PULP: A Multi-Core Platform for Micropower In-Sensor Analytics

OPRECOMP SUMMER SCHOOL

03.09.2019



¹Department of Electrical, Electronic and Information Engineering

Davide Rossi

davide.rossi@unibo.it



ETH zürich ²Integrated Systems Laboratory

A Very Short Review on CMOS power (and power minimization)



Summary of Power Dissipation Sources

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

- α switching activity
- C_L load capacitance
- V_{swing} voltage swing
- *f* − frequency

- I_{DC} static current
- *I*_{leak} leakage current

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



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
- Micro-Architecture
- 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 (n²)
 - ... 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 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-transistions



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 Vdd 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 Vdd = 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_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 curent)
- .. and many more







Near-Sensor Processing



Computing fot the Internet of Things



		Near-S	Sensor Pro	cessing	
		INPUT (COMPUTATION	IAL OUTPUT	COMPRESSION
	Image	BANDWIDTH	DEMAND	BANDWIDTH	FACTOR
	Tracking: [*Lagroce2014]	80 Kbps	1.34 GOPS	0.16 Kbps	500x
	Voice/Sound	k			
	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
	Extreme	ly compact o	utput (single	index, alarm, s	<mark>ignature)</mark>
	Comput	ational powe	<mark>r of ULP μCo</mark>	ntrollers is not	enough
	Parallel	worloads			
	PULP				Davide Rossi 06.09.2019 16

PULP: pJ/op Parallel ULP computing



Parallel + Programmable + Heterogeneous ULP computing 1mW-10mW active power



Parallel Ultra Low Power



Minimum Energy Operation



Near-Threshold Computing (NTC):

- Don't waste energy pushing devices in strong inversion
- Recover performance with parallel execution
- Aggressively manage idle power (switching, leakage)
- Manage Process and temperature variations in NT





Targe [et Worl MOPS	kload]	1-Core Energy Efficiency (ideal) [MOPS/mW]	4-Cores Energy Efficiency (ideal) [MOPS/mW]	Ratio	
	100		43	55	1.3x	
	200		33	50	1.5x	
	400		18	43	2.4x	J



*Measured on our first prototype



Going faster allows to integrate system power over a smaller period

The main constraint here is the power envelope



Building PULP



Building PULP

SIMD + MIMD + sequential



1 Cycle Shared Multi-Banked L1 Data Memory + Low Latency Interconnect

"GPU like" shared memory \rightarrow low overhead data sharing

- Near Threshold but parallel \rightarrow Maximum Energy efficiency when Active
- + strong power management for (partial) idleness



D. Rossi *et al.*, «PULP: A Parallel Ultra Low Power Platform for Next Generation IoT Applications," in *HOT CHIPS* 2015.

Near threshold FDSOI technology



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



Near Threshold + Body Biasing Combined



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!



Selective, Fine Grained Body Biasing

- 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 (in LVT) RBB (in RVT)





D. Rossi *et. al.*, «A 60 GOPS/W, −1.8V to 0.9V body bias ULP cluster in 28nm UTBB FD-SOI technology», in Solid-State Electronics, 2016.

PULPv1



CHIF	PFEATURES
Technology	28nm FDSOI (RVT)
Chip Area	3mm ²
# Cores	4xOpenRISC
I\$	4x1kbyte (private)
TCDM	16 kbyte
L2	16 kbyte
BB regions	6
VDD range	0.45V-1.2V (SRAM: 0.55V-1.2V)
VBB range	-1.8V - +0.9V
Perf. Range	1 MOPS-1.9GOPS
Power Range	100 μW-127 mW
Peak Efficiency	60 GOPS/W@20 MOPS, 0.55V

PULPv1 issues:

- 1) World record energy efficiency (60 MOPS/mW), but @ too small performance (20 MOPS)
- 2) SRAM limits voltage scalability (very well known problem...)



D. Rossi *et. al.*, «A 60 GOPS/W, −1.8V to 0.9V body bias ULP cluster in 28nm UTBB FD-SOI technology», in Solid-State Electronics, 2016.

