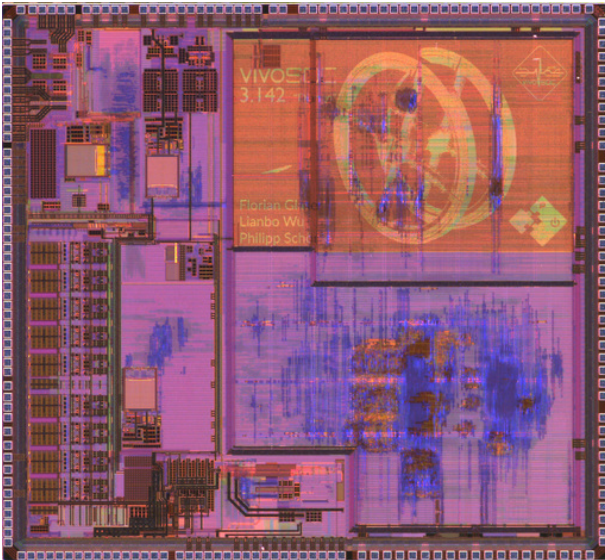


PULP PLATFORM

Open Source Hardware, the way it should be!

# RISC-V for IoT, the **PULP** experience



Frank K. Gürkaynak <kgf@ee.ethz.ch>

Digital Circuits and Systems Group  
ETH Zürich

 <http://pulp-platform.org>

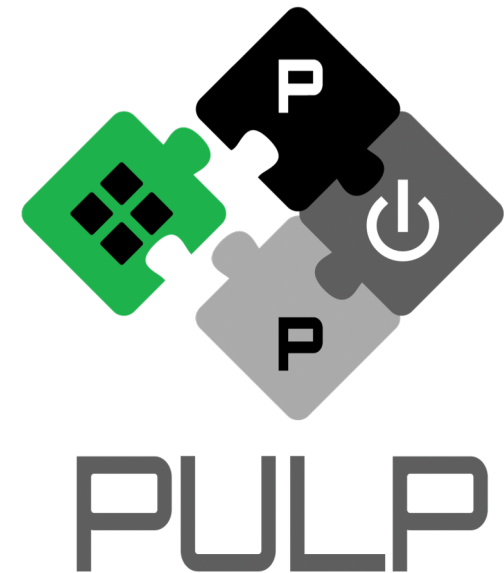
 [@pulp\\_platform](https://twitter.com/pulp_platform)

**ETH** zürich



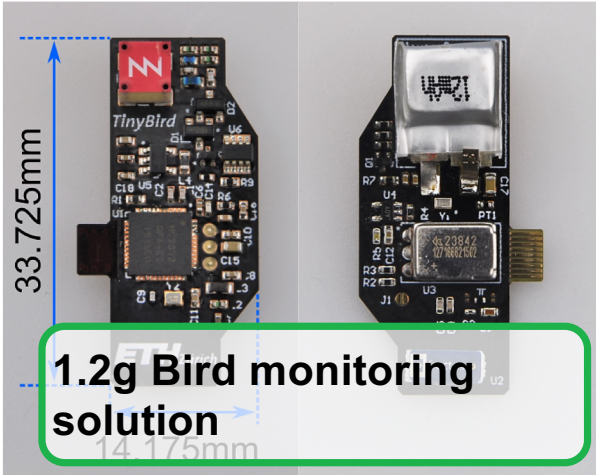
# How did we start in 2013?

- We wanted to design energy efficient computing systems
  - Equally efficient for IoT and HPC over a wide range
- Key points
  - Parallel processing
  - Near threshold computing
  - Efficient switching between operating modes
  - Making best use of technology
  - Heterogeneous acceleration
- Parallel Ultra Low-Power (**PULP**) platform was born



# IoT design @ Digital Circuits & Systems Group

**64mW autonomous  
Nano Drone**



**1.2g Bird monitoring  
solution**

**Sensor for Avalanche  
studies**



**Human touch  
powered sensor**



**LORA – GNSS  
localization, big & small**



**Wearable bio-signal  
processing systems**





# Who is behind PULP?



ETH Z



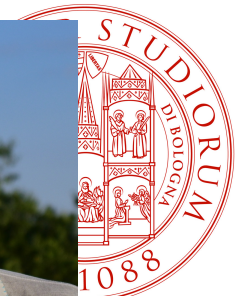
Frank



Prof. Luca Benini



Davide Rossi



STUDIORUM  
DI BOLOGNA



In total about 60 people work  
on projects related to PULP  
in Zurich and Bologna

<https://pulp-platform.org/team.html>



Europe

na  
world

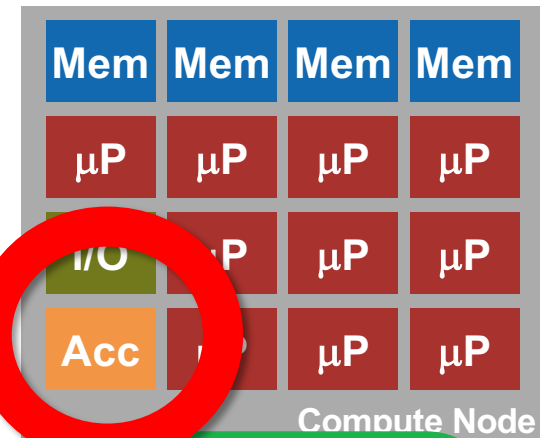
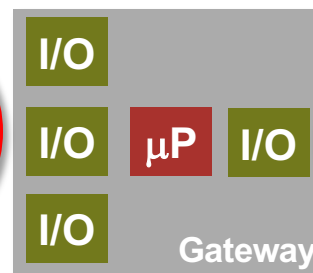
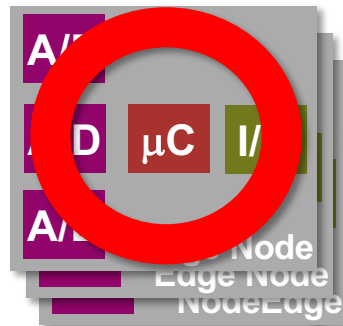


Chief Architect in STMicroelectronics (2009-2012)



