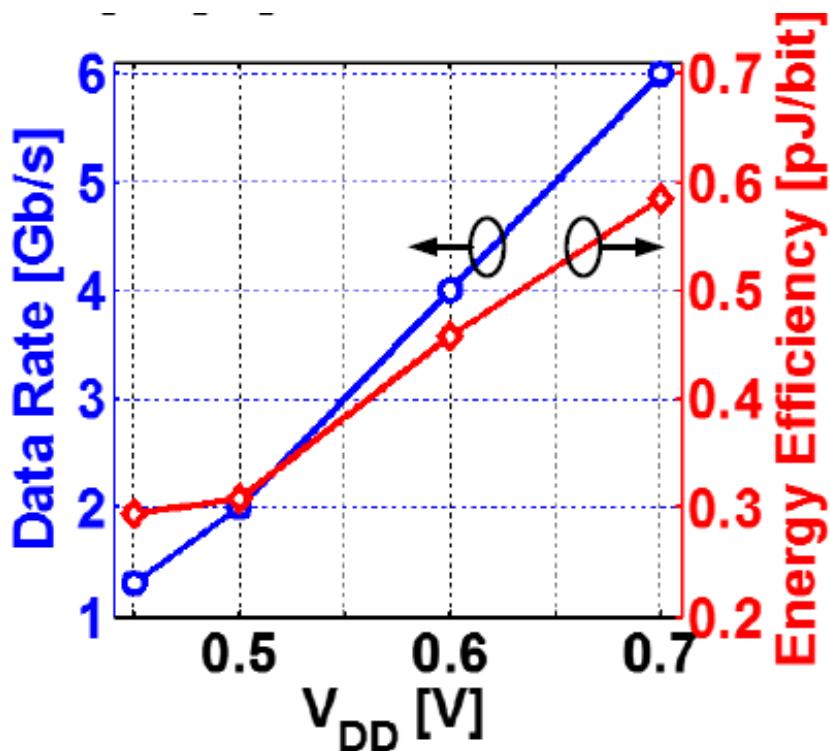
Slave Port interface to set-up sensor interface parameters



Integrating & communicating?

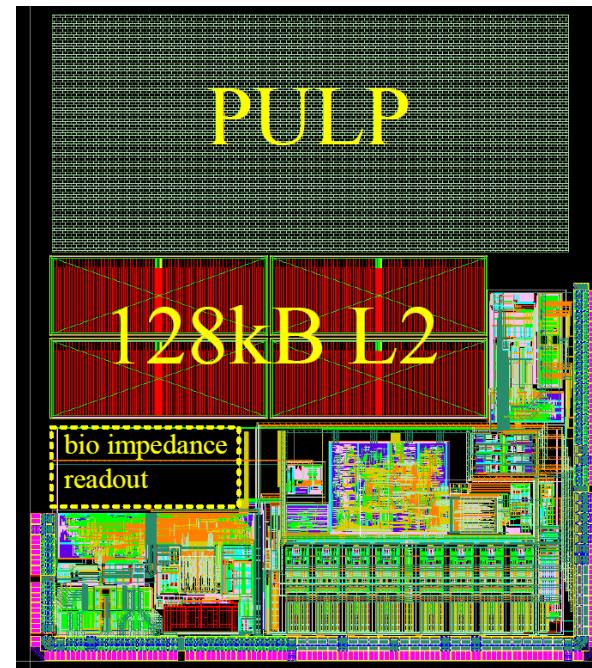
ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Sub pJ/b Chip2chip link



Source-synchronous, pseudo-differential, unterminated, Voltage Mode, 200mVpp, 1/8 rate CLK, self-calibrating PLL-based phase generator [ISSCC15]

VivoSoC2 – full AFE integration



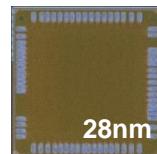
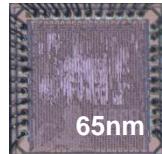
Multiple Biomedical Interfaces:

8 ExG channels with lead-off detection, Photoplethysmography (PPG), Bio-impedance monitoring, temperature sensing, 6 channels nerve blocker and stimulator, standard digital interfaces

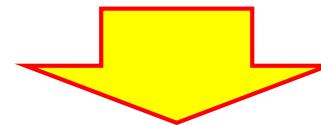
Open Source Parallel ULP computing for the IoT

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

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



PULP
Parallel Ultra Low Power



Processor &
Hardware IPs

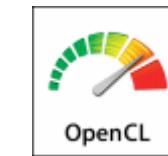


Compiler
Infrastructure

MicroPython
Python for microcontrollers



OpenMP



OpenVX
Programming
Model



PULP
 Parallel Ultra Low Power

[https://github.co](https://github.com)

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!

[About](#)
[Release Plan](#)
[Resources](#)
[Download](#)

EE Times

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](#)

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).

Conclusions



- IoT a Challenge and an opportunity
 - Computing Everywhere!
- Energy efficiency requirements: pj/OP and below
 - Technology scaling alone is not doing the job for us
 - Ultra-low power architecture and circuits are needed, but not sufficient in the long run
 - Most promising technologies: 3D integration, low-leakage, non-volatile silicon-compatible mem-computing devices, but adoption rate will be **SLOW**
- Approximate computing to the rescue
 - Holistic approach needed, math, number systems, algorithms, tools, runtime, architecture, circuits, devices → **from approximate to transprecision computing**
- Non-Von-Neumann + Transprecision... **Jackpot potential!**
- **Open Source HW & SW approach for building an innovation ecosystem**

Thanks for your attention!!!



www.pulp-platform.org

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

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