The Memory bottleneck

PULPv1 POWER BREAKDOWN @ BEST ENERGY POINT:



PULPv1 issues:

- 1) World record energy efficiency (60 MOPS/mW), but @ too small performance (20 MOPS)
- 2) SRAM limits voltage scalability (very well known problem...)
- 3) SRAM forms a huge bottleneck for energy efficiency (>60% of total power)



ULP (NT) Bottleneck: Memory

- "Standard" 6T SRAMs:
 - High VDDMIN
 - Bottleneck for energy efficiency
 - >50% of energy can go here!!!
- 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)
 - High technology portability
 - Major area overhead 4x → 2.7x with controlled placement





A. Teman *et.al.*, 'Power, Area, and Performance Optimization of Standard Cell Memory Arrays Through Controlled Placement', in ACM TDAES, May 2016

PULPv2

		CHIP	FEATURES
<u> </u>		Technology	28nm FDSOI (LVT)
		Chip Area	3mm ²
		# Cores	4xOpenRISC
		I\$ (SCM)	4x1kbyte (private)
	L1 MEM (SRAM)	TCDM	32 + 8 Kbyte
2		L2	64 kbyte
	BBMUXes	BB regions	10
M	CORES	VDD range	0.3-1.2V (0.5-1.2V)
		VBB range	0V-2V
	I\$ (SCM)	Perf. Range	1 MOPS - 4 GOPS
		Power Range	10μW - 300 mW
		Peak Efficiency	192 GOPS/W@0.5V



D. Rossi *et al.*, "Energy-Efficient Near-Threshold Parallel Computing: The PULPv2 Cluster," in *IEEE Micro*, Sep./Oct. 2017.

I\$: a Look Into 'Real Life' Applications

SCM-BASED I\$ IMPROVES EFFICIENCY BY ~2X ON SMALL BENCHMARKS, BUT...

Applications on PULP



- 2) Capacity miss (with small caches)
- 3) Jumps due to runtime (e.g. OpenMP, OpenCL) and other function calls



OpenMP for PULP

- OpenMP on PULP:
 - Lightweight implementation on top of a bare-metal runtime (custom GCC libgomp)
 - A subset of OpenMP 3.0 supported (e.g. no tasking)
 - Power management embedded in the runtime (transparent to the end-user)
- Architectural Implications
 - #pragmas are translated into runtime function calls
 - Function calls cause I\$ cache pollution
 - Call to OpenMP functions feature intrinsic overhead





A. Marongiu et. al., "An OpenMP Compiler for Efficient Use of Distributed Scratchpad Memory in MPSoCs," in IEEE Transactions on Computers, Feb. 2012.

The Solution: Shared I\$

- Share instruction cache
 - OK for data parallel execution model
 - Not OK for task parallel execution model, or very divergent parallel threads
- Architectures
 - SP: single-port banks connected through a readonly interconnect
 - Pros: Low area overhead
 - Cons: Timing pressure, contention
 - MP: Multi-ported banks
 - Pros: High efficiency
 - Cons: Area overhead (several ports)
 - Hierarchical Cache:
 - Pros: P&R friendly
 - Cons: Slightly slower than the others.

Private I\$ (traditional)



Shared Single-Port I\$





I. Loi, *et.al.*, "The Quest for Energy-Efficient I\$ Design in Ultra-Low-Power Clustered Many-Cores," in IEEE Transactions on Multi-Scale Computing Systems, In Press.

The Solution: Shared I\$

Multi-Port Shared I\$



- MP: Multi-ported banks
 - Pros: High efficiency
 - Cons: Area overhead (several ports)

Hierarchitcal Shared I\$



- Hierarchical Cache:
 - Pros: P&R friendly
 - Cons: Slightly slower than the others.



I. Loi, *et.al.*, "The Quest for Energy-Efficient I\$ Design in Ultra-Low-Power Clustered Many-Cores," in IEEE Transactions on Multi-Scale Computing Systems, In Press.