# Too much to do and not enough resources

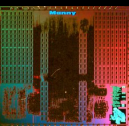
**SIGNAL** → **DATA** → **KNOWLEDGE**



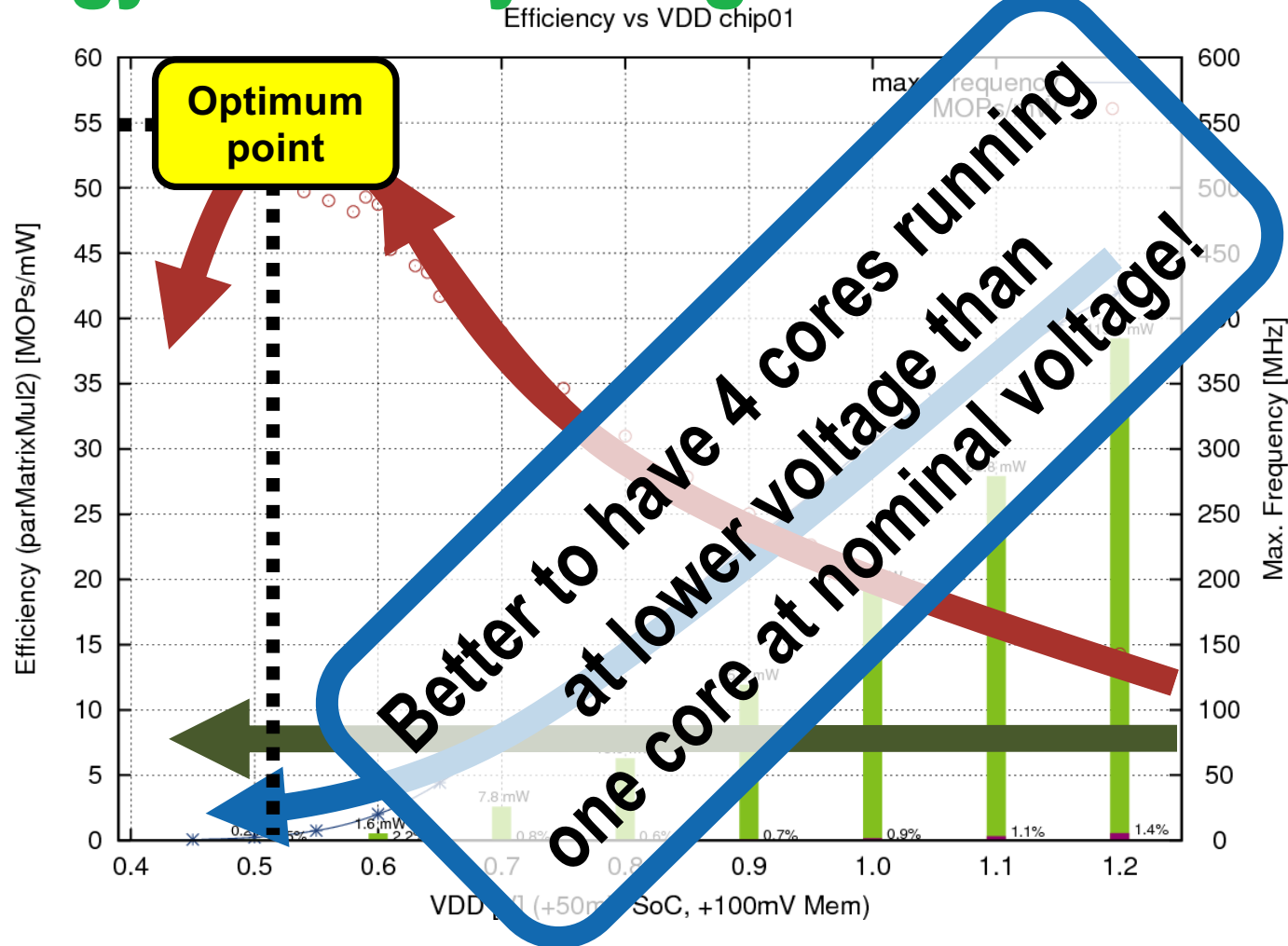
Energy efficient systems that can process more on edge nodes

Less data to transmit, large savings on energy needed for data transmissions

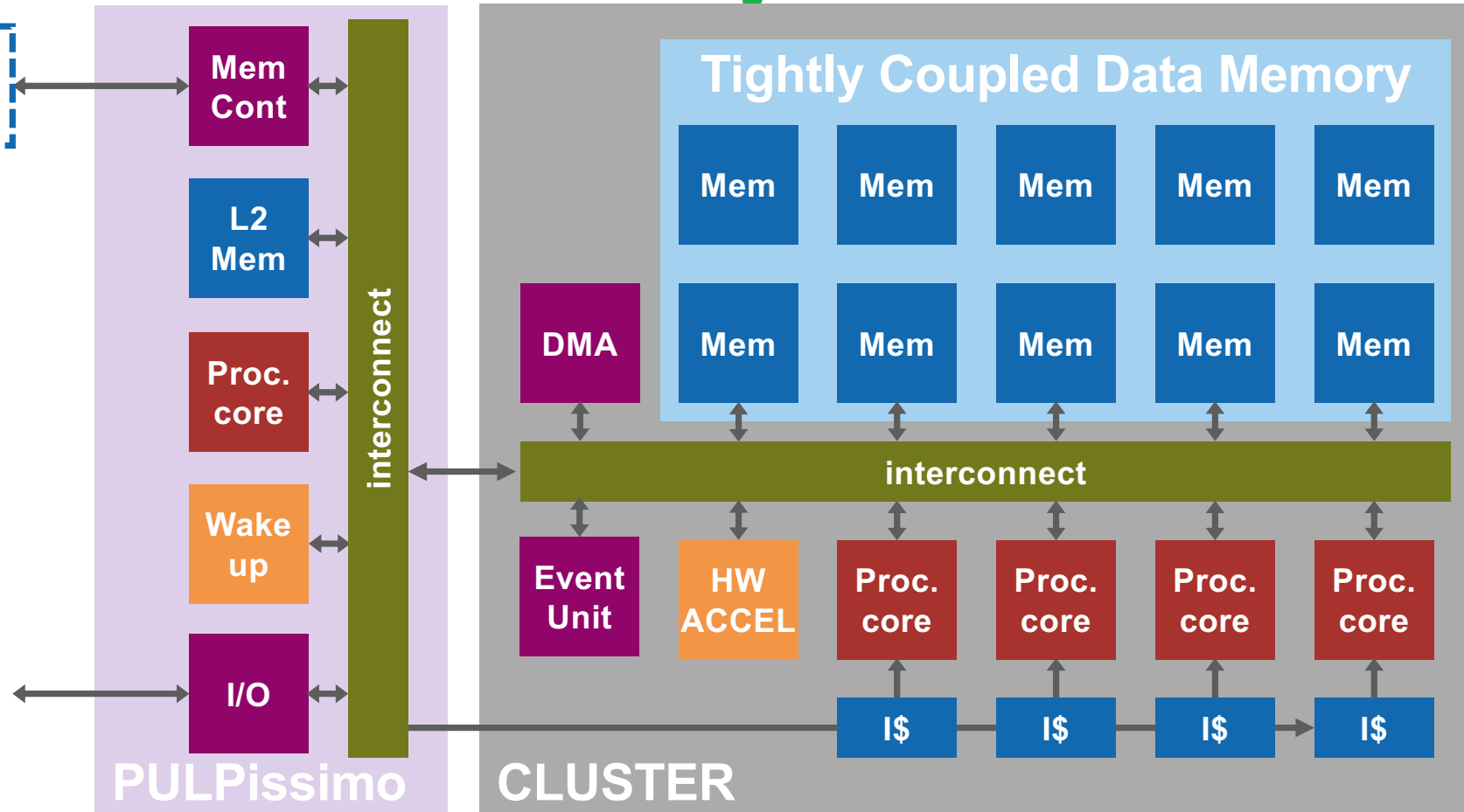
Accelerators for more efficient computing



