

Parallel digital signal processing in a mW power envelope: how and why?

Davide Rossi
DEI-UNIBO



Multithermand AdG
Multiscale Thermal Management of Computing Systems

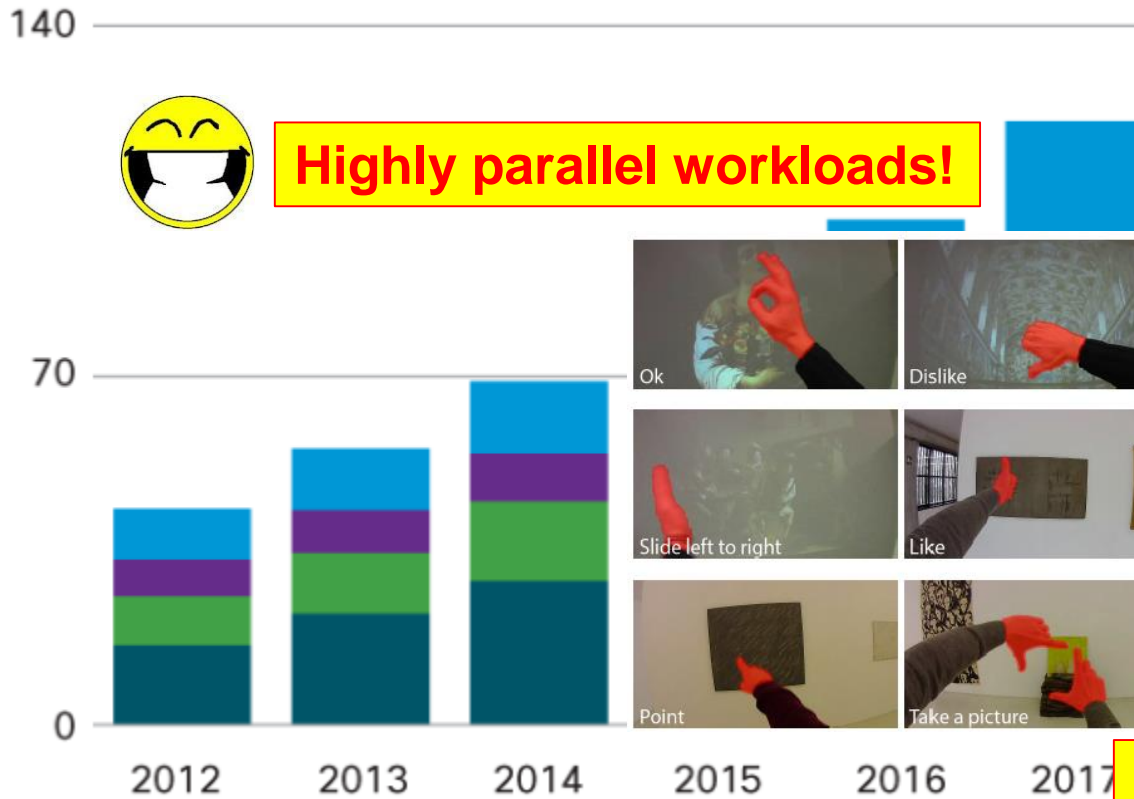


IOT or Data Deluge?

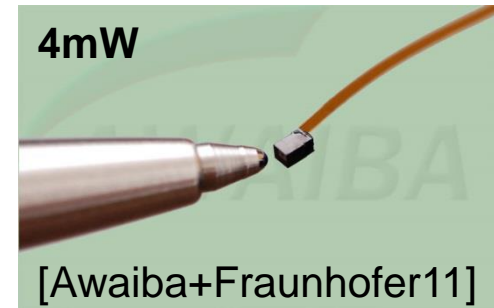


23% CAGR 2012-2017

Exabytes per Month



- Web/Data (24.2%, 18.9%)
- File Sharing (15.7%, 8.1%)
- Managed IP Video (21.8%, 21.0%)
- Internet Video (38.3%, 52.0%)



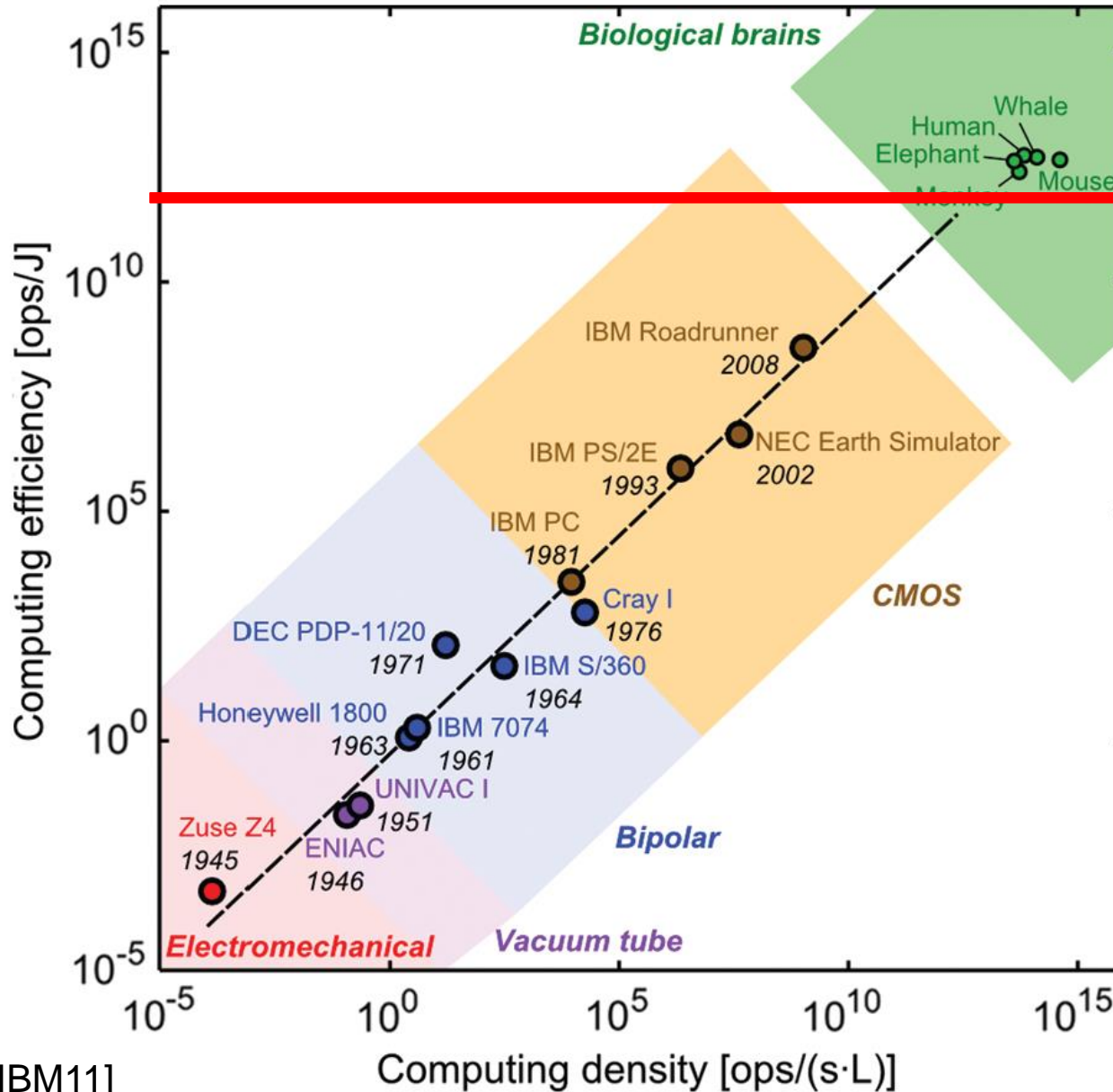
CV is the energy bottleneck

Source: Cisco VNI, 2013

The percentages within parenthesis next to the legend denote the relative traffic shares in 2012 and 2017.

**In-situ stream processing & fusion + visual intelligence is a must!
- In a few mW power envelope!!**

How efficient?



10¹²ops/J



1pJ/op



1GOPS/mW

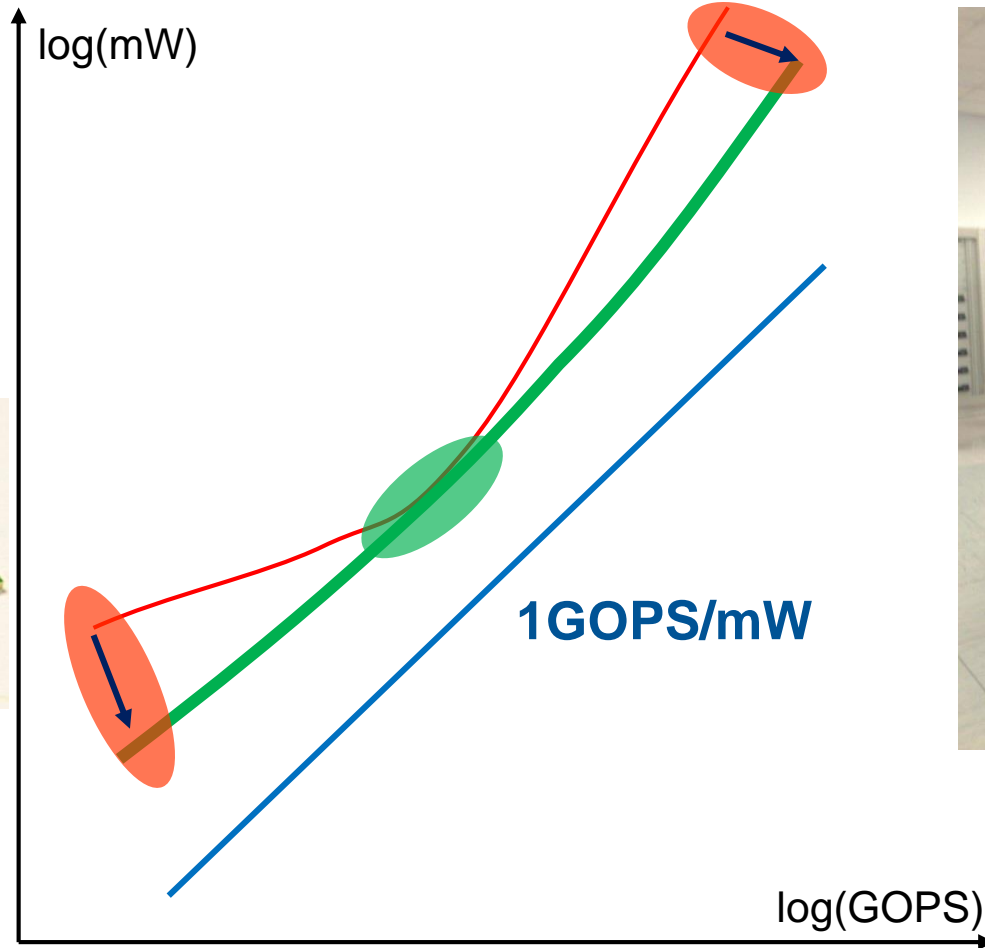
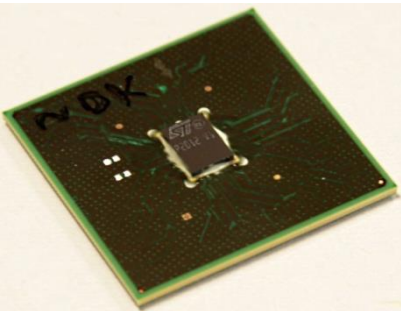
[RuchIBM11]

The challenge of Energy Proportionality



0,003GOPS/mW – 30KW

0,04GOPS/mW



From KOPS (10^3) to EOPS (10^{18})!

A Very Short Review on CMOS power and power minimization



Summary of Power Dissipation Sources

$$P \sim \alpha \cdot (C_L + C_{CS}) \cdot V_{swing} \cdot V_{DD} \cdot f + (I_{DC} + I_{Leak}) \cdot V_{DD}$$

- α – switching activity
- C_L – load capacitance
- C_{CS} – short-circuit capacitance
- V_{swing} – voltage swing
- f – frequency
- I_{DC} – static current
- I_{leak} – leakage current

$$P = \frac{\text{energy}}{\text{operation}} \times \text{rate} + \text{static power}$$

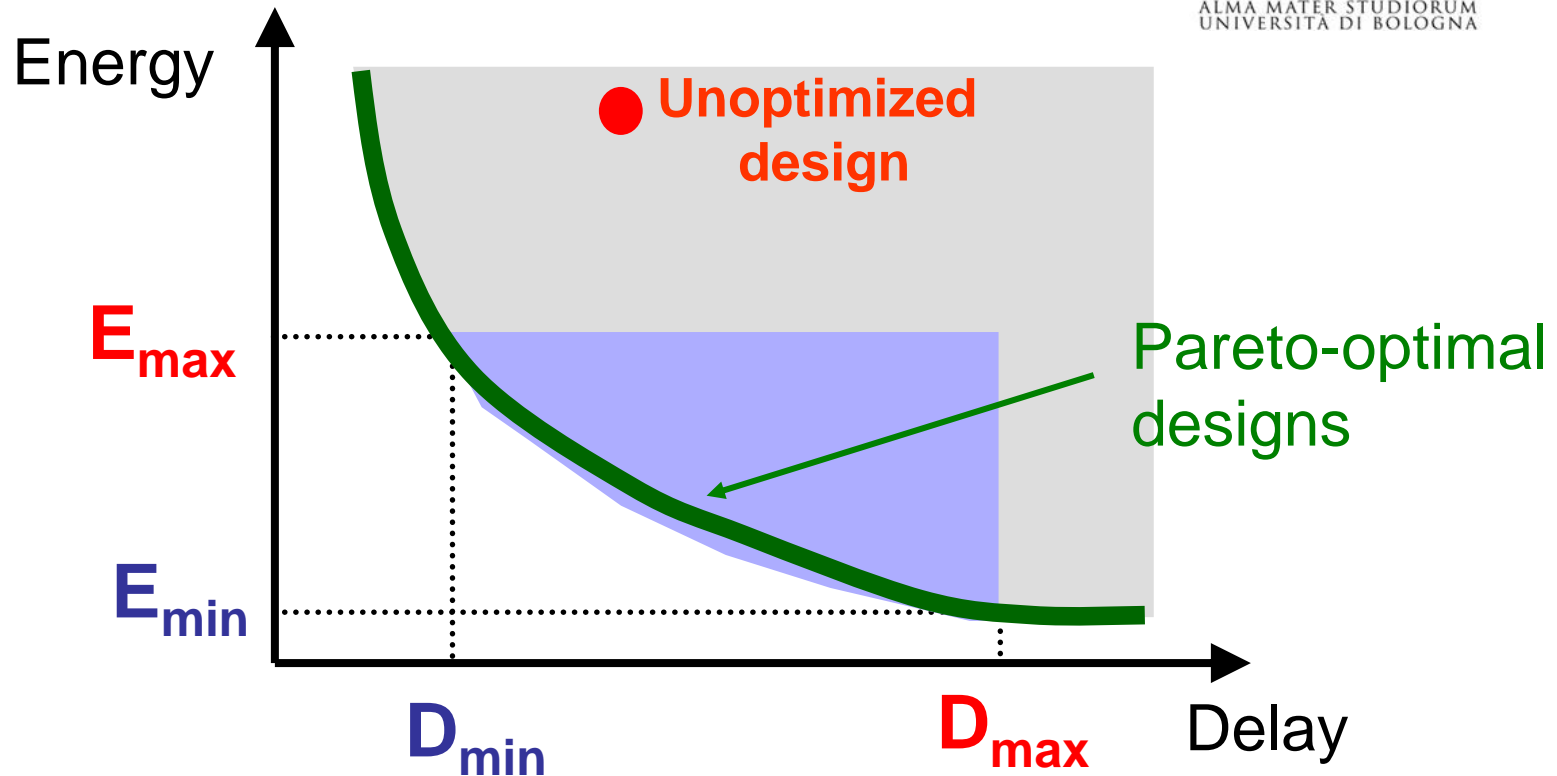
The Traditional Design Philosophy

- Maximum performance is primary goal
 - Minimum delay at circuit level
- Architecture implements the required function with target throughput, latency
- Performance achieved through optimum sizing, logic mapping, architectural transformations.
- Supplies, thresholds set to achieve maximum performance, subject to reliability constraints

The New Design Philosophy

- Maximum performance (in terms of propagation delay) is too power-hungry, and/or not even practically achievable
- Many (if not most) applications either can tolerate larger latency, or can live with lower than maximum clock-speeds
- Excess performance (as offered by technology) to be used for energy/power reduction

Trading off speed for power



In energy-constrained world, design is trade-off process

- ◆ Minimize energy for a given performance requirement
- ◆ Maximize performance for given energy budget

Reducing power @ all design levels

- ◆ Algorithmic level
- ◆ Compiler level
- ◆ Architecture level
- ◆ Organization level
- ◆ Circuit level
- ◆ Silicon level

- ◆ Important concepts:
 - Lower Vdd and freq. (even if errors occur) / dynamically adapt Vdd and freq.
 - Reduce circuit
 - Exploit locality
 - Reduce switching activity, glitches, etc.



Algorithmic level

- ◆ The best indicator for energy is
.... **the number of cycles**

- ◆ Try alternative algorithms with lower complexity
 - E.g. quick-sort, $O(n \log n) \Leftrightarrow$ bubble-sort, $O(n^2)$
 - ... but be aware of the 'constant' : $O(n \log n) \Rightarrow c^*(n \log n)$

- ◆ Heuristic approach
 - Go for a good solution, not the best !!

Biggest gains at this level !!

Compiler level

- ◆ Source-to-Source transformations
 - loop trafo's to improve locality
- ◆ Strength reduction
 - E.g. replace $\text{Const} * A$ with Add's and Shift's
 - Replace Floating point with Fixed point
- ◆ Reduce register pressure / number of accesses to register file
 - Use software bypassing
- ◆ Scenarios: current workloads are highly dynamic
 - Determine and predict execution modes
 - Group execution modes into scenarios
 - Perform special optimizations per scenario
 - DFVS: Dynamic Voltage and Frequency Scaling
 - More advanced loop optimizations
- ◆ Reorder instructions to reduce bit-transitions

Architecture level

- ◆ Going parallel
- ◆ Going heterogeneous
 - tune your architecture, exploit SFUs (special function units)
 - trade-off between flexibility / programmability / genericity and efficiency
- ◆ Add local memories
 - prefer scratchpad i.s.o. cache
- ◆ Cluster FUs and register files (see next slide)
- ◆ Reduce bit-width
 - sub-word parallelism (SIMD)

Organization (micro-arch.) level

◆ Enabling V_{dd} reduction

- Pipelining
 - cheap way of parallelism
- Enabling lower freq. \Rightarrow lower V_{dd}
- Note 1: don't pipeline if you don't need the performance
- Note 2: don't exaggerate (like the 31-stage Pentium 4)

◆ Reduce register traffic

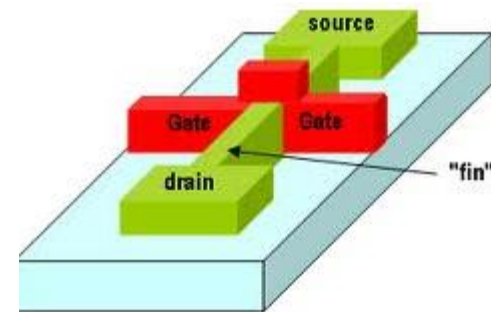
- avoid unnecessary reads and write
- make bypass registers visible

Circuit level

- ◆ Clock gating
- ◆ Power gating
- ◆ Multiple Vdd modes
- ◆ Reduce glitches: balancing digital path's
- ◆ Exploit Zeros
- ◆ Special SRAM cells
 - normal SRAM can not scale below $V_{dd} = 0.7 - 0.8$ Volt
- ◆ Razor method; replay
- ◆ Allow errors and add redundancy to architectural invisible structures
 - branch predictor
 - caches
- ◆ .. and many more ..