Synthetic Benchmarks



PULP

I. Loi, *et.al.*, "The Quest for Energy-Efficient I\$ Design in Ultra-Low-Power Clustered Many-Cores," in IEEE Transactions on Multi-Scale Computing Systems, In Press.

Real-Life Applications



PULP

I. Loi, *et.al.*, "The Quest for Energy-Efficient I\$ Design in Ultra-Low-Power Clustered Many-Cores," in IEEE Transactions on Multi-Scale Computing Systems, In Press.



HW Synchronizer: Impact on OpenMP primitives

- Cost of OpenMP runtime reduced by more than one order of magnitude
- Better scalability with number of cores



F. Glaser, et. al., "Hardware-Accelerated Energy-Efficient Synchronization and Communication for Ultra-Low-Power Tightly Coupled Clusters," DATE 2019.

Extending RISC-V for NSP

<32-bit precision -> SIMD2/4 opportunity

- 1. HW loops and Post modified LD/ST
- 2. Bit manipulations
- 3. Packed-SIMD ALU operations with dot product
- 4. Rounding and Normalizazion
- 5. Shuffle operations for vectors
 - V1 Baseline RISC-V RV32IMC HW loops
 - V2 Post modified Load/Store Mac
 - **V3** SIMD 2/4 + DotProduct + Shuffling
 - Bit manipulation unit Lightweight fixed point

Small Power and Area overhead



M. Gautschi et al., "Near-Threshold RISC-V Core With DSP Extensions for Scalable IoT Endpoint Devices," in IEEE TVLSI, Oct. 2017.

ISA Extensions improve performance

for (i = 0; i < 100; i++)
 d[i] = a[i] + b[i];</pre>

Baseline

 mv
 x5, 0

 mv
 x4, 100

 Lstart:
 Auto-incr load/store

 lb
 x2, 0(: mv
 x5, 0

 lb
 x3, 0(: mv
 x4, 100

 addi
 x10,x1
 Lstart:

 addi
 x11,x1
 lb
 x2, 0(

 add
 x2, x3
 lb
 x3, 0(
 lb
 x2, 0(x10!)

 sb
 x2, 0(:
 addi
 x4, x4
 add
 x2, x3
 add
 x2, x3, x2

 addi
 x4, x5
 bne
 x4, x5, Lstart
 pv.add.b
 x2, x3, x2

 Lend:
 sw x2, 0(x12!)

11 cycles/output 8 cycles/output 5 cycles/output 1,25 cycles/output

PULP

What About FP?

- The adoption of floating point formats requiring a lower number of bits (*smallFloats*) can reduce execution time and energy consumption
 - Simpler Logic (smaller pj/op)
 - Vectorization (smaller execution time)
- SW support:
 - *flexFloat* library for emulation of smallFloat format
 - Automatic exploration tool to tune precision of individual FP operations for the required application accuracy
- Hardware support:
 - smallFloat scalar and vector operations
 - Casting operations to move data through different FP types

Tagliavini et. A.I, «A Transprecision Floating-Point Platform for Ultra-Low Power Computing», DATE 2018

Process & Temperature Variations in NTC

l

Body Bias-Based Performance Monitoring and Control

Goal:

- Compensate the effects of external factors like ambient temperature
- Reduce PVT Margings at design time
- Improve Energy Efficiency

→ Mixed Hardware/Software control loop exploiting on-chip frequency measurement (PMB)

ETH zürich

ÉCOLE POLYTECHNIQUE Fédérale de lausanne

D. Rossi et al., "A Self-Aware Architecture for PVT Compensation and Power Nap in Near Threshold Processors," in IEEE Design & Test, Dec. 2017.