# Energy efficiency is higher at near threshold



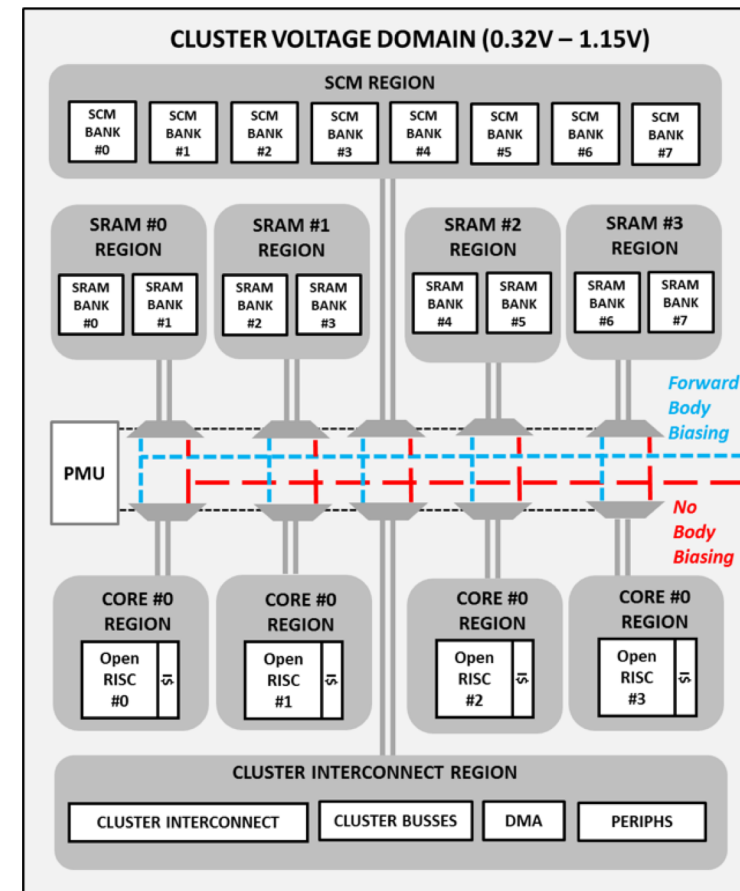
# Cluster based PULP systems





# ST28 FDSOI and GF22FDX designs with BB

- SoC partitioned in separate clock, power and body bias regions
- Cluster 1 Vdd, 10 BB regions
  - **Boost mode**: active + FBB
  - **Normal mode**: active + NO BB
  - **Idle mode**: clock gated + NO BB (in LVT) RBB (in RVT)
- SoC has 3 Vdd regions
  - Cluster, L2, Always on, IOs



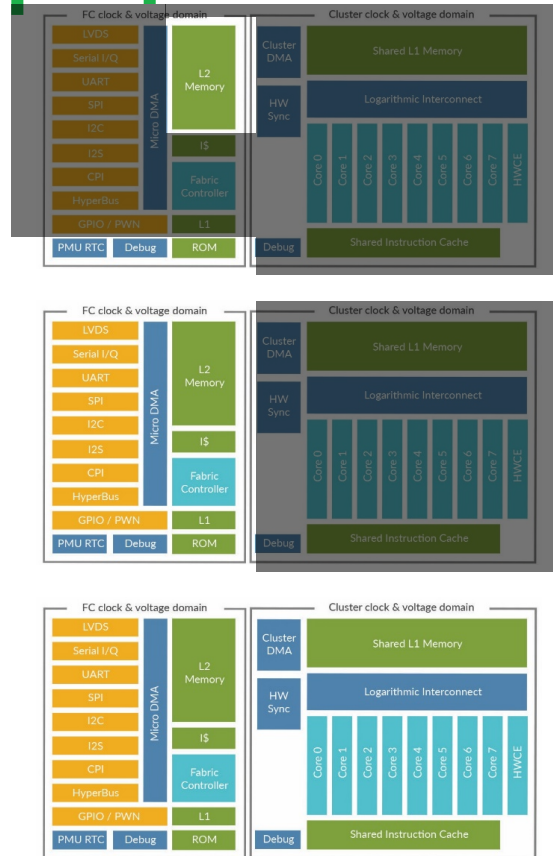
D. Rossi et. al., «A 60 GOPS/W, -1.8V to 0.9V body bias ULP cluster in 28nm UTBB FD-SOI technology», in Solid-State Electronics, 2016.