Silicon level

- ◆ Higher V_t ($V_{\text{threshold}}$)
- ◆ Back Biasing control
 - see thesis Maurice Meijer (2011)
- ◆ SOI (Silicon on Insulator)
 - silicon junction is above an electr. insulator (silicon dioxide)
 - lowers parasitic device capacitance
- ◆ Better transistors: Finfet
 - multi-gate
 - reduce leakage (off-state current)
- ◆ .. and many more



Two Ideas to Remember

...with their caveats



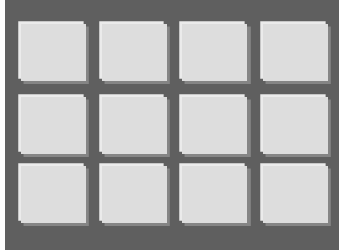
Go Simple+Parallel



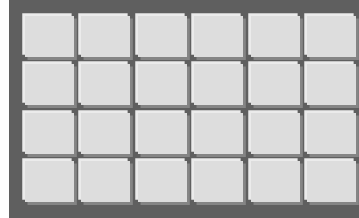
ALMA MATER STUDIORUM
UNIVERSITA DI BOLOGNA

13mm, **100W**, 48MB Cache, 4B Transistors, in 22nm

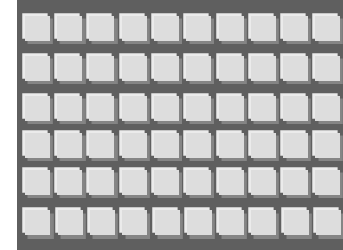
12 Cores



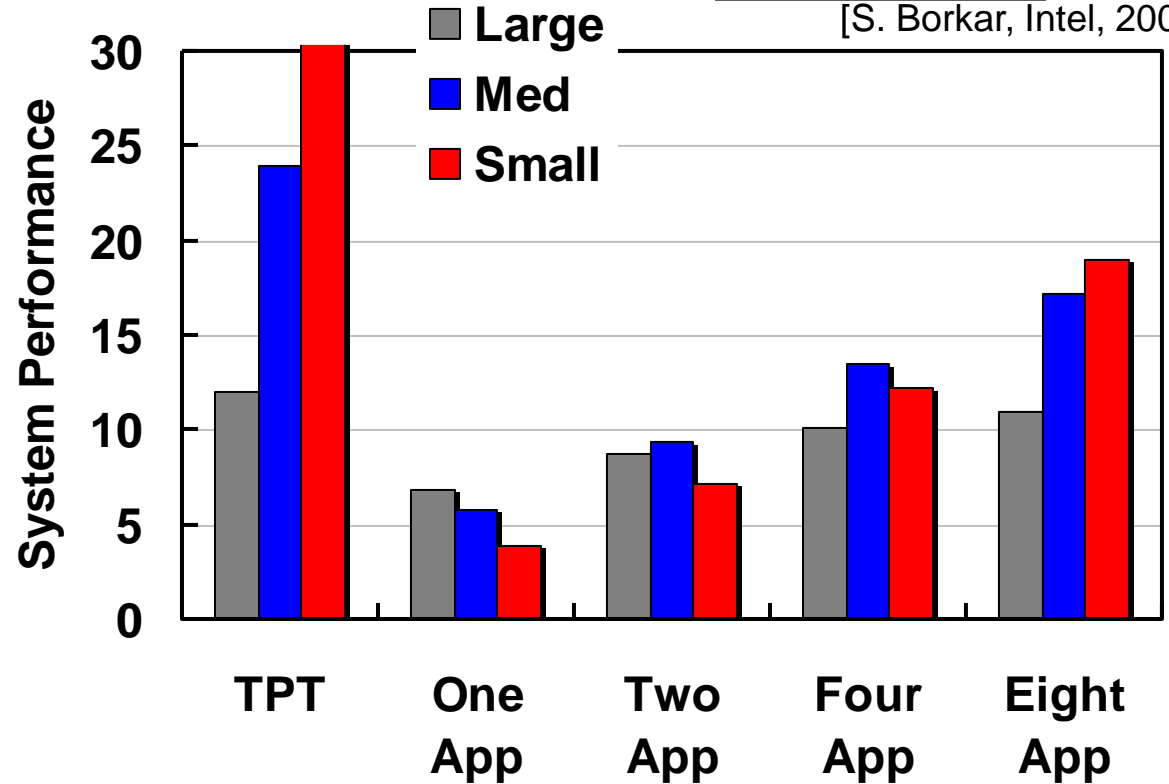
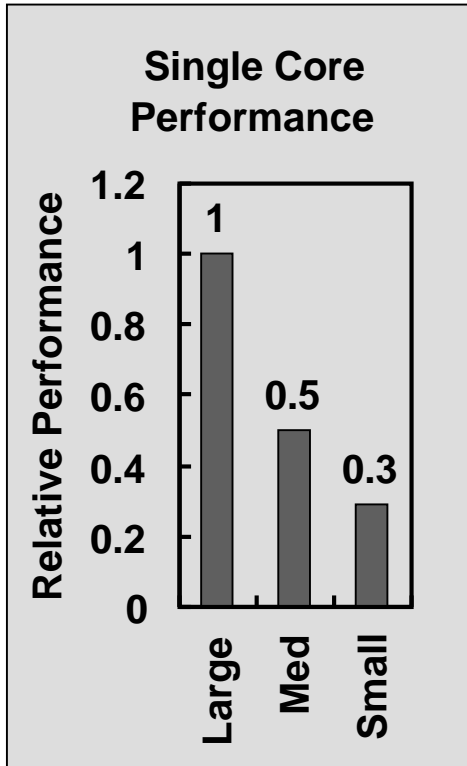
48 Cores



144 Cores



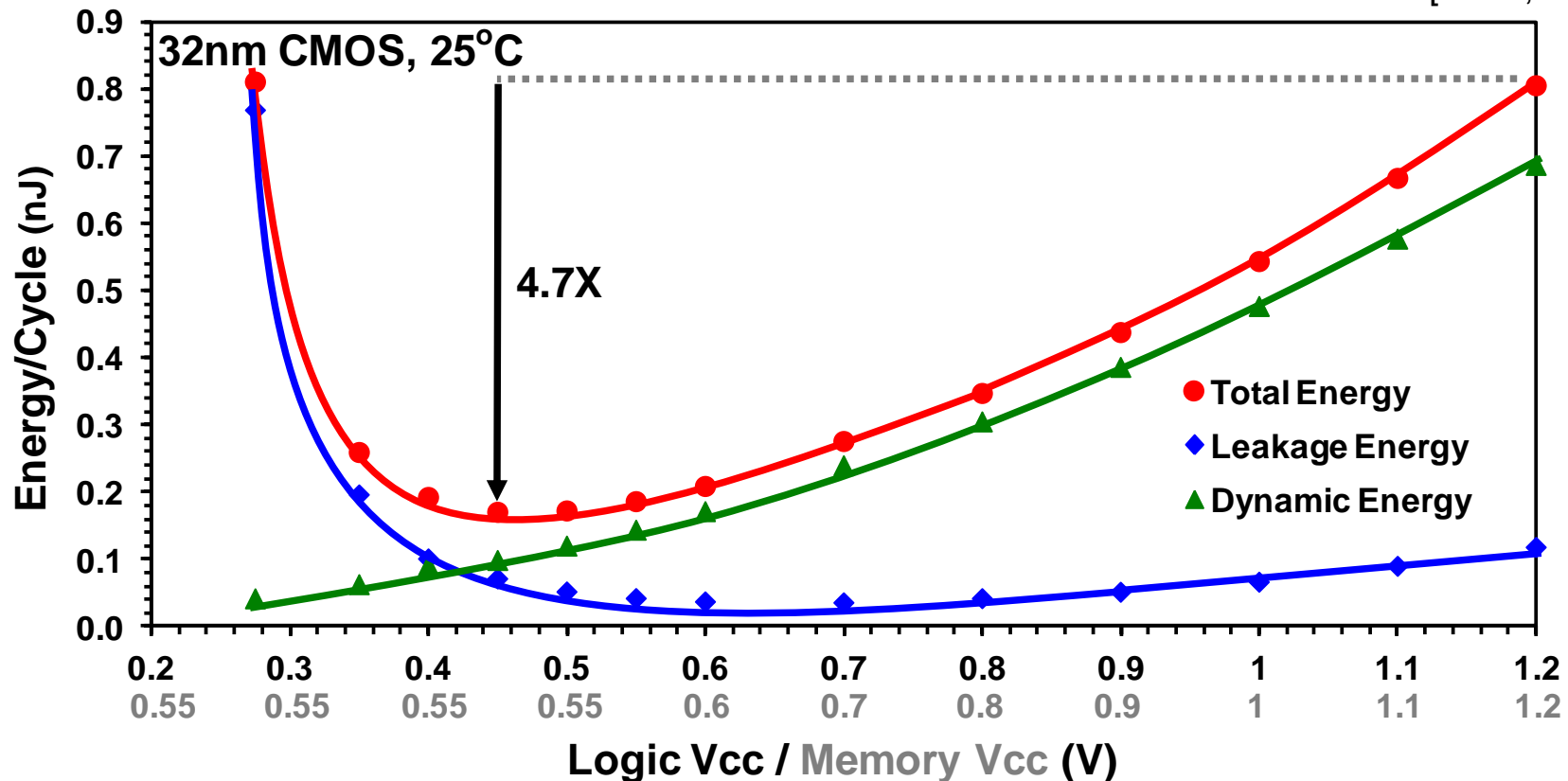
[S. Borkar, Intel, 2006]



Lower Voltage Supply



[V. De, Intel, 2013]



...but be careful with Leakage and its variability!

Introducing PULP

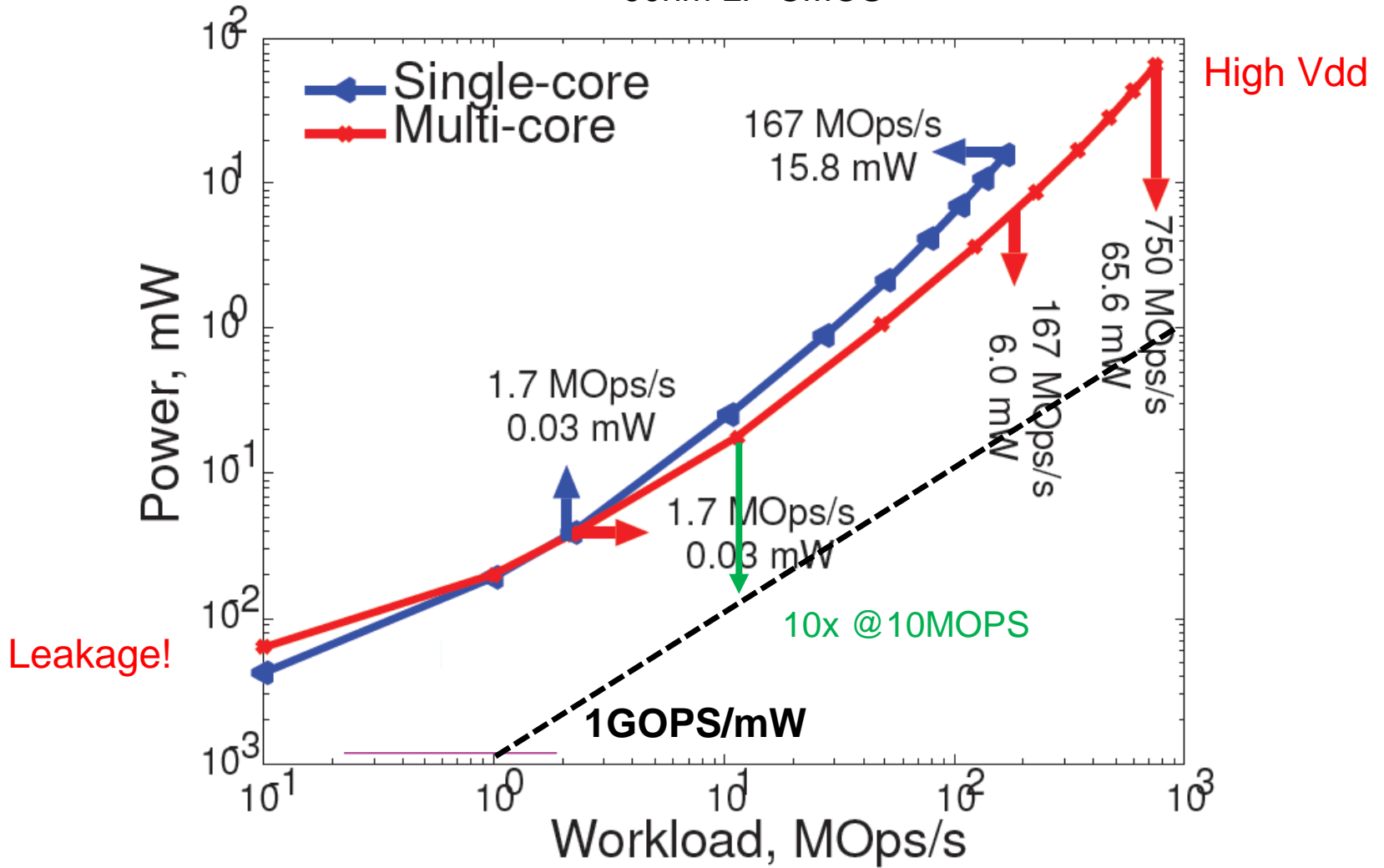


NTC Multicore?



2010 - With EPFL (Atienza, Burg)

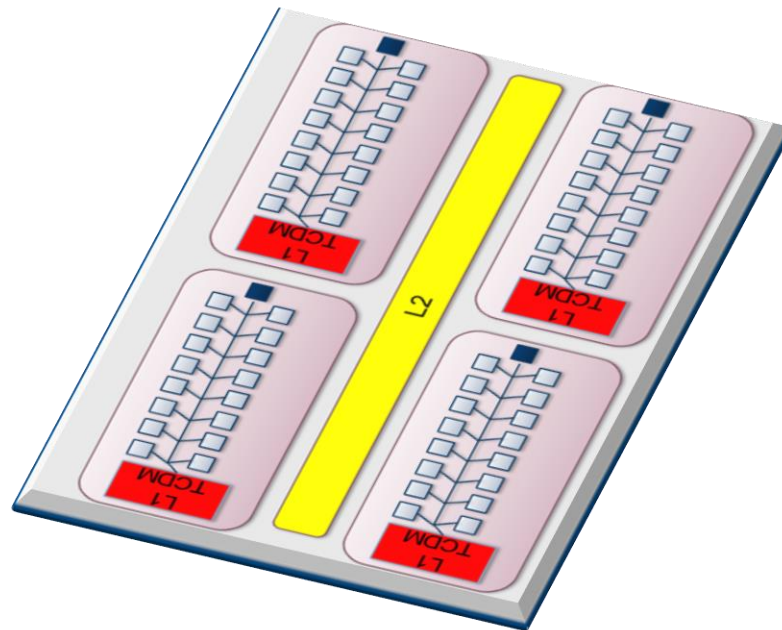
90nm LP-CMOS



A pJ/OP Parallel ULP Computing Platform



- **pJ/OP** is traditionally* the target of ASIC + uCntr
 - ➔ Scalable: [KOPS, TOPS], 32bit architecture
 - ➔ Flexible: OpenMP, OpenVX
 - ➔ Open: Software & **HW**



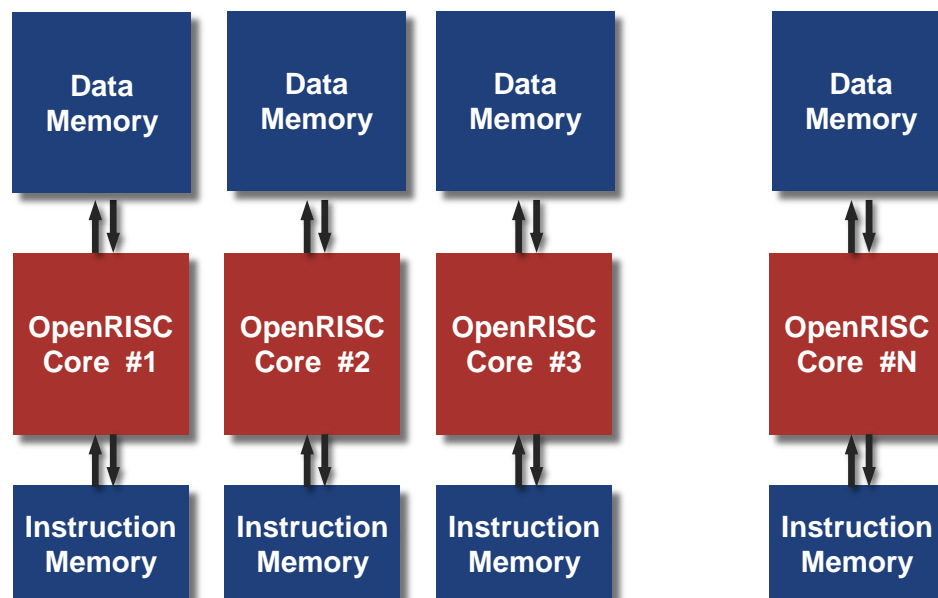
***1.57TOPS/W**: Kim et al., “A 1.22TOPS and 1.52mW/MHz Augmented Reality Multi-core Processor with Neural Network NoC for HMD Applications”, ISSCC 2014

Making PULP



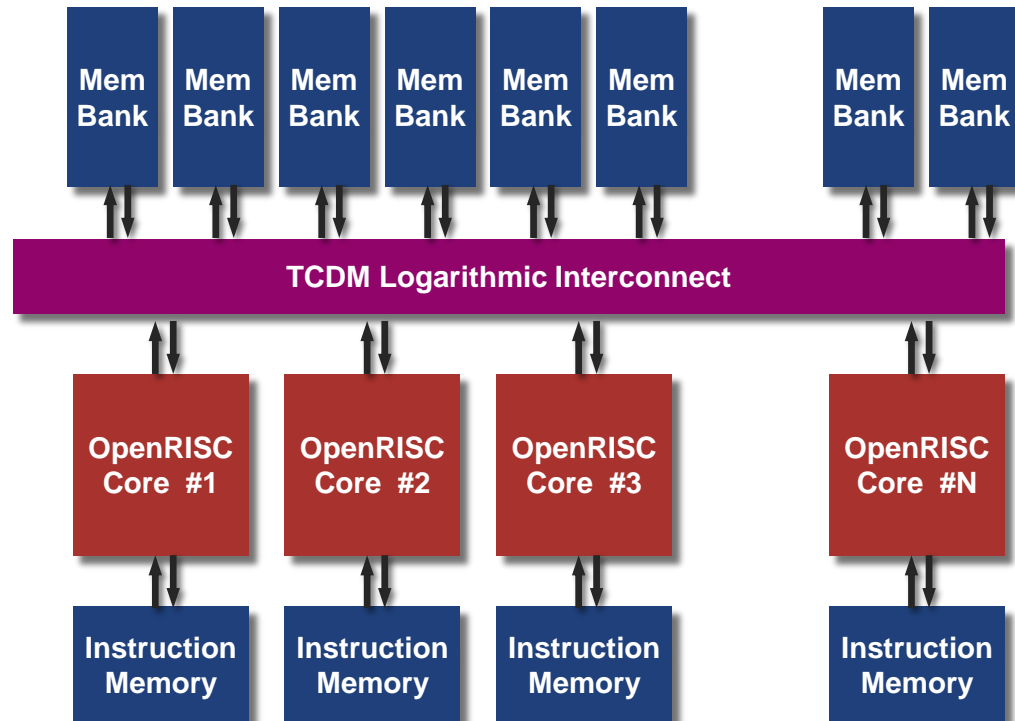
Start with an simple+efficient processor (~1IPC)