Davide Rossi | 06.09.2019 | 42

life.augmented

Dynamic PVT Compensation

A Full SoC Perspective

Putting All together: Mr. Wolf IoT Processor

- The SoC is an advanced microcontroller based on an ultra-low-power 32-bit RISC-V processor
 - 512kB of L2 Memory
- Rich set of peripherals:
 - QSPI, SPI, I2C, I2S
 - HyperRam + HyperFlash
 - Camera Interface
 - JTAG (Debug), GPIOs, PWM, ROM
 - RTC, Interrupt controller
- Autonomous IO DMA Subsystem
- Advanced power management
 - 2 Switchable Power Domains
 - 2 low-power FLLs (IO, SoC, Cluster)
 - On-chip DCDC and LDO
 - State-Retentive L2 Memory
- Parallel Programmable Accelerator:
 - 8x DSP enhanced RISC-V processors
 - 64kB of Shared L1 Memory
 - 2x Shared Floating Point Units
 - Shared latch-based Instruction Cache
 - High performance DMA(L2<->L1)
 - Event Unit supporting fast synchronization among cores

Mr. Wolf Chip Results

Technology	CMOS 40nm LP
Chip area	10 mm²
VDD range	0.8V - 1.1V
Memory Transistors	576 Kbytes
Logic Transistors	1.8 Mgates
Frequency Range	32 kHz – 450 MHz
Power Range	72 μW – 153 mW

Power Managent (DC/DC + LDO)	VDD [V]	Freq.	Power
Deep Sleep	0.8	n.a.	72 µW
Ret. Deep Sleep	0.8	n.a.	76.5 - 108 μW
SoC Active	0.8 - 1.1	32 kHz - 450 MHz	0.97 - 38 mW
Cluster Active	0.8 - 1.1	32 kHz - 350 MHz	1.6 - 153 mW

A. Pullini, D. Rossi, I. Loi, G. Tagliavini and L. Benini, "Mr.Wolf: An Energy-Precision Scalable Parallel Ultra Low Power SoC for IoT Edge Processing," in IEEE Journal of Solid-State Circuits, vol. 54, no. 7, pp. 1970-1981, July 2019.

Mr. Wolf XPULP Extensions Performance

■ RISCY (no opt) ■ RISCY (compiler opt.) ■ Intrinsic functions

Efficient Synchronization and Power Management

GOALS:

- Reduce cost of parallelization of data-parallel programming models (OpenMP, OpenCL)
- Shut-down (clock gating) processors during idle periods

RESULTS

~15x latency and energy reduction for a barrier

30%-50% latency and energy reduction for a mutex

F. Glaser et.al. "Hardware-Accelerated Energy-Efficient Synchronization and Communication for Ultra-Low-Power Tightly Coupled Clusters", DATE 2019

ARCHITECTURE

Mr. Wolf Parallel Speed-up

Mr. Wolf Computing Performance (Fixed+Float)

On 2D Matrix Multiplication

| 25.7.2018 | 50

Mr. Wolf Computing Efficiency (Fixed+Float)

On 2D Matrix Multiplication

Recovering Efficiency Through Flexible Customization

Closing The Accelerator Efficiency Gap with Agile Customization

Recovering Even More Efficiency

Fixed function accelerators have limited reuse... how to limit proliferation?

Learn to Accelerate

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

Image recognition [RussakovskyIMAGENET2014] Speech recognition [HannunARXIV2014]

Flexible acceleration: learned CNN weights are "the program"

PULP CNN Performance

Average performance and energy efficiency on a 32x16 CNN frame

PULPv3 ARCHITECTURE, CORNER: tt28, 25°C, VDD= 0.5V, FBB = 0.5V

Sub pJ/OP? Approximate Computing

ULP

From Approximate to Transprecision: YodaNN¹

- Approximation at the algorithmic side \rightarrow Binary weights
- BinaryConnect [Courbariaux, NIPS15]
 - Reduce weights from 12-bit to a binary value -1/+1
 - Stochastic Gradient Descent with Binarization in the Forward Path

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

Learning large networks is still an issue with binary connect...