# Scaling proportional to computing demand

Duty Cycling

Coarse Grain Classification

Full Blown Analysis



1 to 50  $\mu$ W

## MCU sleep mode

- ✓ Low quiescent LDO
- ✓ Real Time Clock 32kHz only
- ✓ L2 Memory partially retentive

0.5 to 5 mW

## MCU active mode

- ✓ Embedded DC/DC, high current
- ✓ Voltage can dynamically change
- ✓ One clock gen active, frequency can dynamically change
- ✓ Systematic clock gating

5 to 50 mW

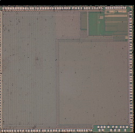
## MCU + Parallel processor active mode

- ✓ Embedded DC/DC, high current
- ✓ Voltage can dynamically change
- ✓ Two clock gen active, frequencies can dynamically change
- ✓ Systematic Clock Gating

Ultra fast switching time from one mode to another  
Ultra fast voltage and frequency change time



Highly optimized system  
level power consumption



# How PULP and RISC-V come together

- Initially we did not want to design our own processors
  - Wanted to use available processors (ARC, ARM.. )
  - It proved difficult to design systems that we could share with our collaborators

- Then we used OpenRISC cores (2013-2015)



- We had to completely redesign and optimize these cores

- We moved to RISC-V starting in 2015

- Adapted the decoder of our optimized OpenRISC core
  - Make use of a growing SW development environment
  - ETH is one of the founding members of the RISC-V foundation





# Our research is not implementing RISC-V cores

- **We develop efficient programmable architectures**
  - Processor cores of various capabilities are required for that
- **We need efficient implementations of cores for our research**
  - To produce relevant results, our cores have to perform **as good as other solutions**
  - We ended up spending quite an effort to make sure they perform really well
- **Processor core alone is not enough**
  - You need peripherals, interconnect solutions, programming support...
- **PULP Platform** provides us a playground for our research
  - And we share it as open source



# RISC-V cores developed by PULP team

32 bit			64 bit
Low Cost Core	Core with DSP support	Core for Streaming	Linux capable Core
<ul style="list-style-type: none"> <li>■ IBEX                             <ul style="list-style-type: none"> <li>■ Zero-riscy</li> <li>■ Micro-riscy</li> </ul> </li> <li>■ 2 options                             <ul style="list-style-type: none"> <li>■ RV32-ICM</li> <li>■ RV32-CE</li> </ul> </li> </ul> <p>ARM Cortex-M0+</p>	<ul style="list-style-type: none"> <li>■ RI5CY                             <ul style="list-style-type: none"> <li>■ RV32-ICMXF</li> <li>■ SIMD</li> <li>■ HW loops</li> <li>■ Fixed point manipulation</li> </ul> </li> </ul> <p>ARM Cortex-M4F</p>	<ul style="list-style-type: none"> <li>■ Snitch                             <ul style="list-style-type: none"> <li>■ RV32-IMAxFD</li> <li>■ Small int core</li> <li>■ Big 64bit FPU</li> <li>■ Extensions for streaming</li> </ul> </li> </ul> <p>novel</p>	<ul style="list-style-type: none"> <li>■ Ariane                             <ul style="list-style-type: none"> <li>■ RV64-ICMAFD</li> </ul> </li> </ul> <p>ARM Cortex-A55</p>

Maintained By LowRISC

Very Mature Core

Brand New Core

Frequent Updates

# RI5CY ISA extensions for performance

```
for (i = 0; i < 100; i++)
    d[i] = a[i] + b[i];
```

## Baseline

```
mv    x5, 0
mv    x4, 100
Lstart:
    lb    x2, 0(x10)
    lb    x3, 0(x11)
    addi  x10, x10, 1
    addi  x11, x11, 1
    add   x2, x3, x2
    sb    x2, 0(x12)
    addi  x4, x4, -1
    addi  x12, x12, 1
    bne   x4, x5, Lstart
```

**11 cycles/output**

## Auto-incr load/store HW Loop

```
mv    x5, 0
mv    x4, 100
Lstart:
    lb    x2, 0(x10!)
    lb    x3, 0(x11!)
    addi  x4, x4, -1
    add   x2, x3, x2
    sb    x2, 0(x12!)
    bne   x4, x5, Lstart
```

**8 cycles/output**

```
lp.setupi 100, Lend
    lb    x2, 0(x10!)
    lb    x3, 0(x11!)
    add   x2, x3, x2
Lend:    sb x2, 0(x12!)
```

**5 cycles/output**

## Packed-SIMD

```
lp.setupi 25, Lend
    lw    x2, 0(x10!)
    lw    x3, 0(x11!)
    pv.add.b x2, x3, x2
Lend: sw x2, 0(x12!)
```