Making PULP



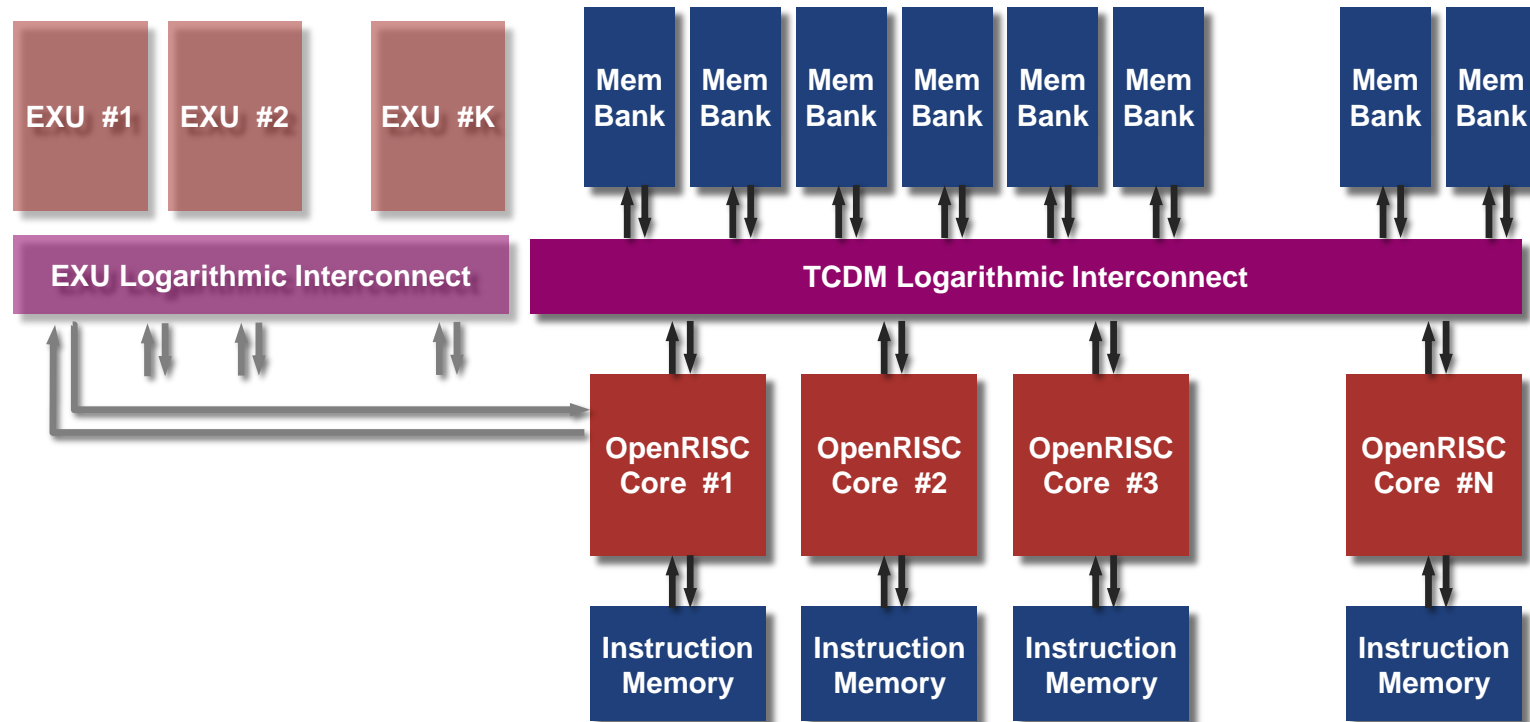
Parallel processors for performance @ NT

Making PULP



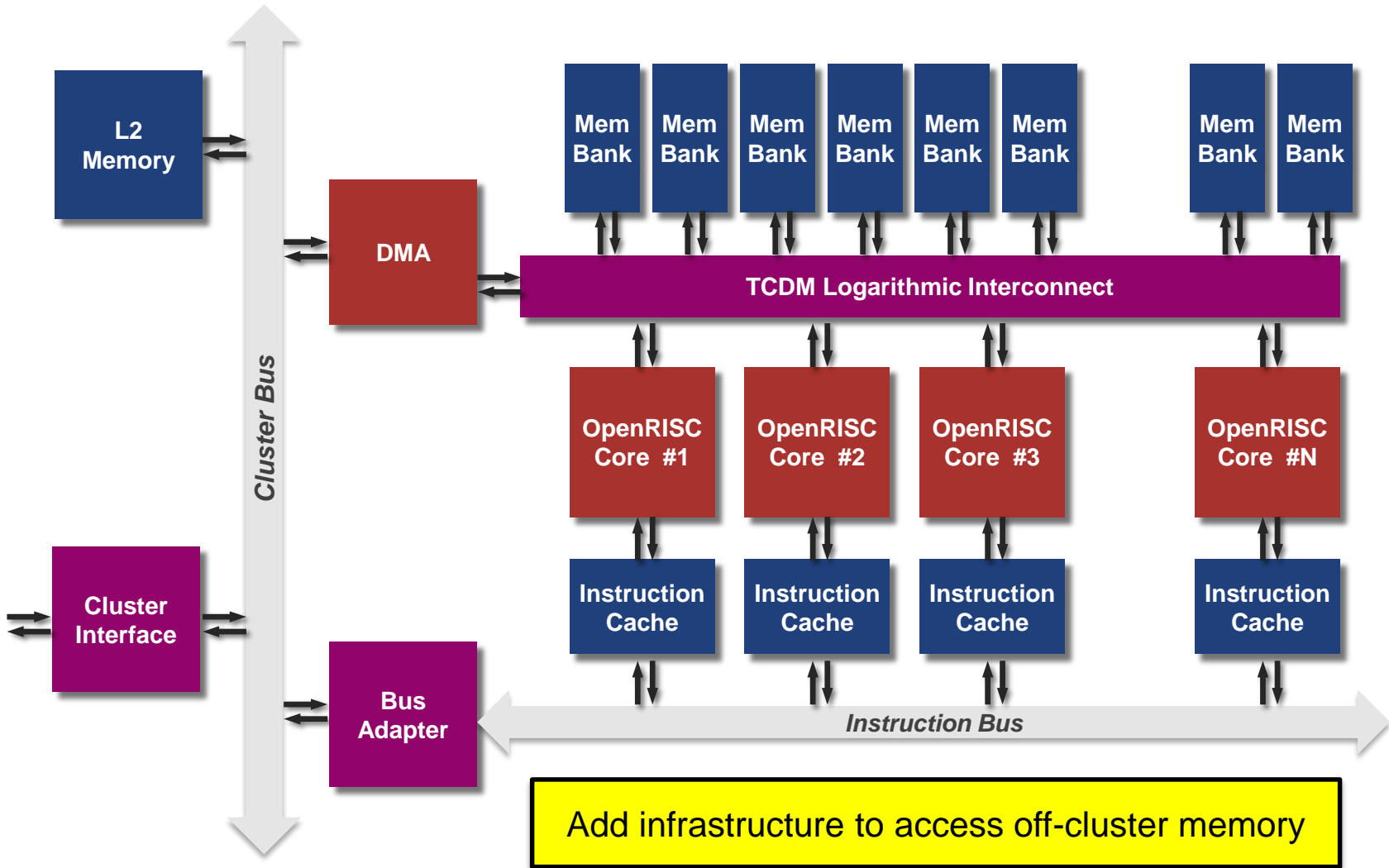
Parallel access to shared memory → Flexibility

Making PULP

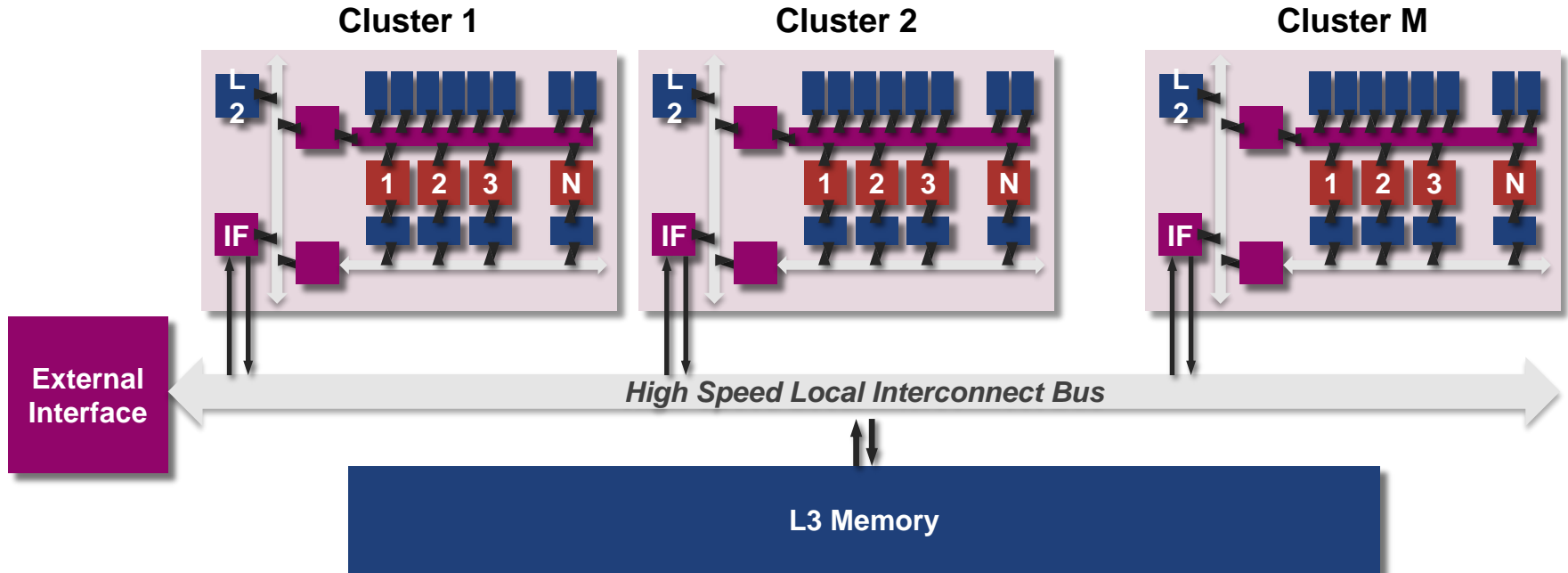


Optional Instruction Extension → Acceleration

Making PULP

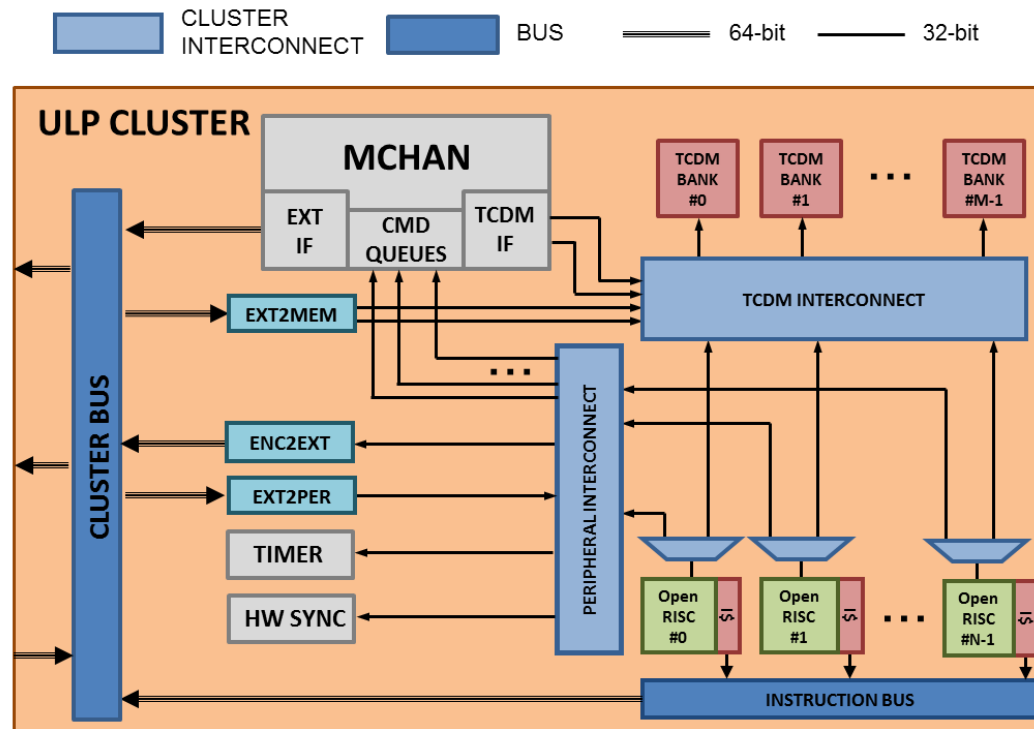


Making PULP



Multiple clusters (f, V_{dd}, V_{bb}) form a PULP system

PULP Cluster



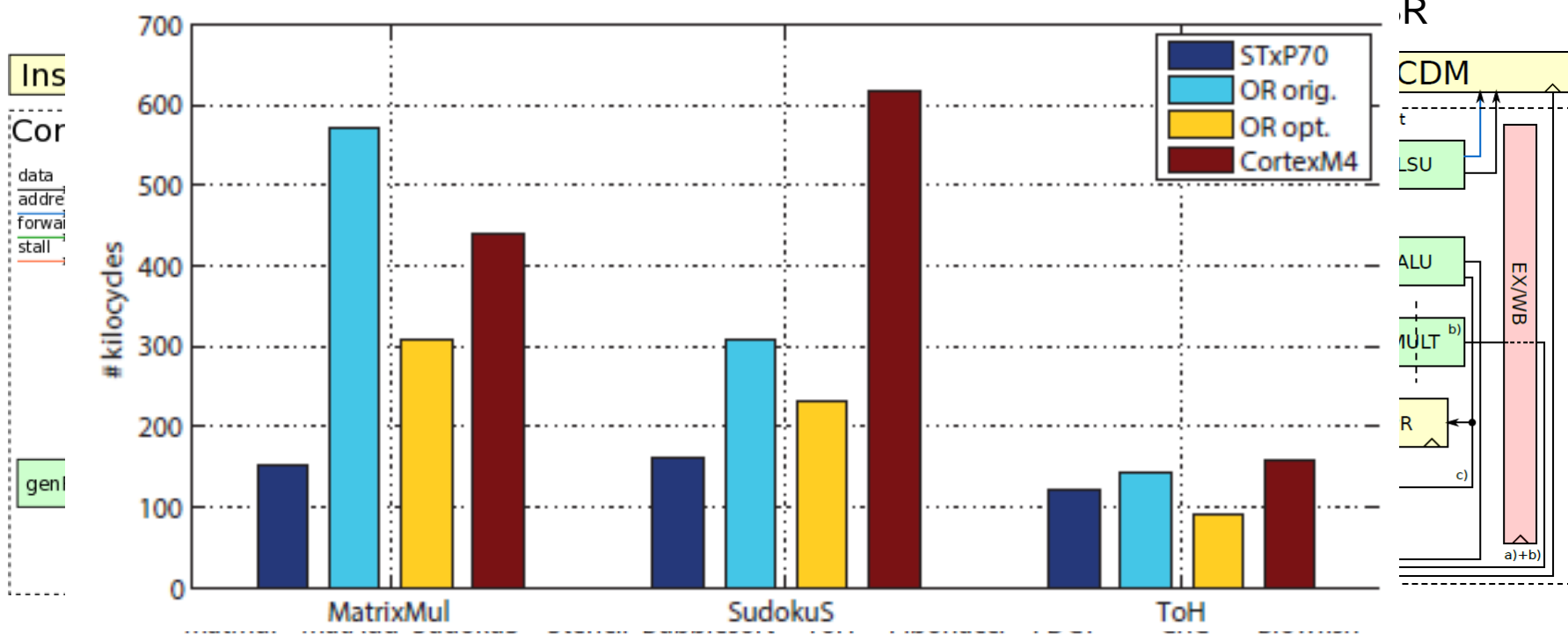
Design choices

- **I\$** → high code locality & simple architecture
- **No D\$** → low locality & high complexity: $B_{pmm}^2_{D\$} / B_{pmm}^2_{DTCDM} < 0,4$
- **Sharing L1** → less copies, easy work-balancing, low T_{clk} overhead in NT
- **Multibank** → smaller energy per access, “almost” multiported

OpenRisc Optimization

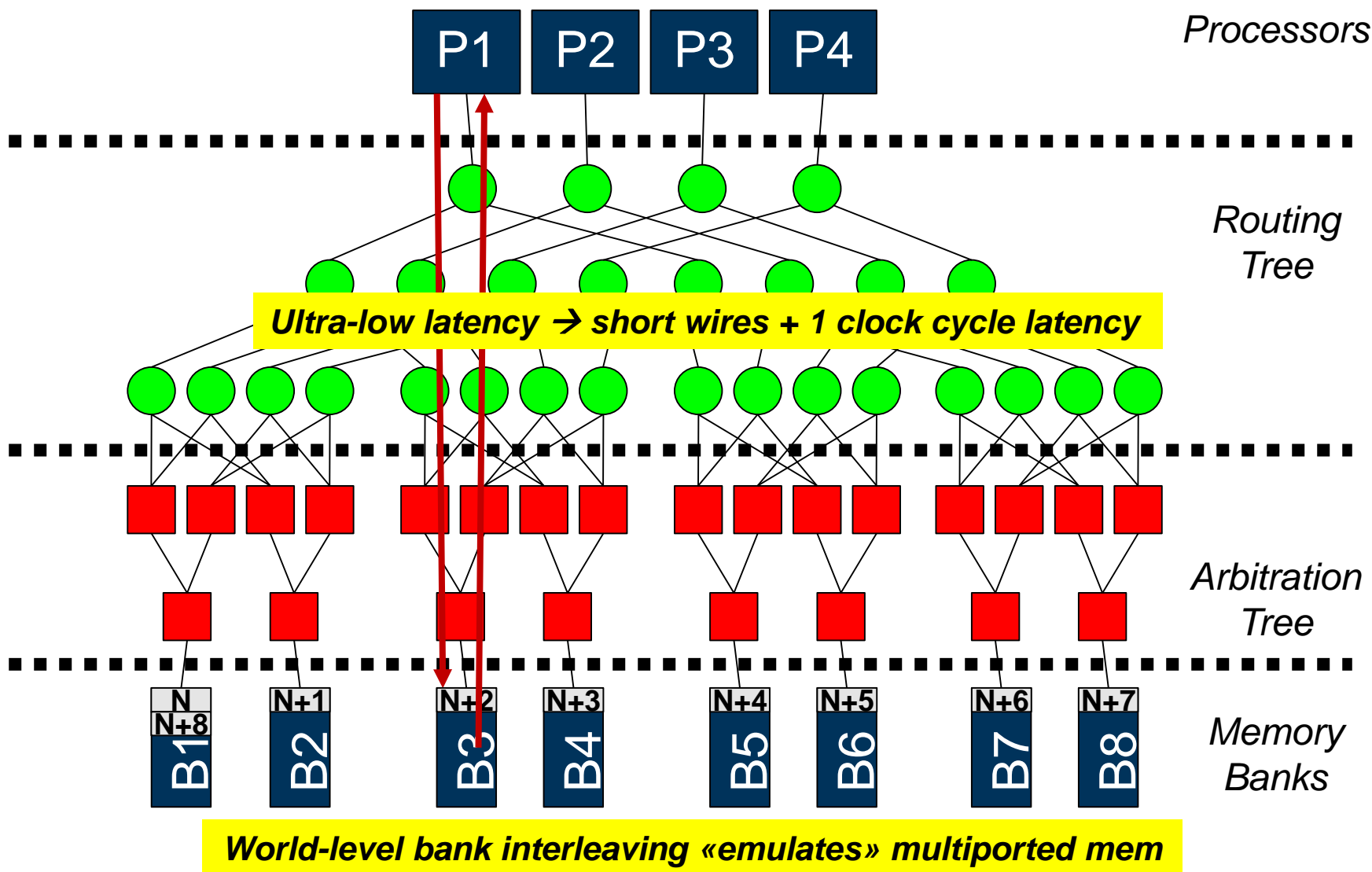


- Superpipelining harmful for energy efficiency
 - Focused speed optimization on the critical path dominated by MEM
- Low pipeline depth → high IPC with simple microarchitecture