- Ultra-optimized HW is possible!
 - Power reduction because of arithmetic simplification (multipliers →Two's complement + muxes)
 - Major arithmetic density improvements
 - Area can be used for more energy-efficient weight storage
 - SCM memories for lower voltage \rightarrow E goes with 1/V²

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

Comparison sith SoA

Publication	Throughput [GOPS]	En.Eff. [GOPS/W]	Supply [V]	Area Effic. [GOPS/MGE]
Neuflow	320	490	1.0	17
Leuven	102	2600	0.5 - 1.1	64
Eyeriss	84	160	0.8 - 1.2	46
NINEX	569 <mark>†</mark>	1800	1.2	51
k-Brain	411	1930	1.2	109
Origami	196 /74	437/ 803	1.2/0.8	90 /34 10x
This Work*	1510 /55	9800/61200	1.2/0.6	1135 /41

*UMC 65nm Technology, post place & route, 25°C, tt, 1.2V/0.6V

Breakthrough for ultra-low-power CNN ASIC implementation

fJ/op in sight: manufacturing in Globalfoundries 22FDX

From Frame-based to Event-based

Back to System-Level

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 inteface
- **Mixed-signal event triggering** with an ULP imager with internal processing AMS capability

A Neuromorphic Approach for doing *nothing* VERY well

GrainCam Imager (FBK)

Pixel-level spatial-contrast extraction

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

SAMPLE

1-BIT MEMORY

BIT LINES

GrainCam Readout

120

100

80

60

40

20

Active

Idle

Power Consumption [uW]

Readout modes:

- IDLE: readout the counter of asserted pixels
- ACTIVE: sending out the addresses of asserted pixels (address-coded representation), according raster scan order

Event-based sensing: output frame data bandwidth depends on the external context-activity

Power Management

M. Rusci, et. al., "A Sub-mW IoT-Endnode for Always-On Visual Monitoring and Smart Triggering," in IEEE IoT Journal, 2017.

Event-Driven Binary Deep Network

Manuele Rusci et. al. «Always-ON visual node with a hardwaresoftware event-based binarized neural network inference engine», Computing Frontiers, 2018.

Results

Scenario	BNN with RGB input	Event-based BNN
Image Sensor Power Consumption	1.1mW @30fps	$100\mu W$ @50fps
Image Size	632446 bits	8192 bits
Image Sensor Energy for frame capture	66.7 μJ	$2 \mu J$
Transfer Time (4bit SPI @50MHz)	3.1 msec	0.04 msec
Transfer Energy (8.9mW @0.7V)	$28 \mu J$	$2 \mu J$
BNN Execution Time (168MHz)	81.3 msec	75.3 msec
BNN Energy consumption (8.9mW @0.7V)	$725 \ \mu J$	671 μJ
Total System Energy for Classification	$820 \ \mu J$	674 μJ

- 3 Classes:
 - Pedestrians
 - Bikes
 - Cars
- **84.6% vs. 81.6% Accuracy**

Statistics per frame	Frame-Based	Event-based
Idle (no motion)		
Sensor Power	1.1mW	$20\mu W$
Avg Sensor Data	19764 Bytes	-
Transfer Time	790µsec	-
Processing Time	3.02 msec	-
Avg Processor Power	1.45mW	0.3mW (sleep)
Detection		
Sensor Power	1.1mW	$60\mu W$
Avg Sensor Data	19764 Bytes	\sim 536 Bytes
Transfer Time	790µsec	$21.4 \mu sec$
Processing Time	3.47 msec	187.6μ sec
Avg Processor Power	1.57mW	0.511mW
Classification		
Sensor Power	2mW	$60\mu W$
Avg Sensor Data	79056 Bytes	1024 Bytes
Transfer Time	3.16 msec	41μ sec
Processing Time	81.3 msec	75.3 msec
Processor Energy	760 μ J	677 μJ

Event-Based Avg Power (mW) + Frame-Based Avg Power (mW)

Outlook and Conclusion

Conclusion

- Near-sensor processing → 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
- CNNs-based visual functions can be squeezed into mW envelope
 - Non-von-Neumann acceleration
 - Very robust to trans-precision computations (deterministic and statistical)
 - fJ/OP is in sight!
- More than CNN is needed (e.g. linear algebra, online optimizazion)
- Open Source HW & SW approach \rightarrow 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

The fun is just at the beginning...