**1,25 cycles/output**





# The extensions translate to real speed-ups

## ■ 8-bit convolution

- Open source DNN library

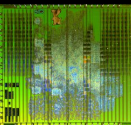
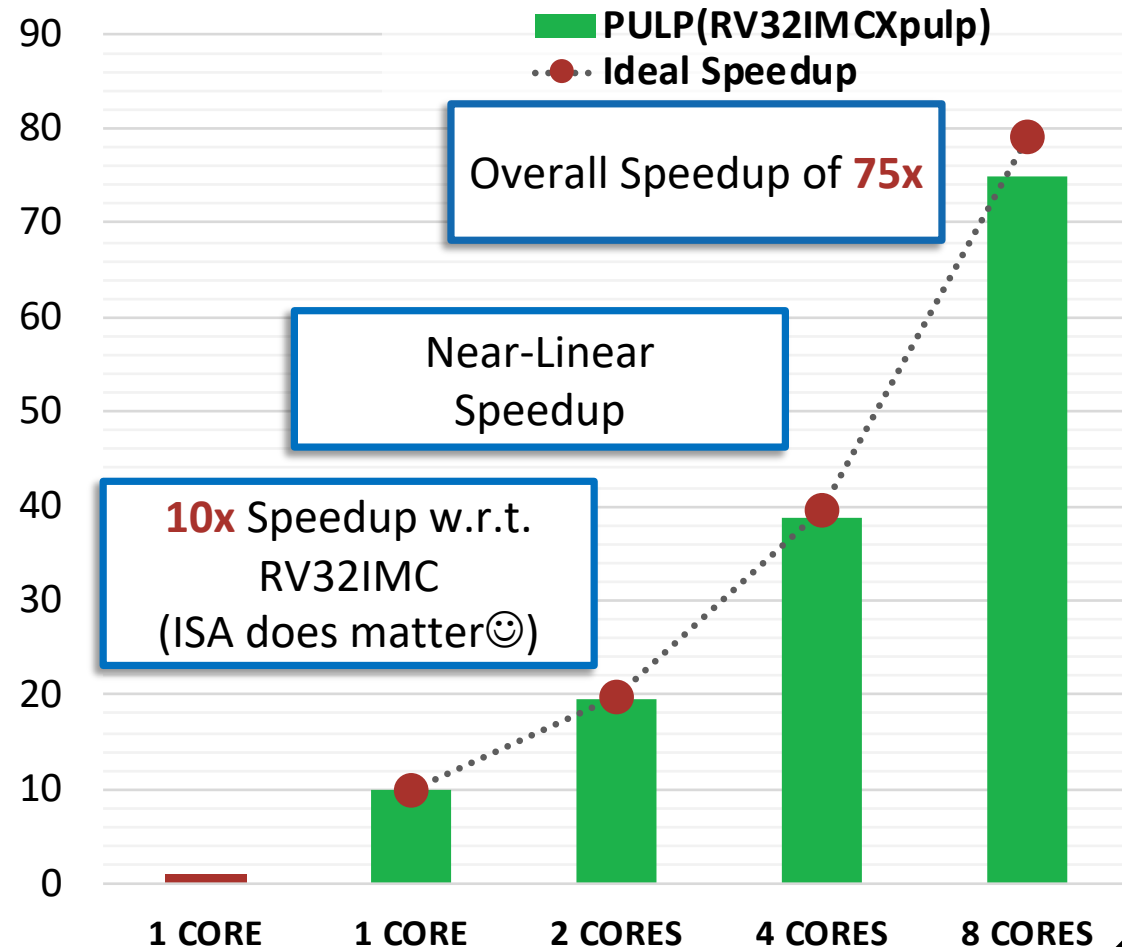
## ■ 10x through xPULP

- Extensions bring real speedup

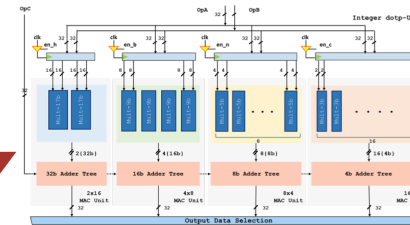
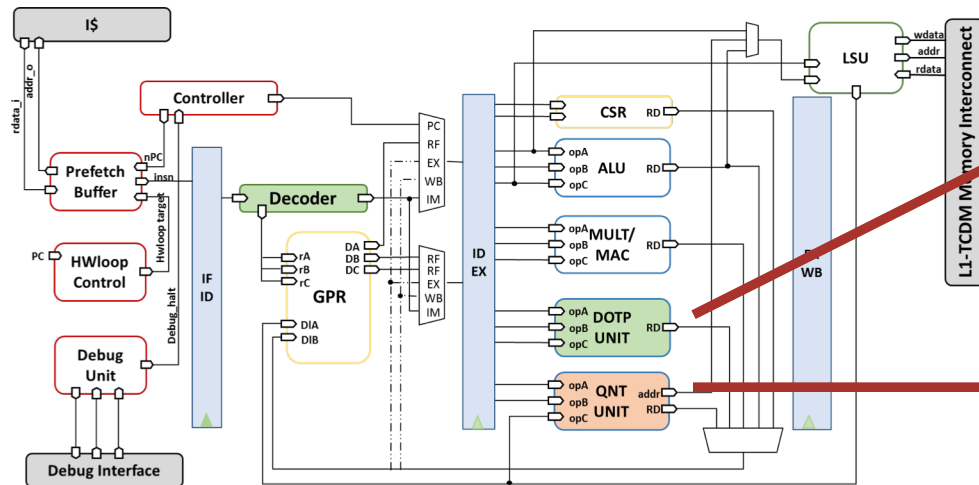
## ■ Near-linear speedup

- Scales well for regular workloads.

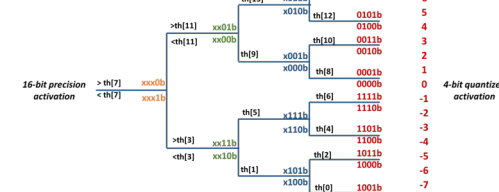
## ■ 75x overall gain



# RISC-V ISA Extensions for extreme quantization



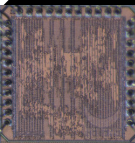
2-bit & 4-bit SIMD DOTP + OP Isolation



QNT UNIT: 2 Quantizations in 9 Cycles

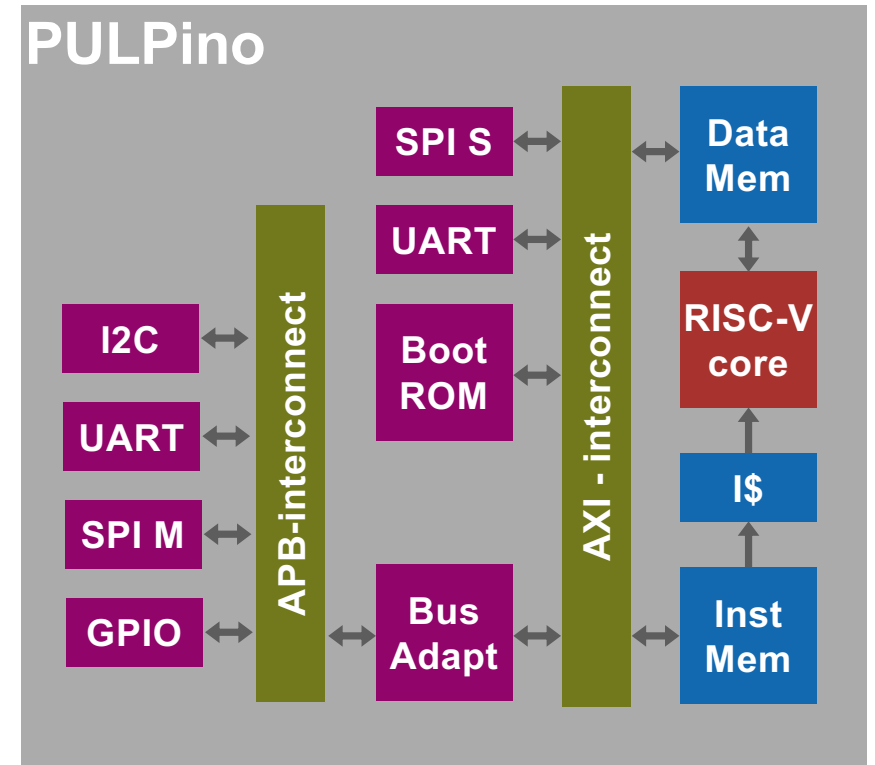
## Overheads (28nm FDX PULPissimo impl.):