50% less energy per operation on average, 5% more area

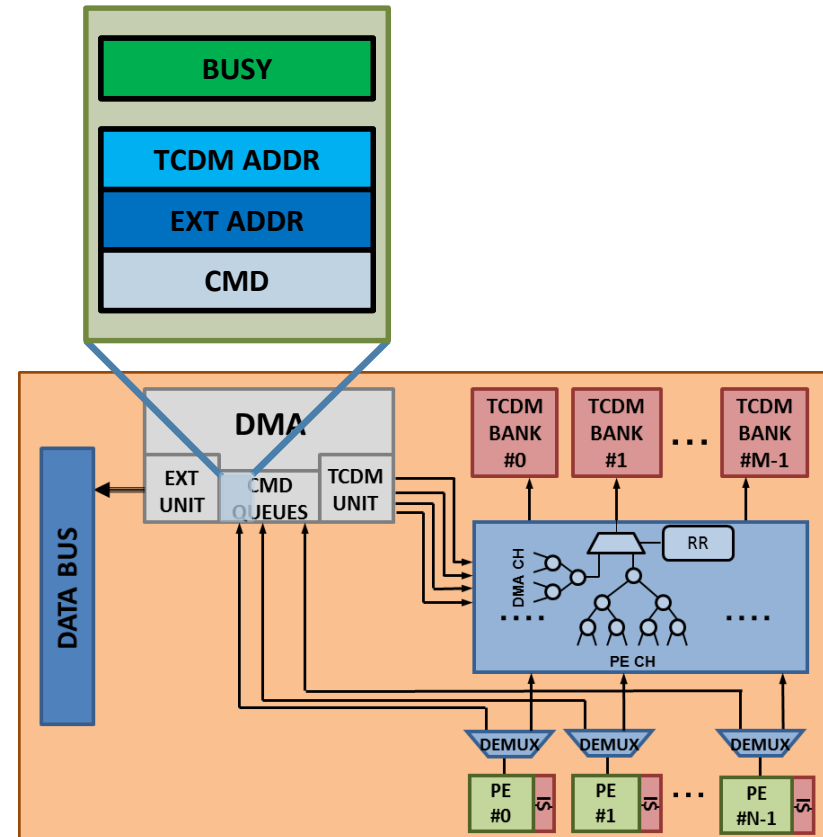
Logarithmic Interconnect



Low latency programming Interface



- Each command queue is dedicated to a core of the cluster: arbitration is made in hardware
 - No need to reserve (lock/unlock) the programming channel
- COREs program DMA through register mapped on the DEMUX
 - The registers belongs to aliases, no need for the processors to calculate (per-core) offsets
- A programming sequence requires
 1. check a free command queue
 2. write address of buffer in TCDM
 3. write address of buffer in L2 memory
 4. Trigger data transfer
 5. Synchronization



Programming Latency: ~10 CLOCK CYCLES!!!

DMA Architecture Overview

CTRL UNIT:

- Arbitration – forwarding and synchronization of incoming requests

TRANSFER UNIT:

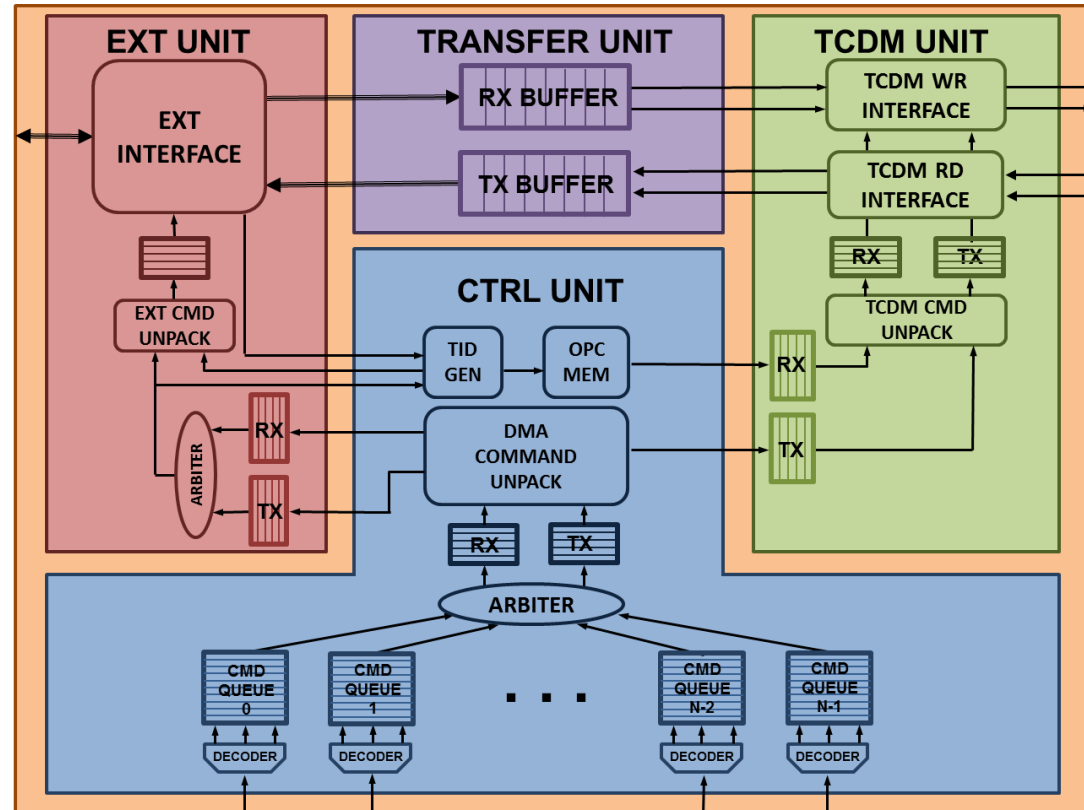
- FIFO Buffers for TX and RX channels

TCDM UNIT:

- Bridge to TCDM protocol – 4 channels (2 RD and 2 WR) 32 bit each

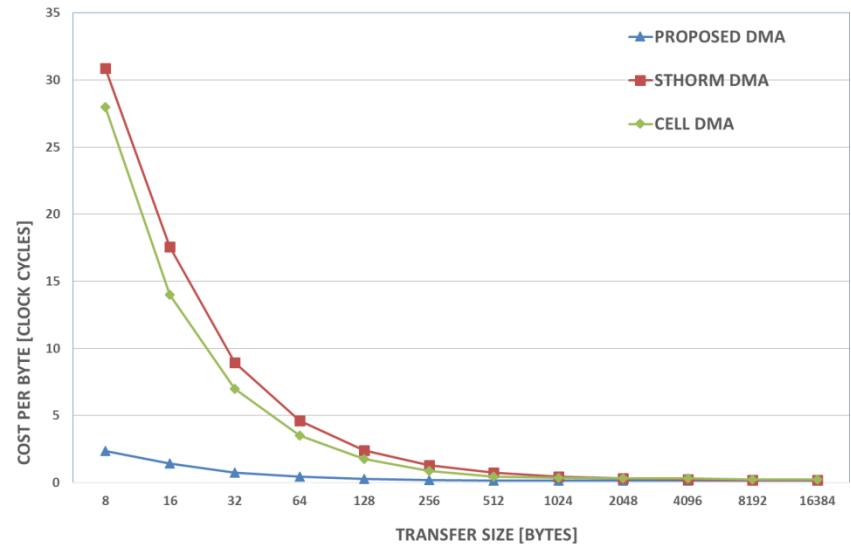
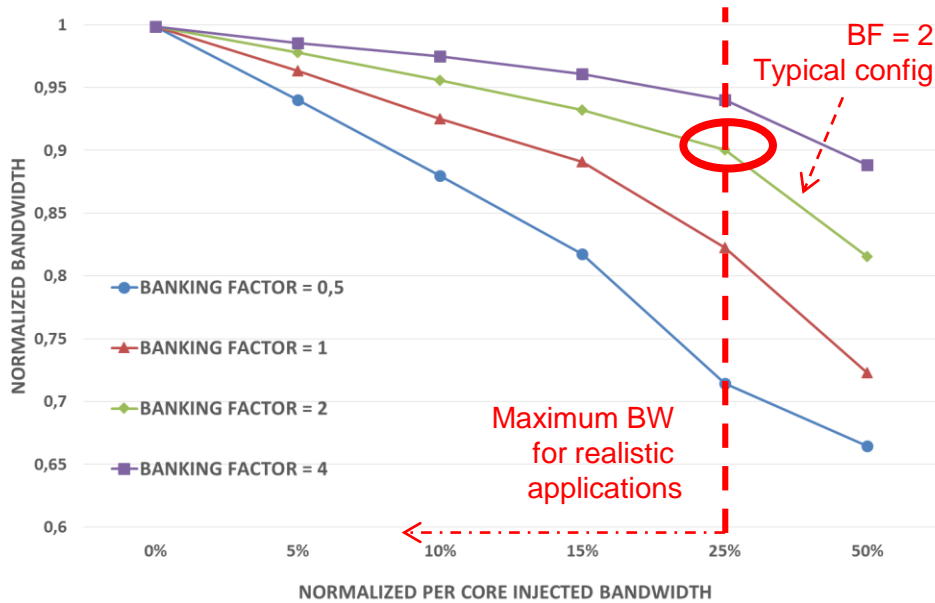
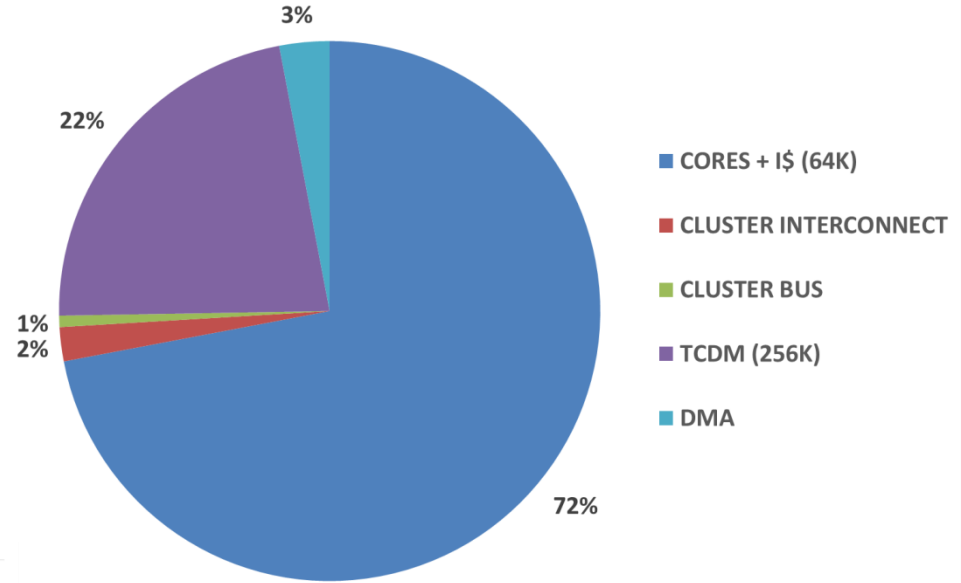
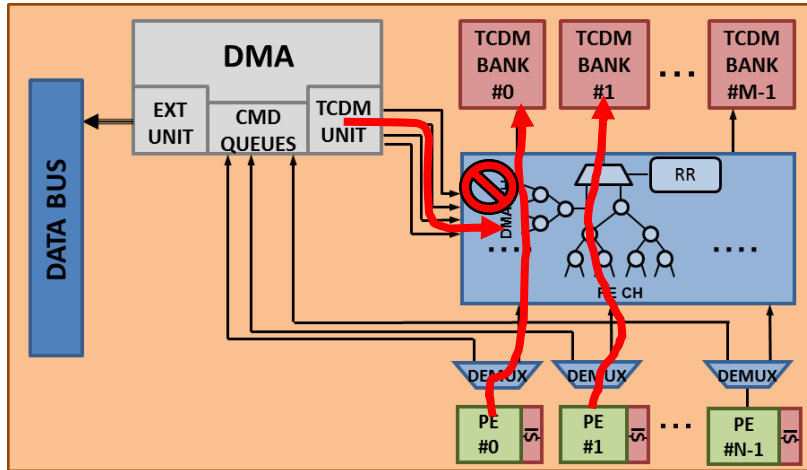
EXT UNIT:

- Bridge to AXI, 64 bit

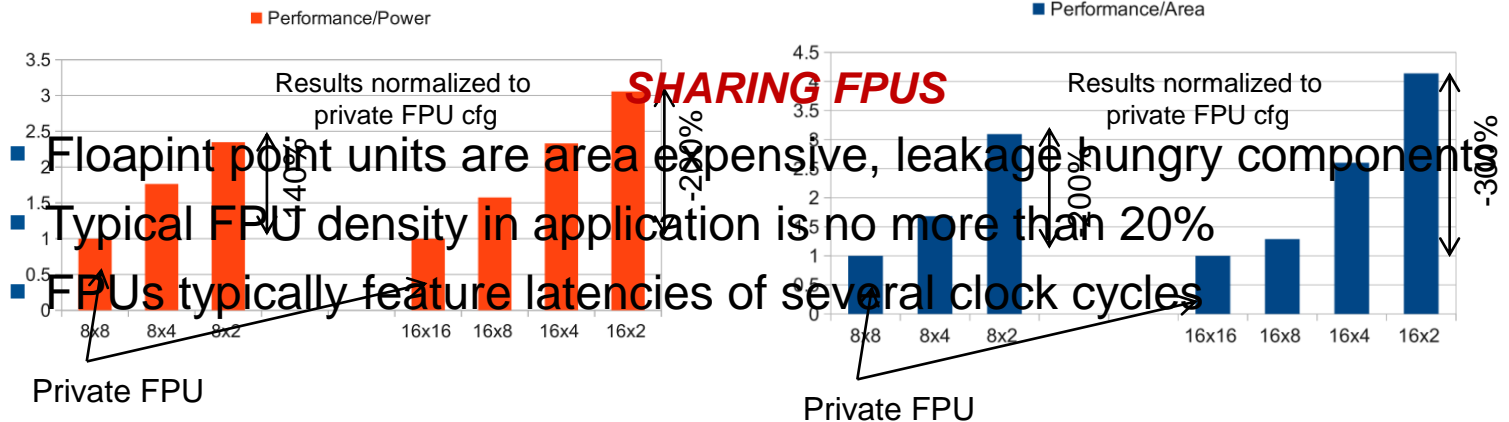
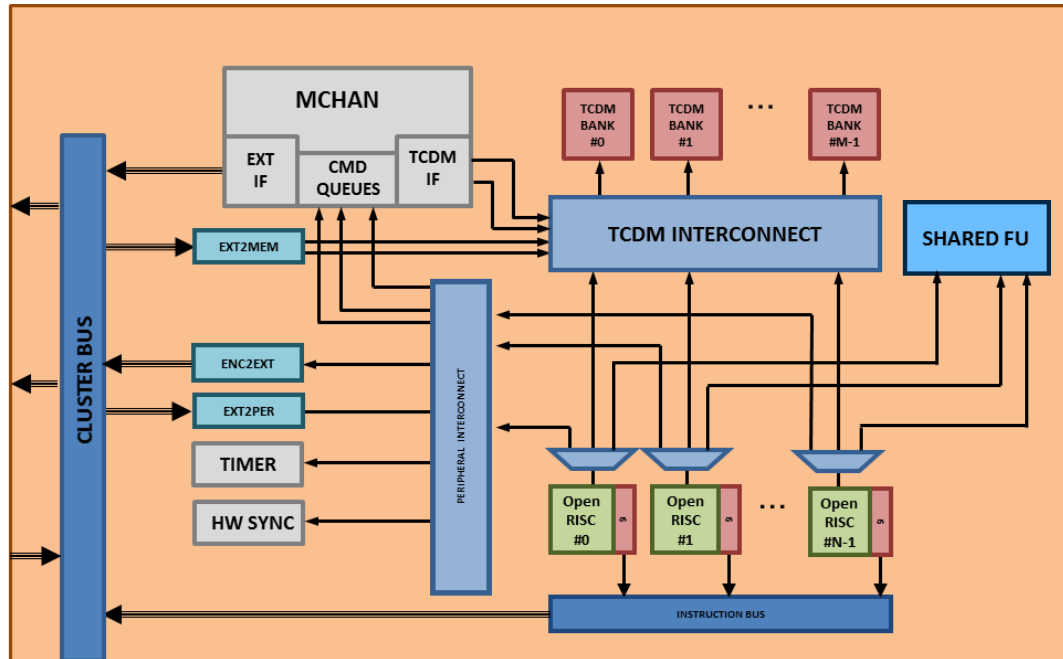


Key idea: only channel packets buffered internally – no DMA transfers!

Cluster DMA



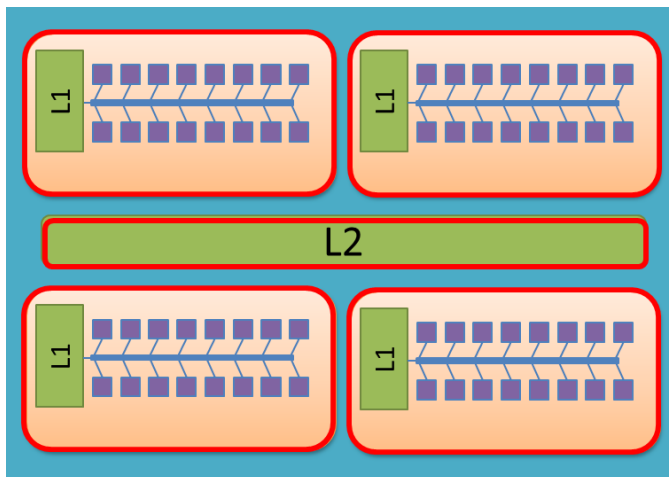
Sharing functional units (instruction set extensions)



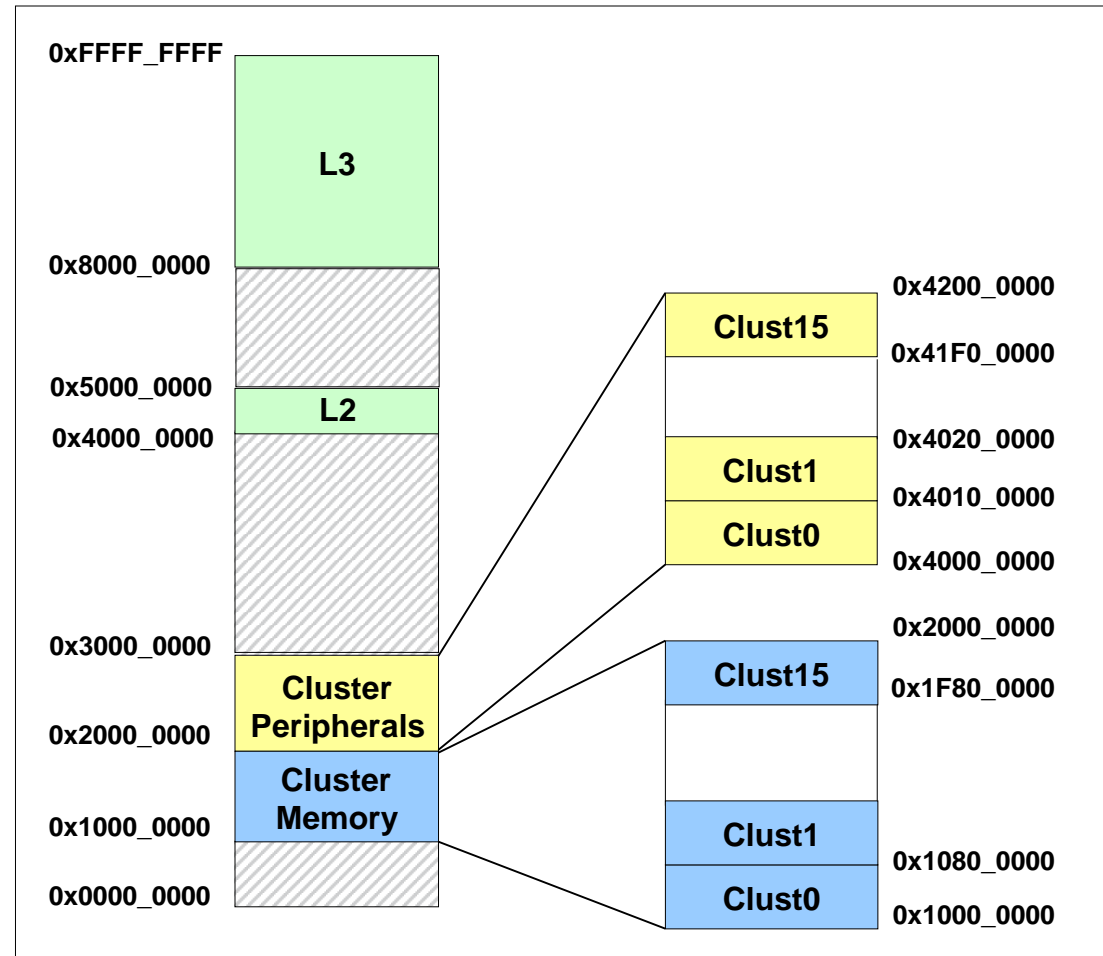
Scaling up



- 32 bit architecture
- ⇒ 4 GB of memory
- Clusters in V_{dd}, CLK domains
- L2 (2D & 3D) ready
- Host VM IF for Heterogeneous Computing



Memory Map:



Programmability: OpenMP

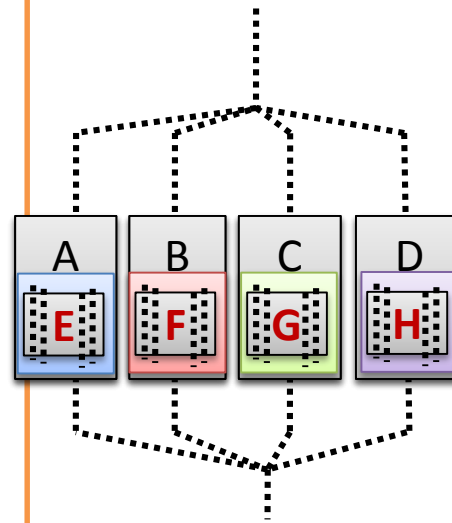


```
while(1)
{
  #pragma omp parallel num_threads(4)
  {
    #pragma sections
    {
      #pragma section
      {
A. #pragma omp parallel num_threads(16)
        ColorScaleConv();
      }
    }
    #pragma section
    {
B.
    }
    #pragma section
    {
C. #pragma omp parallel num_threads(16)
        cvMoments();
    }
    #pragma section
    {
D. #pragma omp parallel num_threads(16)
        cvAdd();
    }
  }
}
```

A powerful abstraction for specifying structured parallelism

```
void ColorScaleConv()
{
  #pragma omp for
  for(i = 0; i < FRAME_SIZE; i++)
  {
    [ALGORITHM]
  }
}
```

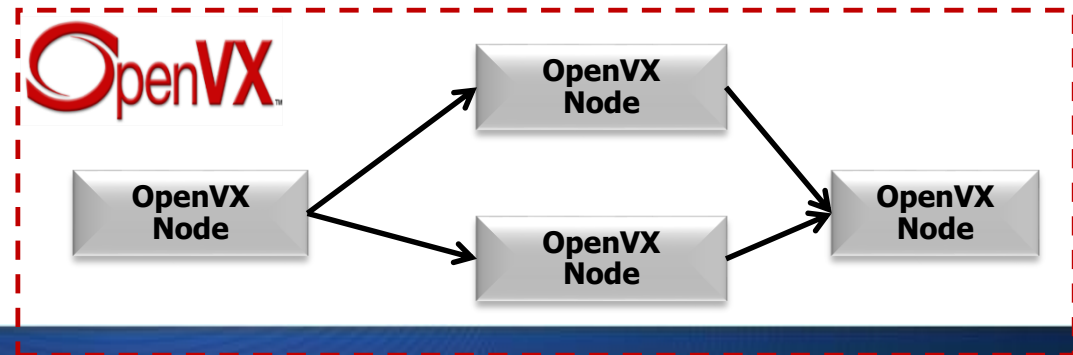
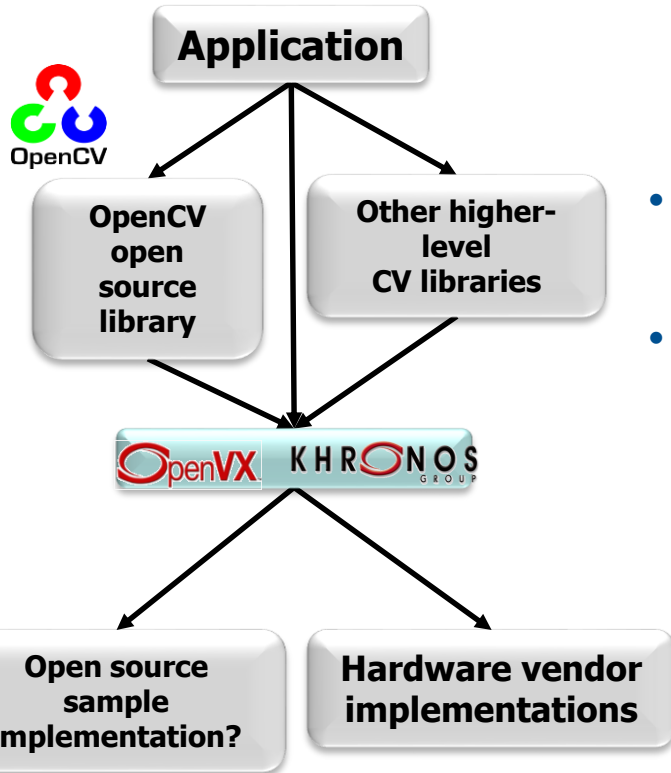
And very suitable for NUMA (cluster-based) systems



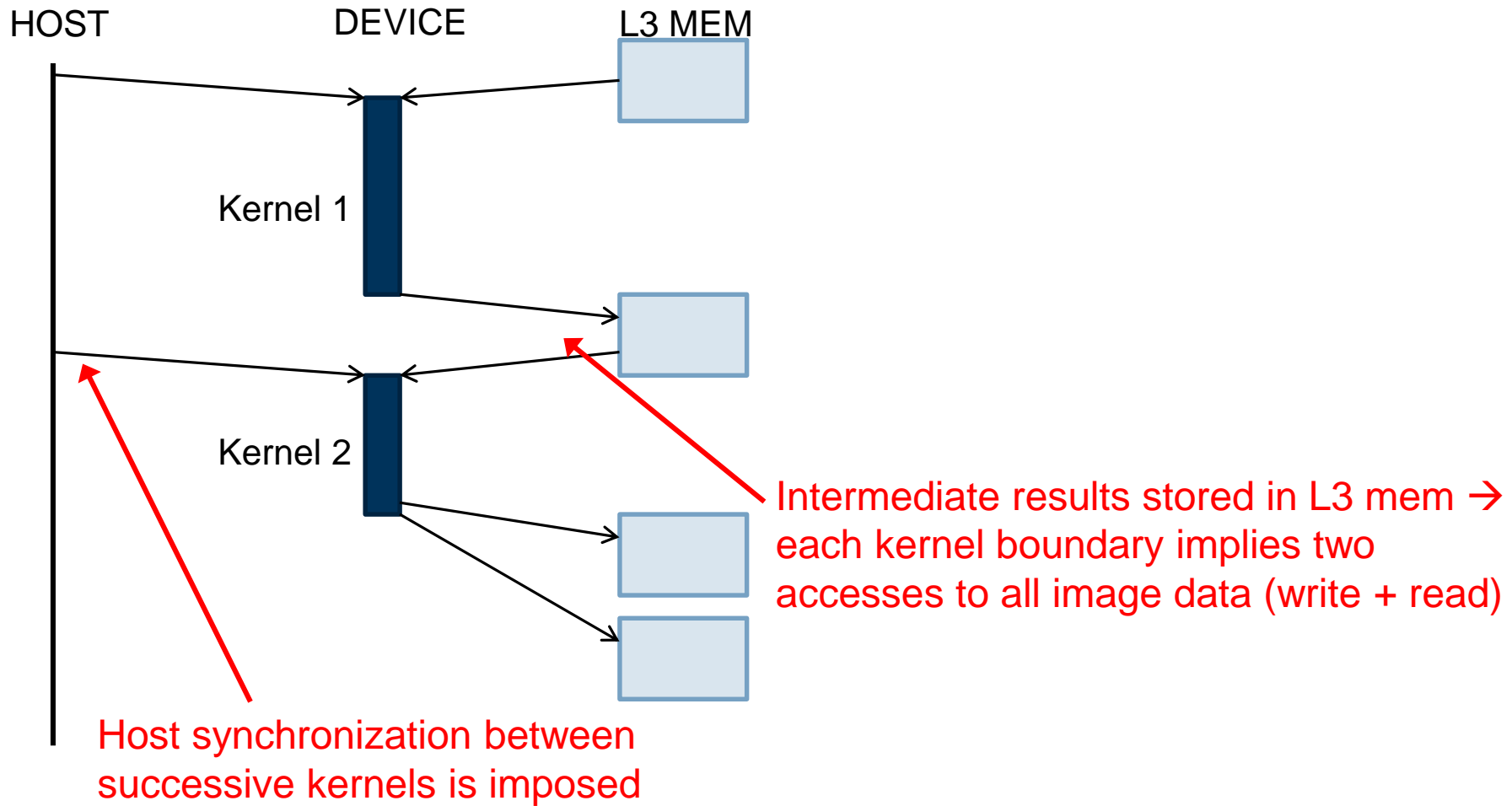
LLVM + OpenCL

+OpenVX → domain specific language

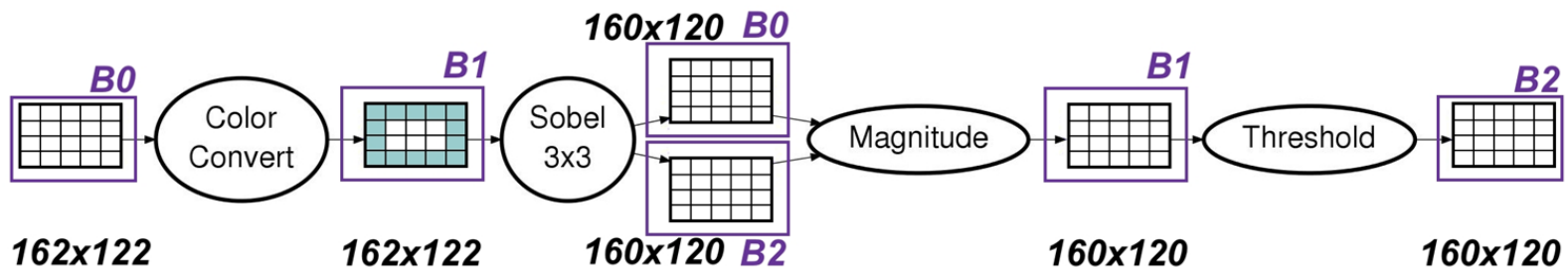
- C-based standard API for vision kernel
 - Defines a set of **standard kernels**
 - Enables hardware vendors to *implement accelerated imaging and vision algorithms*
- Focus on enabling real-time vision
 - On mobile and embedded systems
- **Graph execution model**
 - Each node can be implemented in software or accelerated hardware
 - Data transfer between nodes may be optimized



Leverages OpenCL runtime with OVX nodes treated as OCL kernels



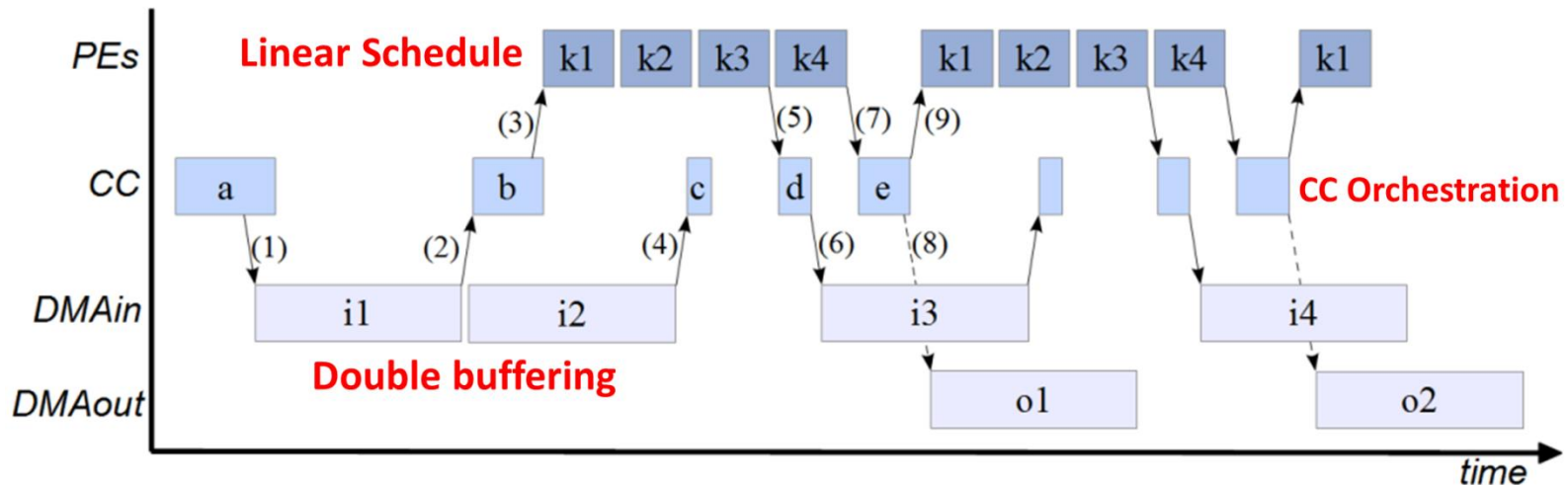
- Smaller image sub-regions (**tiles**) totally reside in TCDM
- All kernels are computed at tile level, no more at image level → *intermediate results are also allocated in TCDM*



- ✓ TCDM is partitioned into 3 buffers (B0, B1, B2)
- ✓ Output tile size is 160x120
- ✓ Sobel3x3 requires image overlapping (1 pixel for each direction), and so the tile size is 162x122
- ✓ ColorConvert does not require tile overlapping, but must provide a 162x122 result tile for the next stage → **tile size propagation**

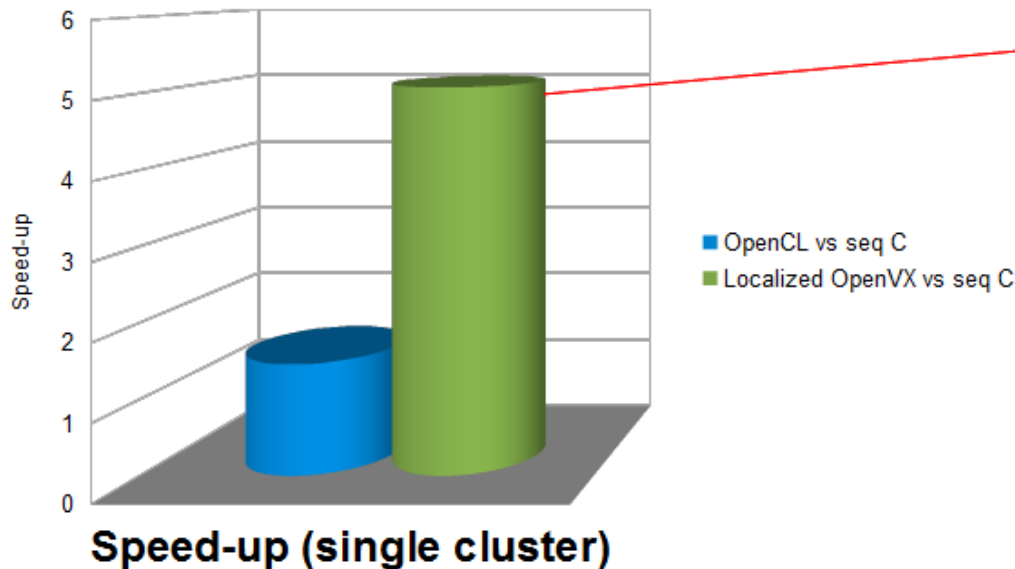
Localized Execution

- Smaller image sub-regions (**tiles**) totally reside in TCDM
- All kernels are computed at tile level, no more at image level → *intermediate results are also allocated in TCDM*

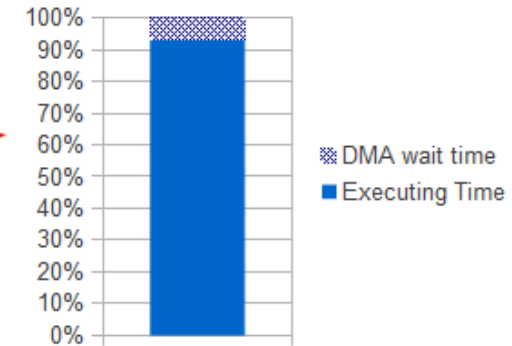


- ✓ TCDM is partitioned into 3 buffers (B0, B1, B2)
- ✓ Output tile size is 160x120
- ✓ Sobel3x3 requires image overlapping (1 pixel for each direction), and so the tile size is 162x122
- ✓ ColorConvert does not require tile overlapping, but must provide a 162x122 result tile for the next stage → **tile size propagation**

- Framework prototype
 - OpenVX support
 - A limited subset of kernel has been implemented
 - Polynomial time (i.e. fast) heuristics suitable for *just-in-time execution*