- Area: ~11% (vs. Ri5CY)
- Timing Overhead: negligible (integrated in PULPissimo)
- 8-bit MatMul power overhead: 1.8% (integrated in PULPissimo)
- GP-app power overhead: 3.5% (integrated in PULPissimo)



# PULPino: Our first single core platform

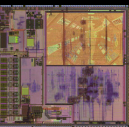
- **Simple design**
  - Meant as a quick release
- **Separate data and inst. mem**
  - Makes it easy in HW
  - Not meant as a Harvard arch.
- **Can use all our 32bit cores**
  - RI5CY, Zero/Micro-Riscy (Ibex)
- **Peripherals from other projects**
  - Any AXI and APB peripherals could be used





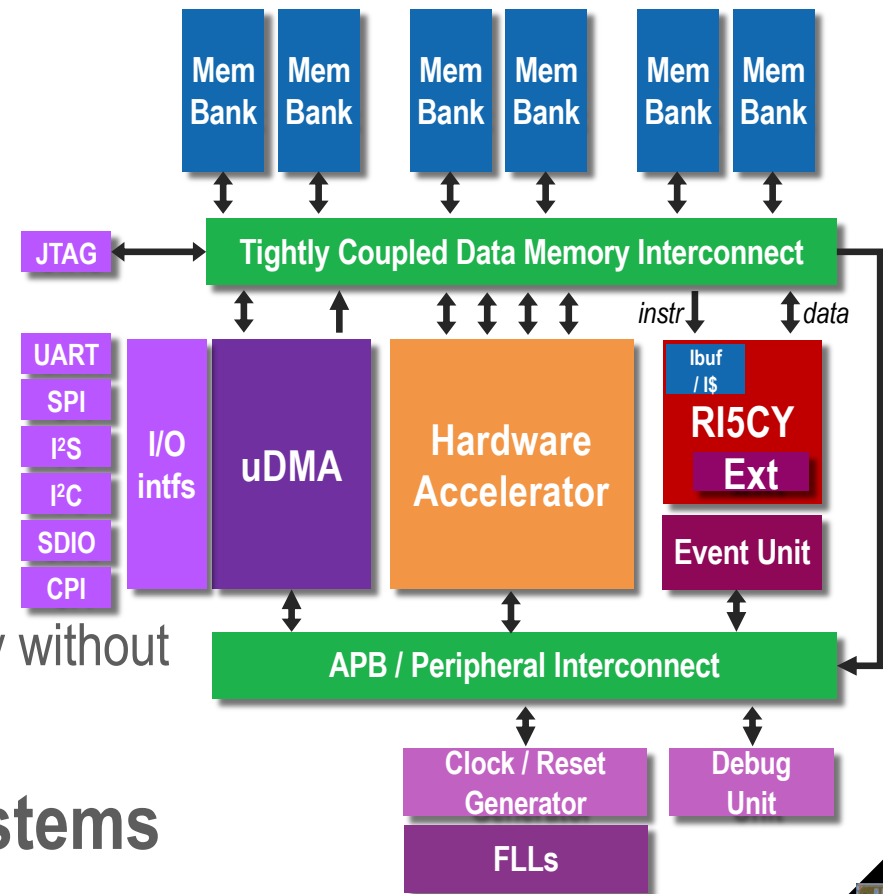
# What kind of acceleration?

- *Standard peripheral talking over AXI/APB (standard)*
- *Instruction set extensions (already discussed)*
- **Shared functional units**
  - Amortizes expensive extensions (FPU/DIV) between multiple units
- **Shared memory accelerators**
  - Our bread and butter, PULPopen, NTX
- **Cluster as an accelerator**
  - HERO, BigPULP, etc



# PULPissimo: The better single core platform

- **Shared memory**
  - Unified Data/Instruction Memory
- **Support for Accelerators**
  - Direct shared memory access
  - Programmed through APB bus
- **uDMA for I/O subsystem**
  - Can copy data directly from I/O to memory without involving the core
- **Used as controller in larger systems**