Efficiency is very high



Cluster PE Efficiency

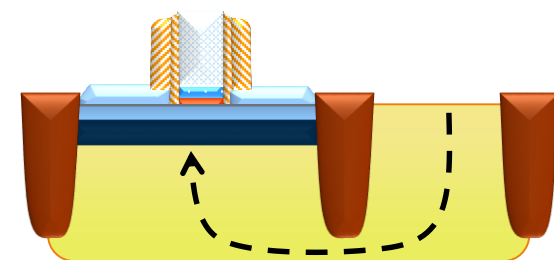
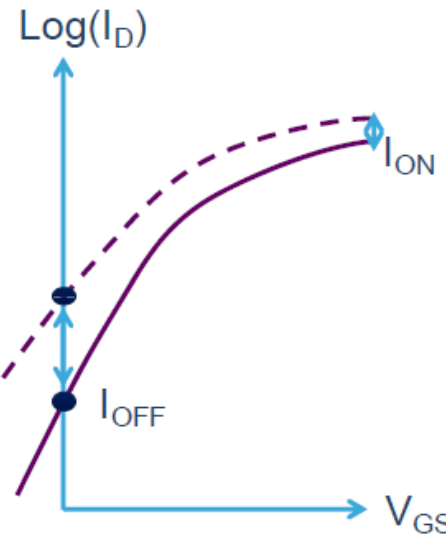
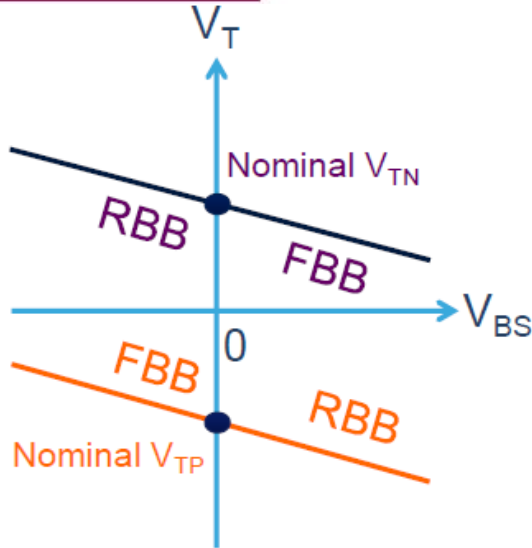
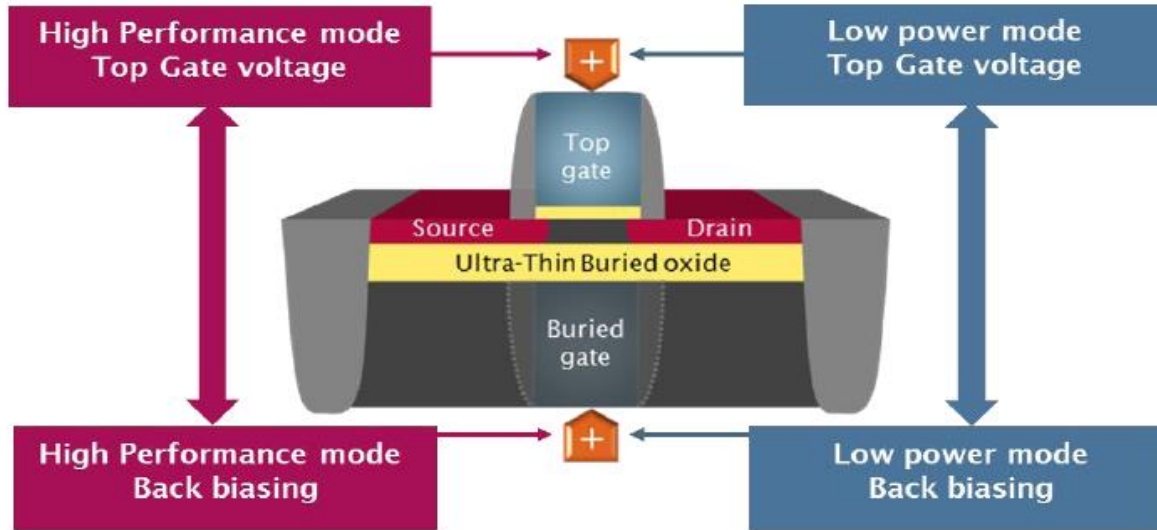
And what about Technology?



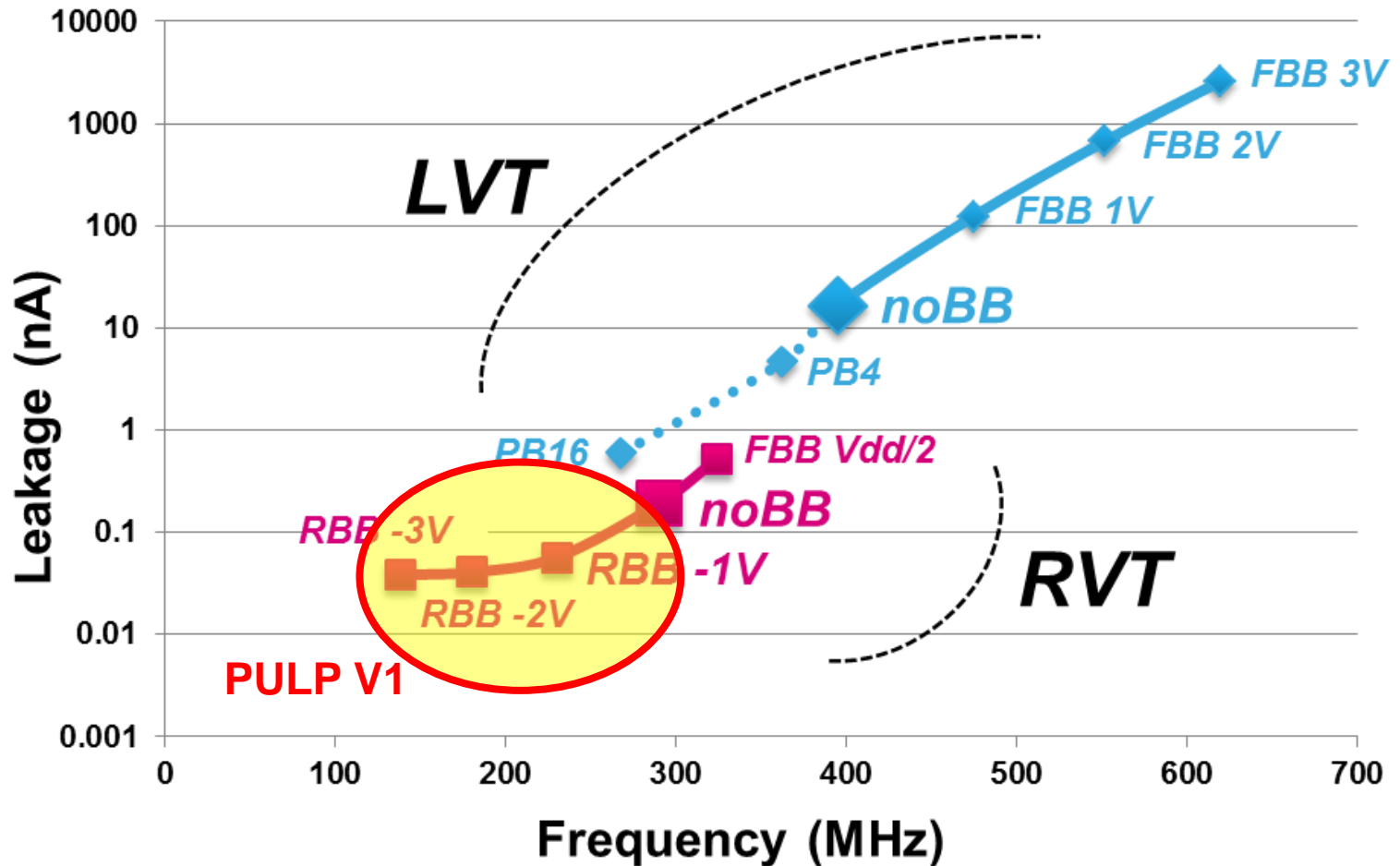
Near threshold FDSOI technology



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA



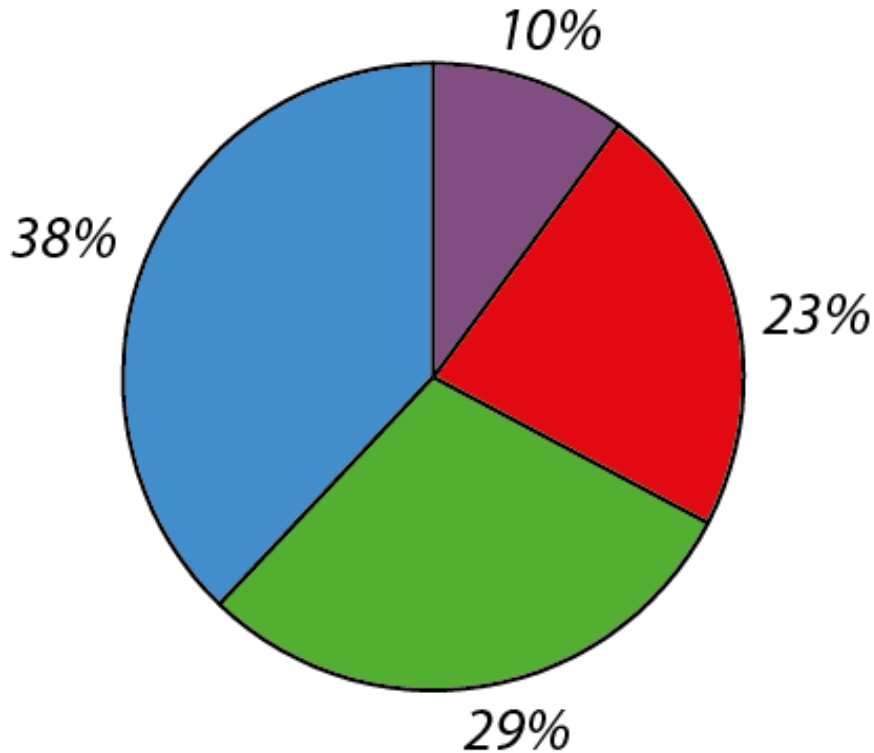
Body bias: Highly effective knob for leakage control!



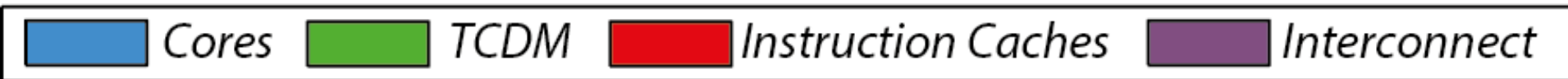
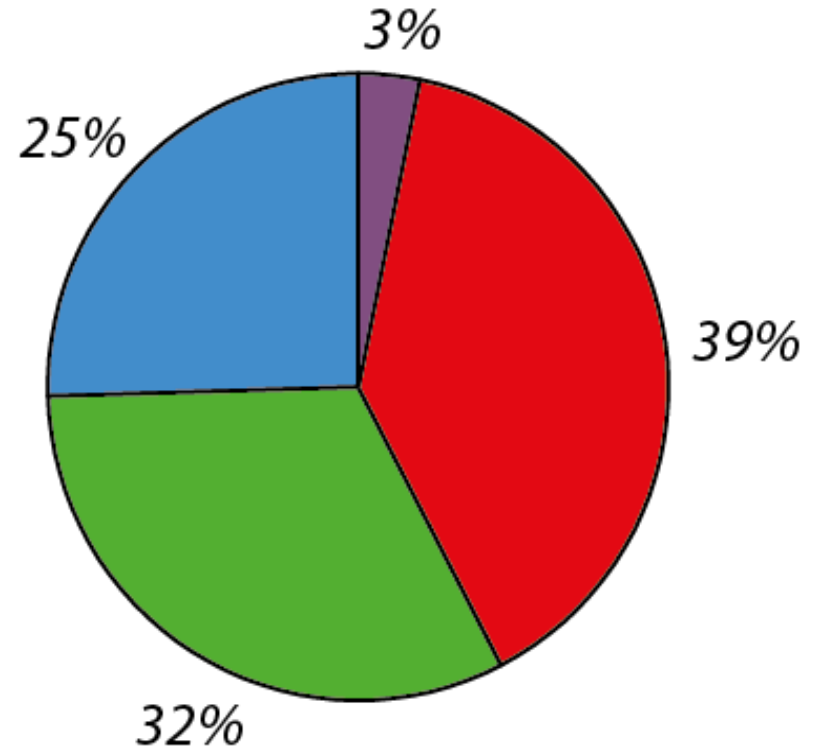
PULP V1: Doing nothing well (with RBB)



Area Breakdown



Power Breakdown



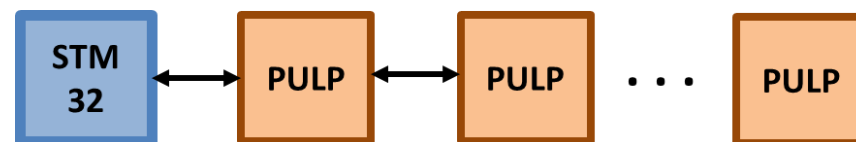
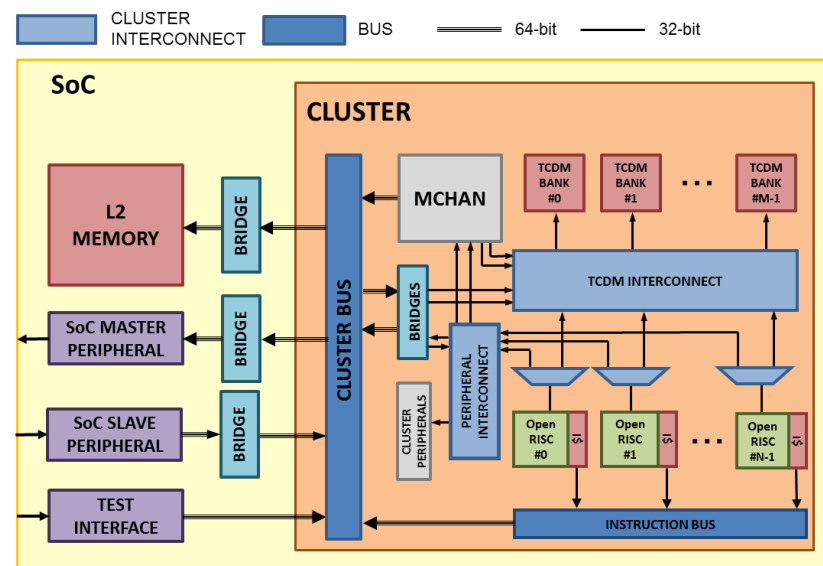
More than 50% of power into memories... this is the next focus area!

Introducing PULP V2



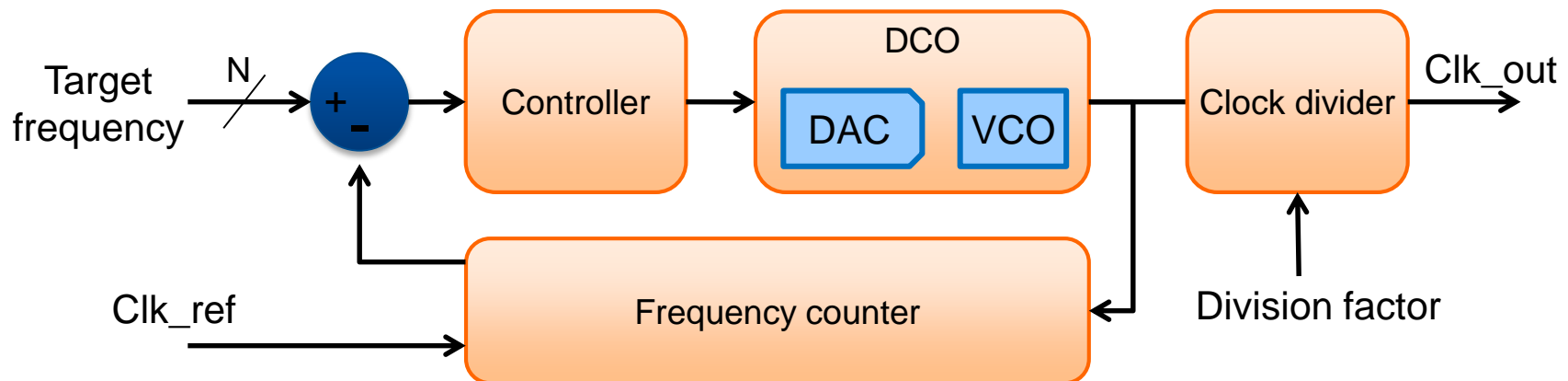
Board/Application-ready chip

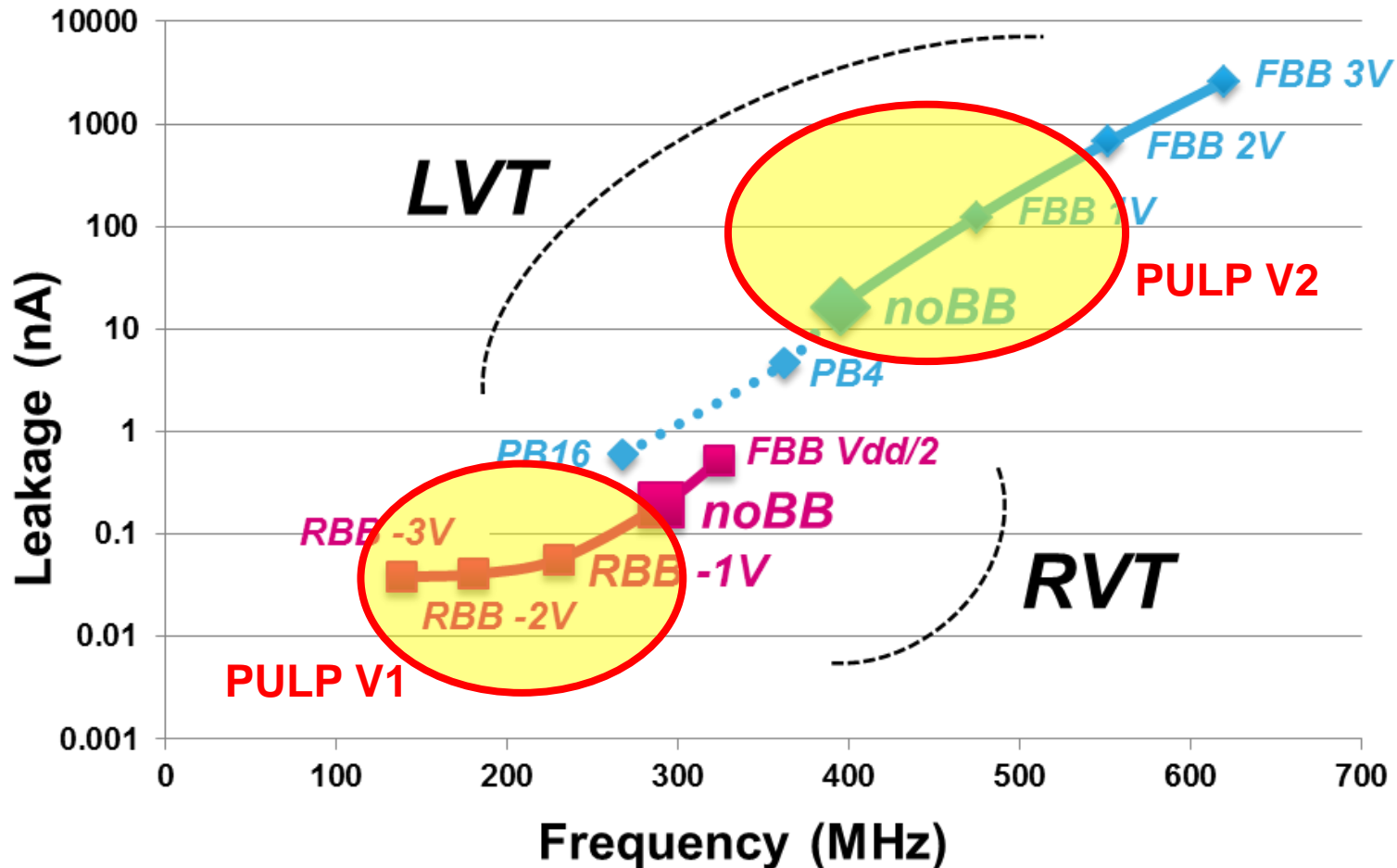
- Implementation of a master and a slave peripheral
 - Standard peripheral (e.g. SPI)
 - Integration with FPGAs or standard low power external memories
- PULP as multi-core accelerator for micro-controller host
 - STM32 host core
 - PULP multi-core accelerator
- Daisy chain of PULP chips
 - Pipeline of parallel processing units
 - Each core performs a stage of computation and forward temporary data to another stage



Standalone mode is also supported!