# Open PULP: our main cluster based system

Ext. Mem

Mem Cont

L2 Mem

Proc. core

Wake up

I/O

interconnect

PULPissimo

Tightly Coupled Data Memory

Mem

Mem

Mem

Mem

Mem

Mem

Mem

Mem

Mem

Mem

DMA

interconnect

Event Unit

HW ACCEL

Proc. core

Proc. core

Proc. core

Proc. core

I\$

I\$

I\$

I\$

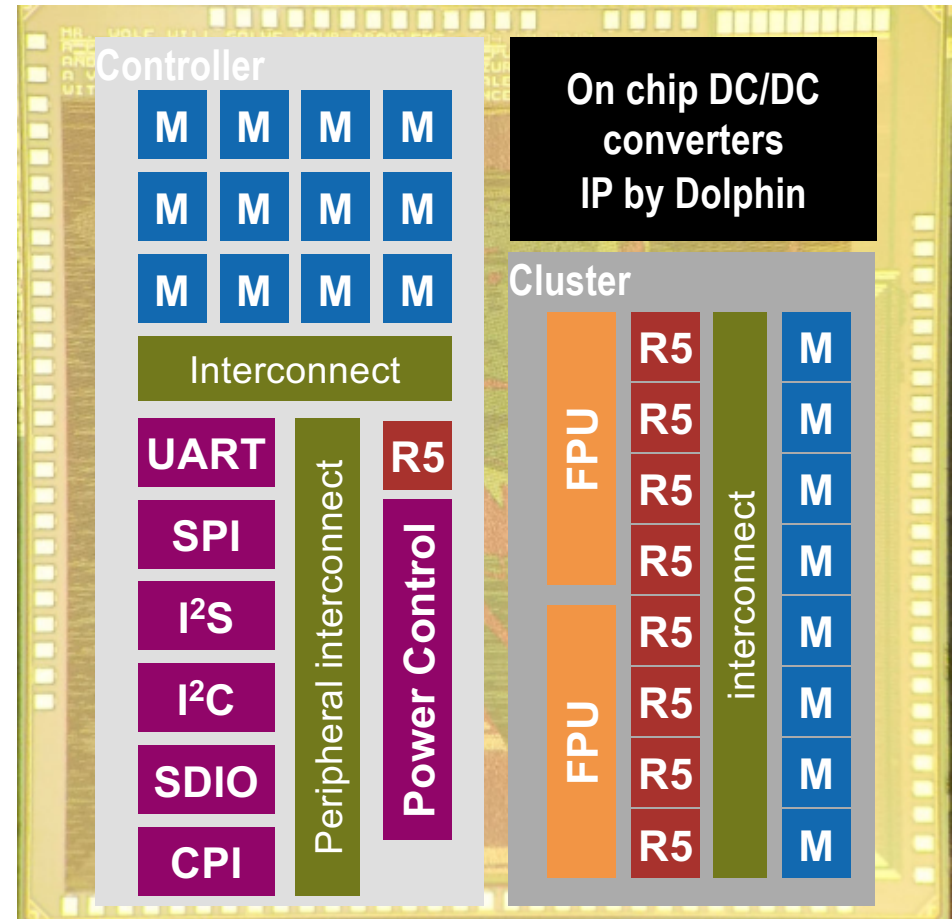
CLUSTER

ETH zürich



# Mr. Wolf (TSMC 40): 8+1 core IoT Processor

- **One cluster with**
  - 8 RISC-V cores
  - 2x shared FPU units
  - 64 kByte of TCDM
- **One controller with**
  - 512 kByte L2 RAM
  - Peripherals
- **On chip voltage regulators**
  - By Dolphin Integration



# PULP uses a permissive open source license

## ■ All our development is on GitHub

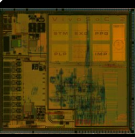
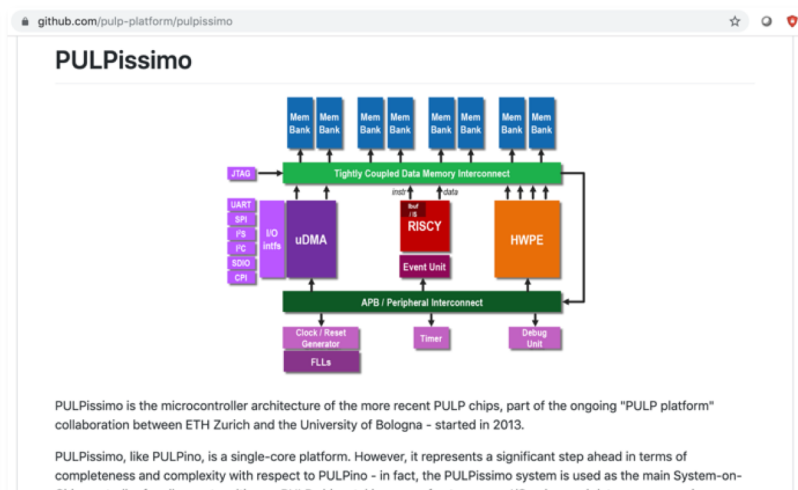
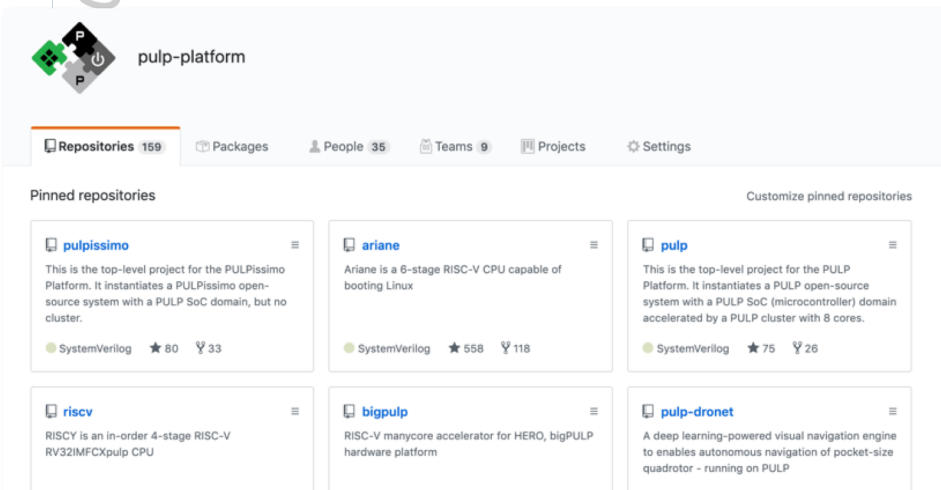
- HDL source code, testbenches, software development kit, virtual platform