- All-digital clock generation based on a Frequency Lock Loop (FLL)
- From 2GHz down to 15KHz (through clock division)
- Frequency step 10MHz (at lowest division factor)
- Small area 3300 μm^2 (50 times smaller than classic PLL)
- Suited to fine-grain GALS architectures
- Frequency reprogramming in less than 180ns
- No frequency overshoot
- 15ps jitter

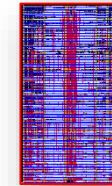
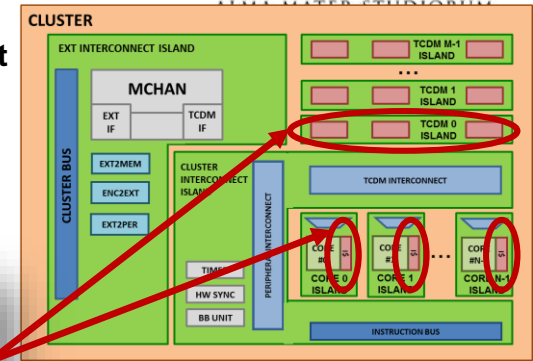




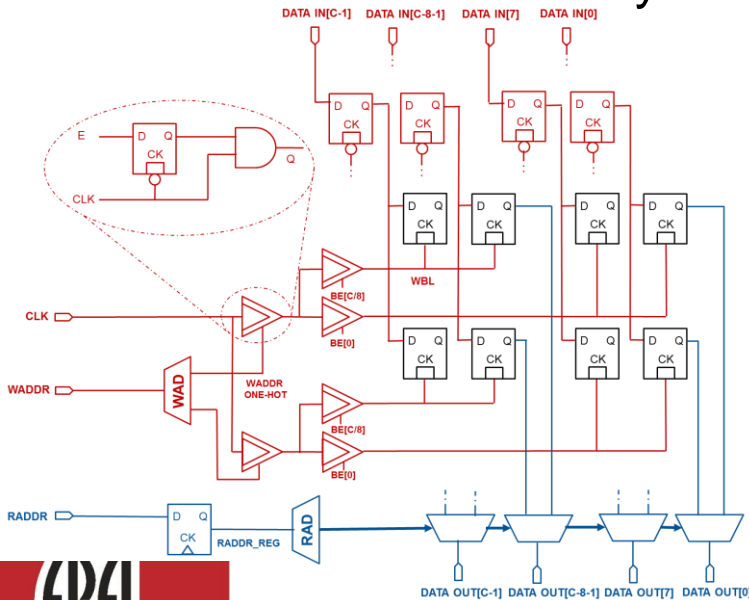
Low-leakage vs. Low voltage (0.3V) → **reactive** or **proactive**?

- Based exclusively on standard cells
- Voltage range identical to the core
 - Only static logic
- Layout based on guided P&R
 - Close to 100% density

64 words x 64 bit
162 x 85 =
13.7k um²



SCMs

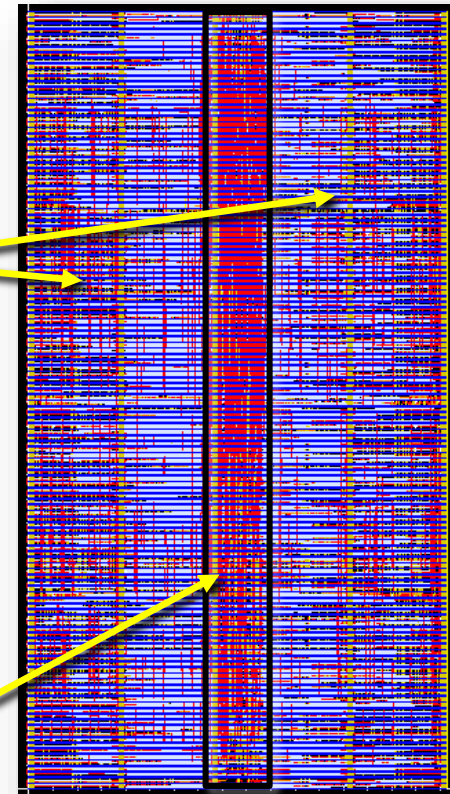


Macro Size:
128 x 32 (4 kb)
86µm x 160µm

Input/Output Delay:
0.3ns/0.7ns @ss0.9V125°
2ns/3.3ns @tt0.3V25° (FBB)

Storage Array + output mux

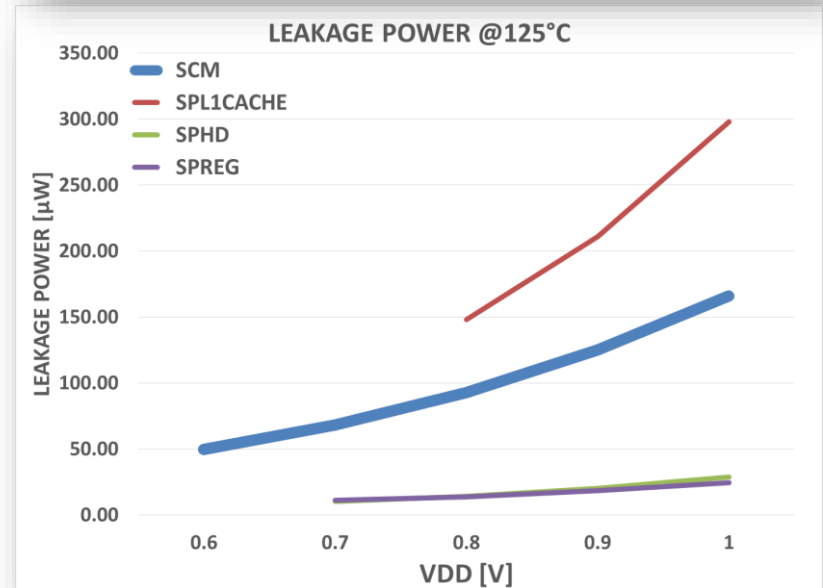
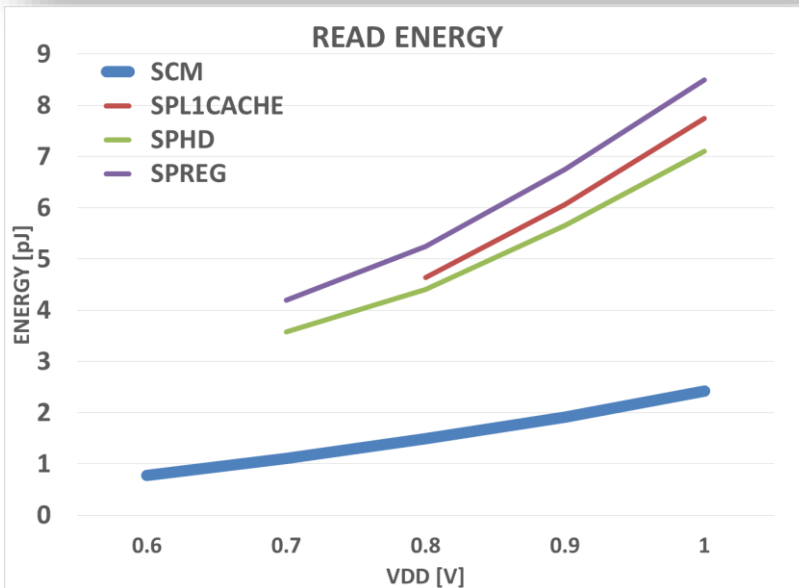
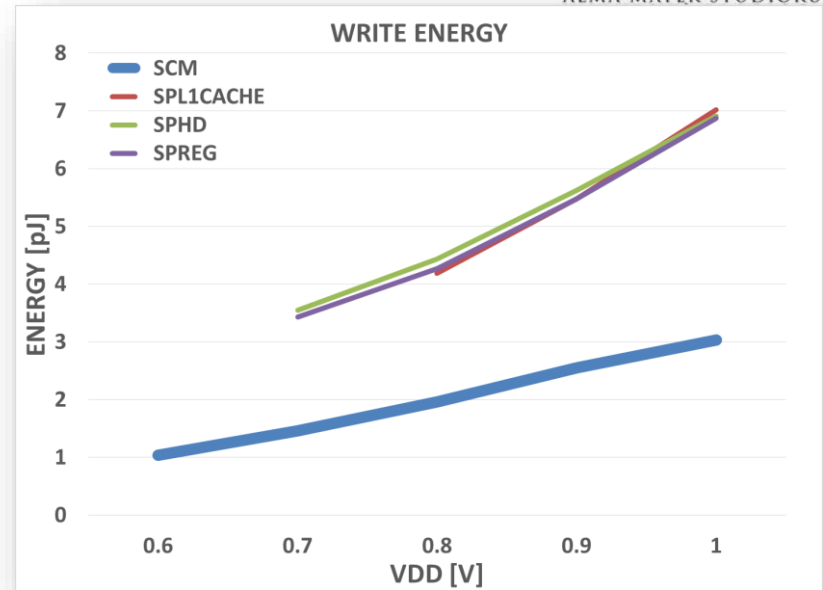
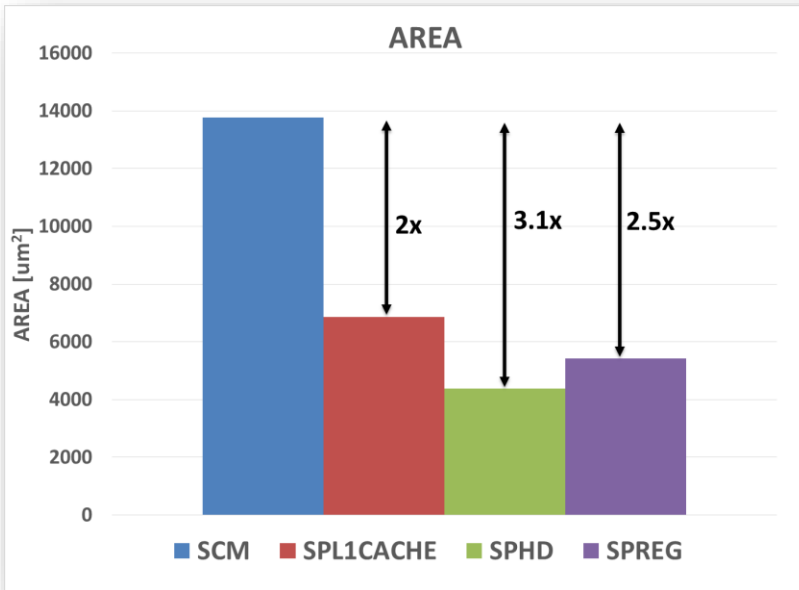
Decoders

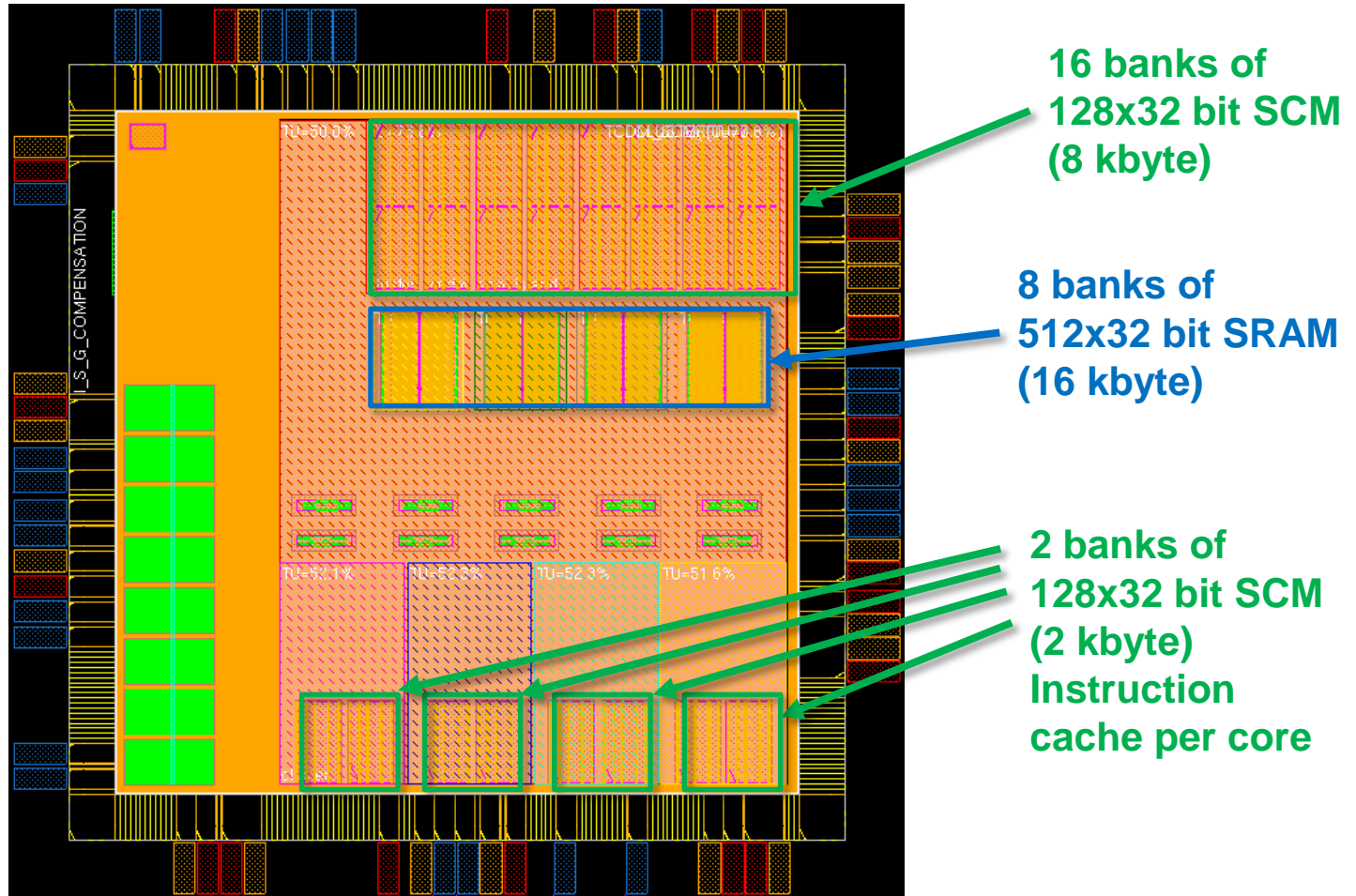


Comparison to 64x64 SRAM



ALMA MATER STUDIORUM
A

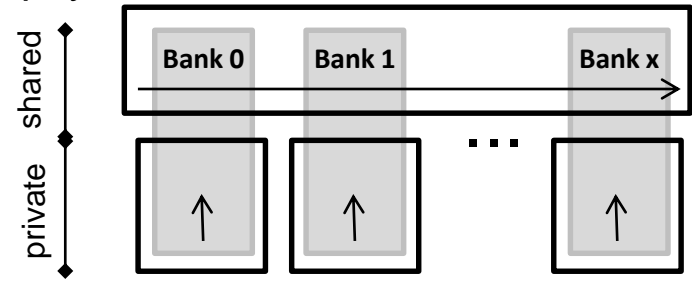




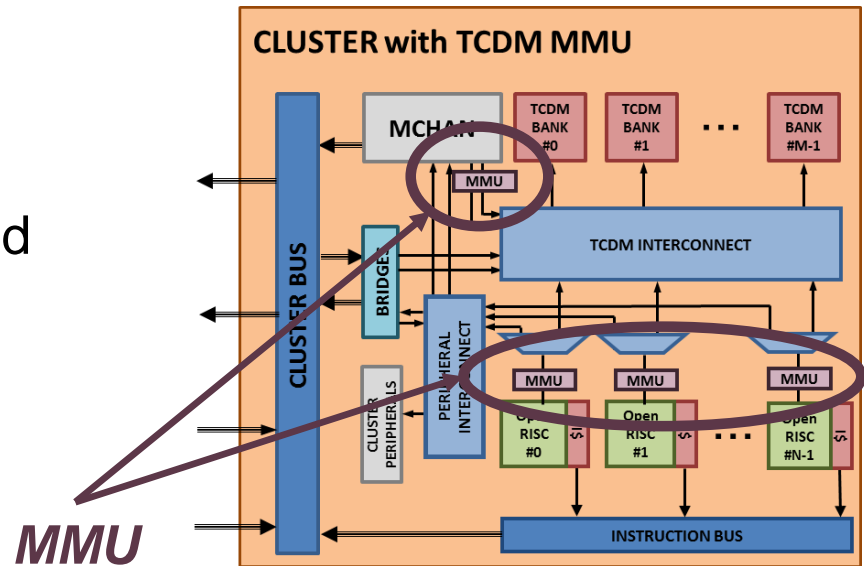
PULP V2 is taping out this Week, PULP V3 is on the drawing board...

- Support for different address mappings:
 - Interleaved: horizontal shutdown and reduces conflicts in shared segments
 - Non-interleaved mapping: private memory avoids conflicts

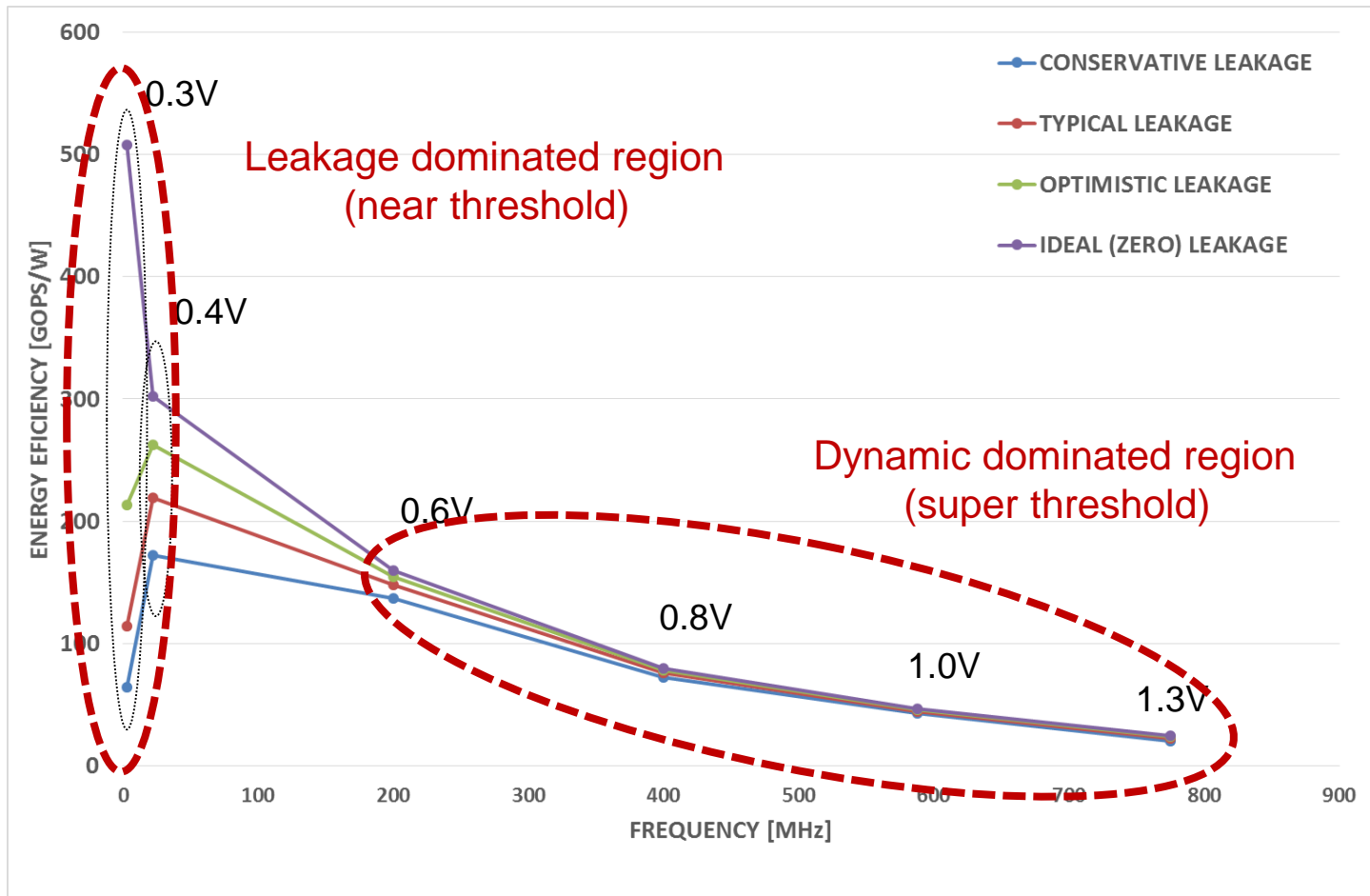
Mapping between logical-physical addresses



- Basic MMU
 - Coexistence of both shared and private memory segments with different address mappings
 - Adapt address mapping is adapted to accommodate partial memory shut-down

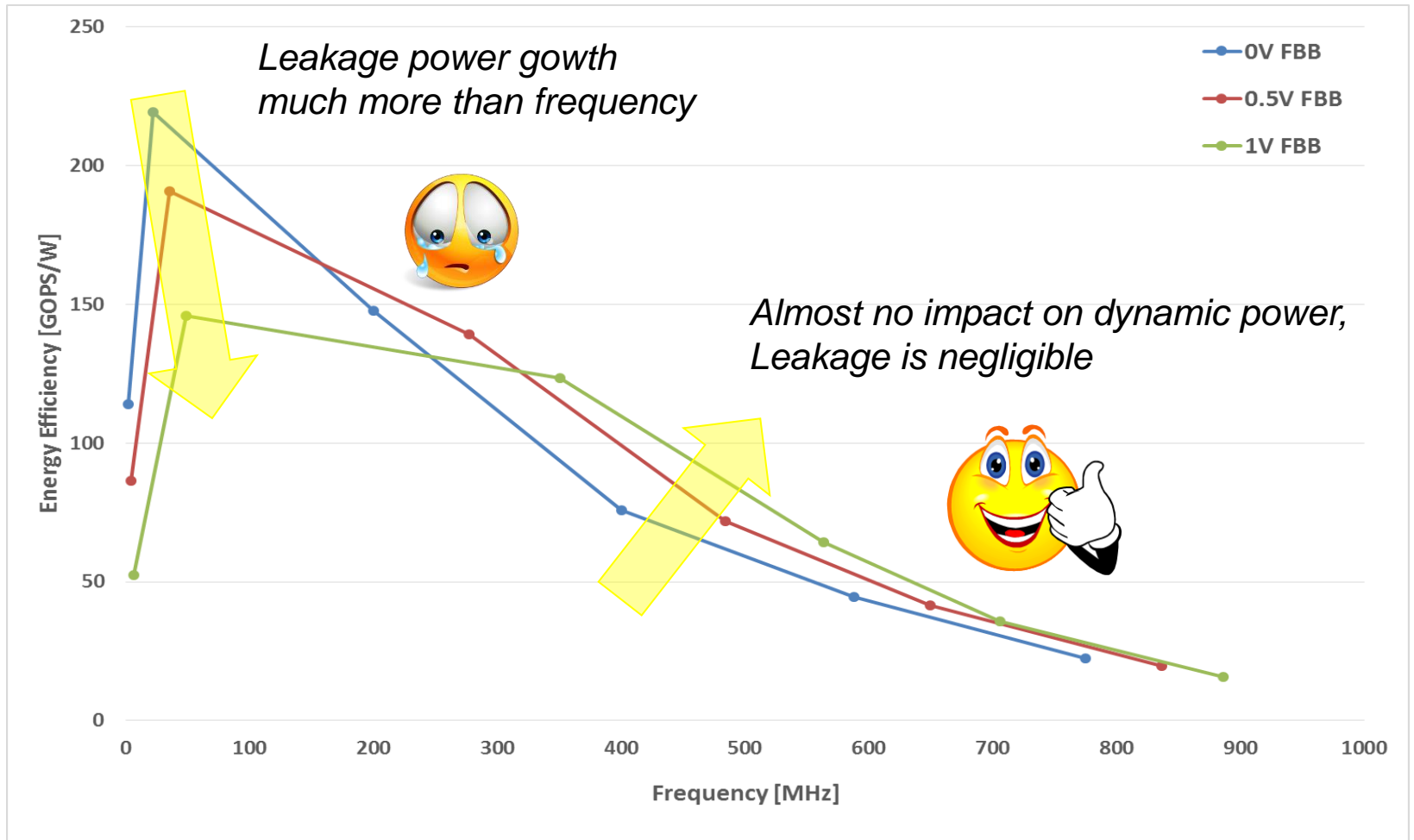


PULP2 Energy Efficiency (Estimated)



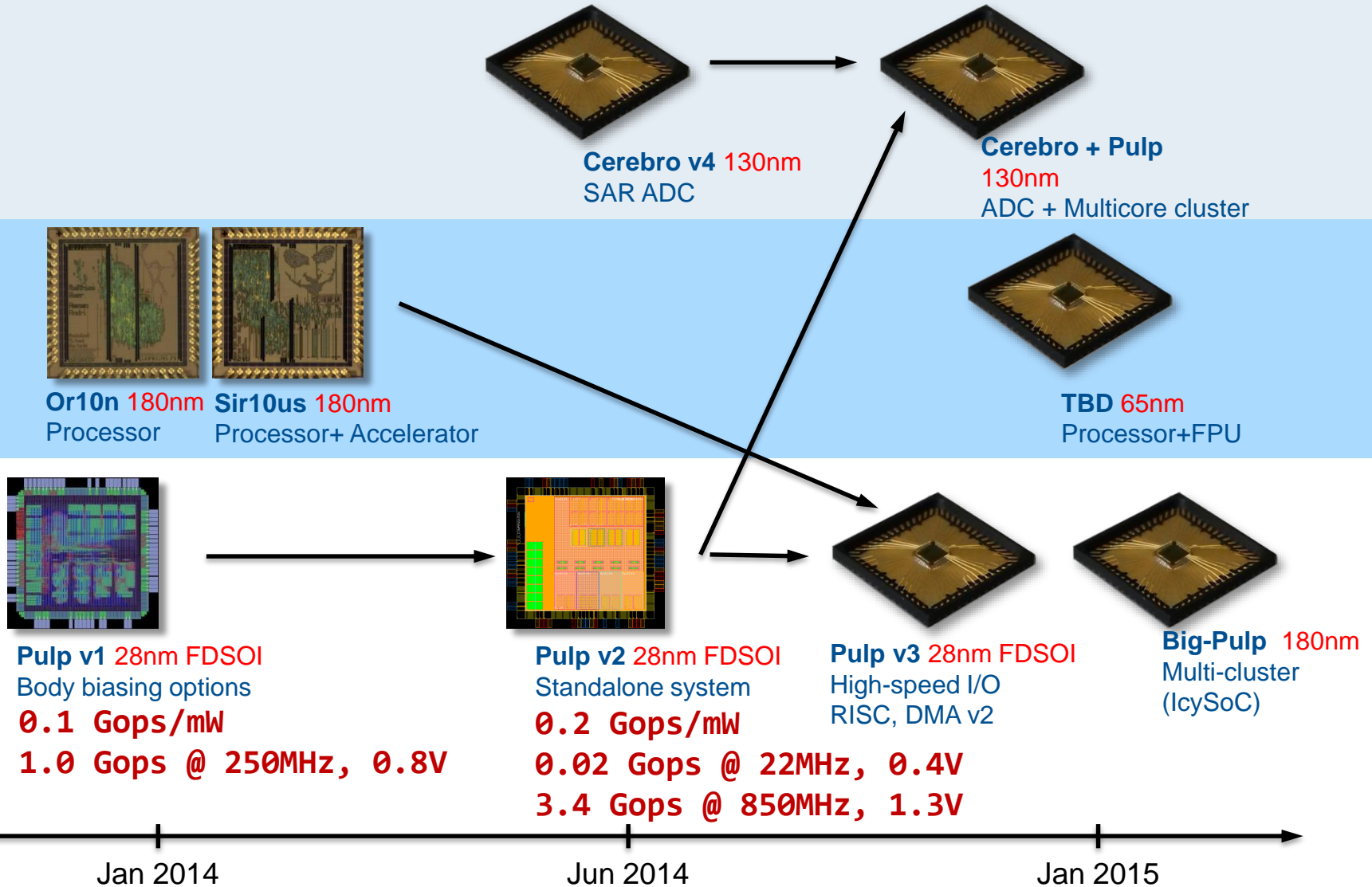
UP TO 220 GOPS/W IN A TYPICAL LEAKAGE SCENARIO!!!

Introducing FBB...



UP TO 1.5X ENERGY EFFICIENCY FOR HIGH WORKLOADS!!!

PULP Family Development



HARDWARE IPs

- PROCESSOR
- INTERCONNECT (LOCAL)
- MEMORY HIERARCHY
- MEMORY CONTROLLER
- HARDWARE
- ...

SOFTWARE

- COMPILER/TOOLCHAIN
- PROGRAMMING MODELS
- TIME

Building an open-source ecosystem for exploring (with silicon!) next-generation parallel computing platforms

SILICON

- OPTIMIZATION FLOW
- IMPLEMENTATION FLOW
- VERIFICATION FLOW
- FULL CUSTOM IPS
- ...



VALIDATION

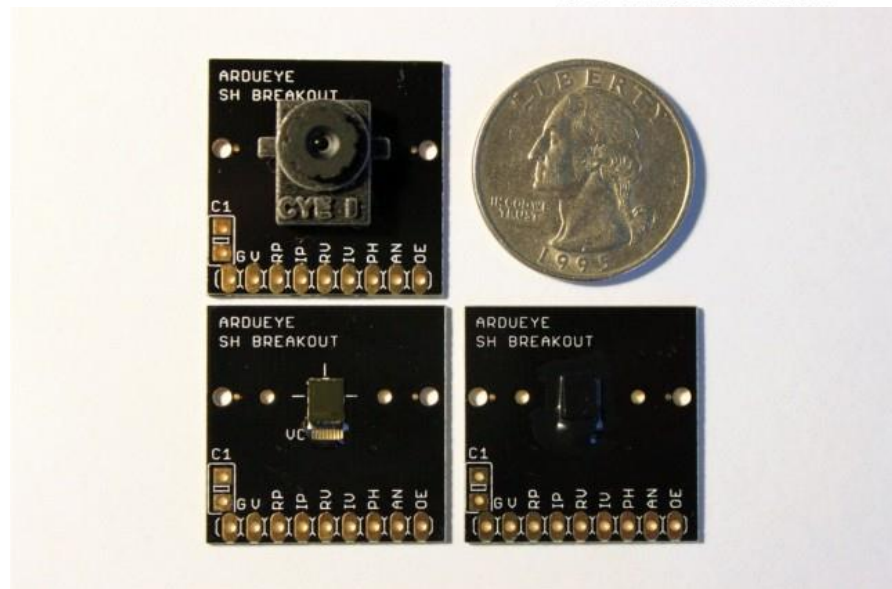
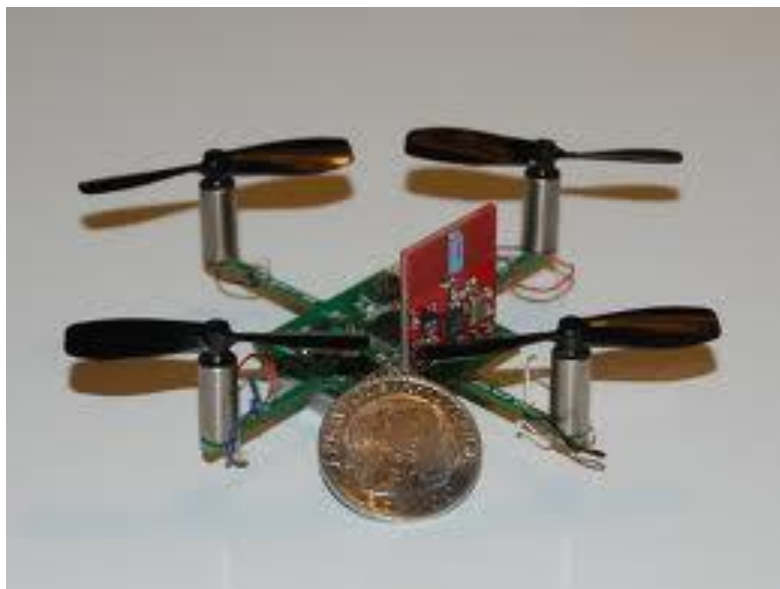
life.augmented

BUT ...

- SUPPORT FOR DEBUG
- SUPPORT FOR PROFILING
- DESIGN FOR TESTABILITY
-

- EMULATION PLATFORM (FPGA)
- BENCHMARKS
- REGRESSION TESTS
- ...

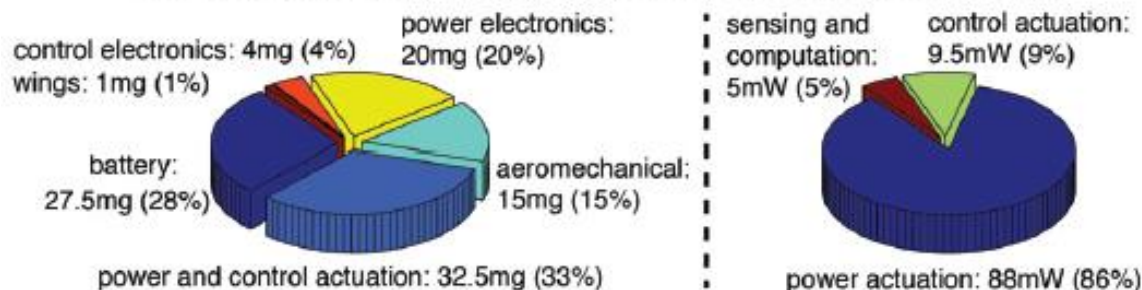
To do what?



112x112pixel (300uW)



component mass (100mg total) and power (102.5mW total)



[Wood13]

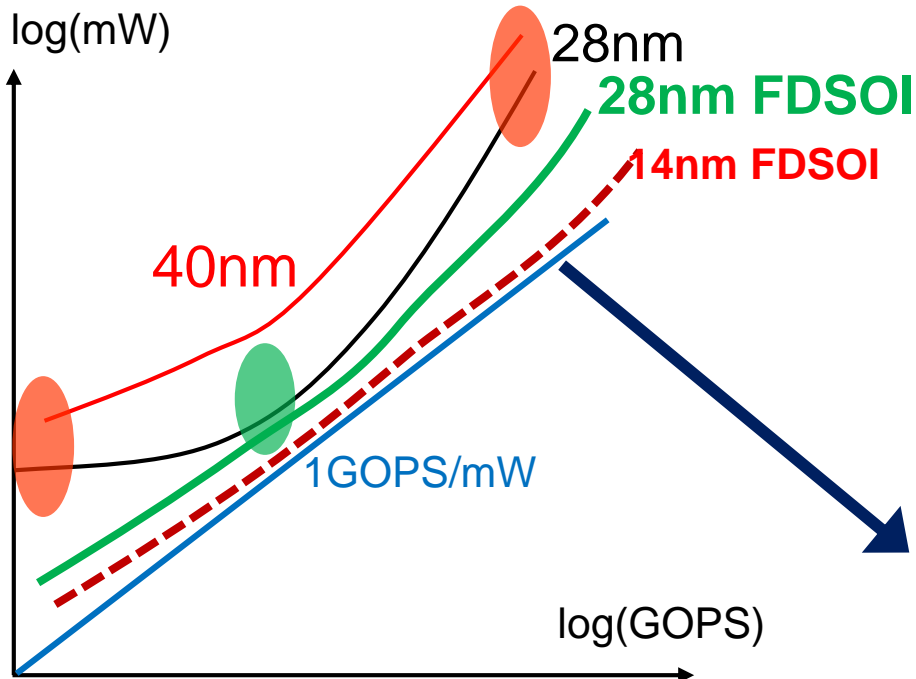
@60fps 0.76MPx/s → with 1KOPS/pixel we need 0.75Gops!

The Grand Challenge: Energy proportionality

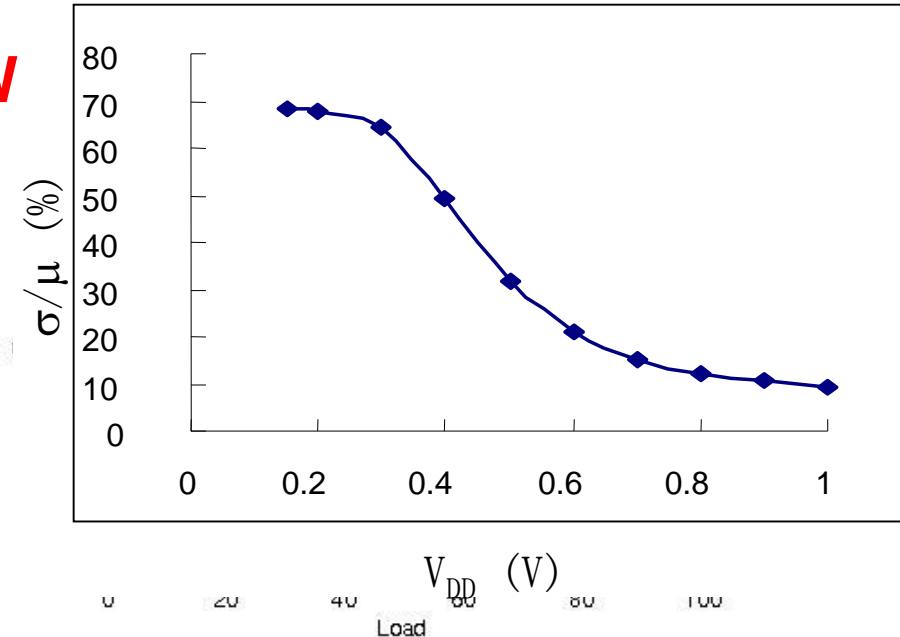


3,200MOPS/W – 30KW

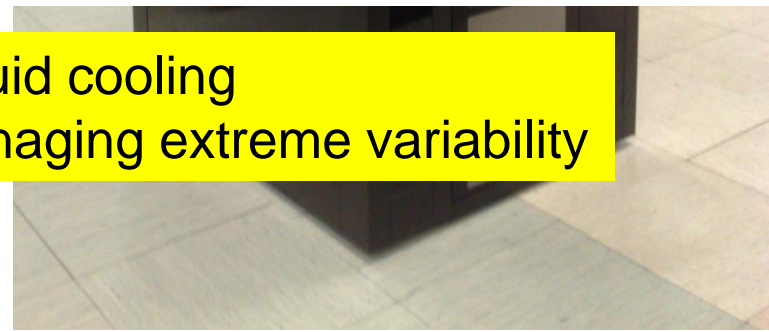
PULP1-2 32b, Kops-Tops @ pJ/W



Goal: reduce «bending up» of energy curve at low & **high** perf!



+ Liquid cooling
+ Managing extreme variability



Thank you!



European Research Council

Multithermand AdG**Multiscale Thermal Management of Computing Systems**