<https://github.com/pulp-platform>



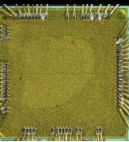
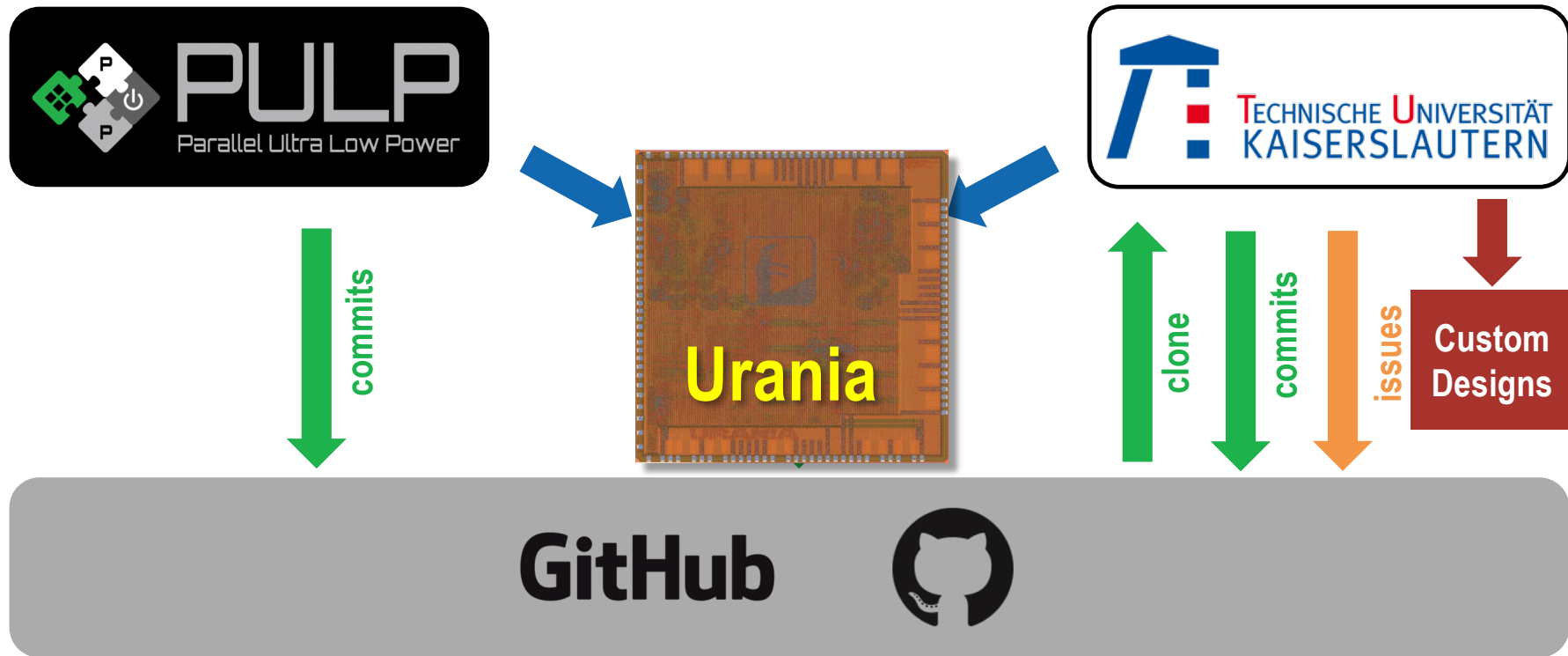
## ■ PULP is released under the permissive Solderpad license

- Allows anyone to use, change, and make products without restrictions.

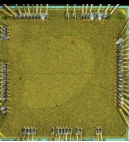
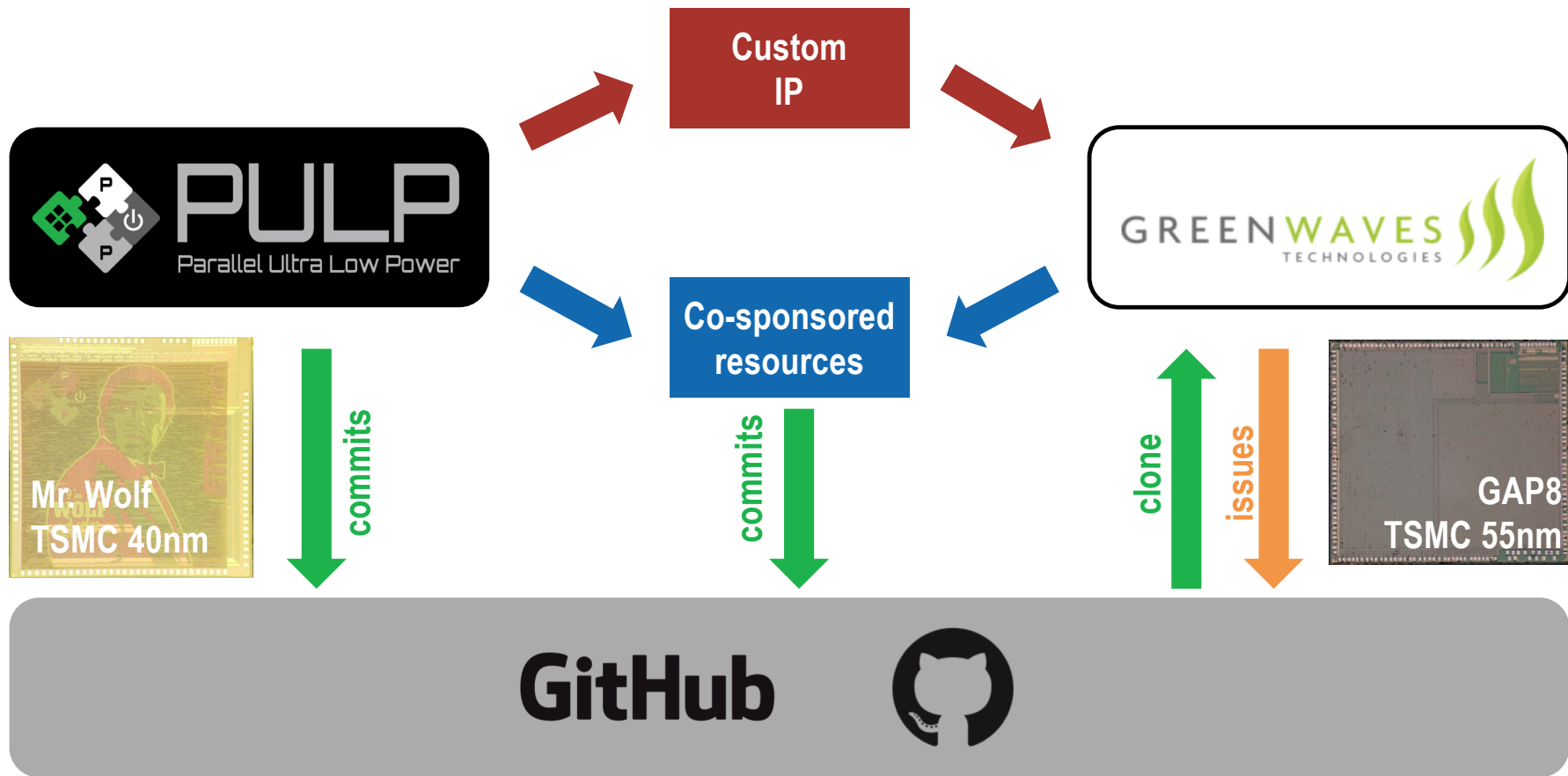




# Open source collaboration scheme explained

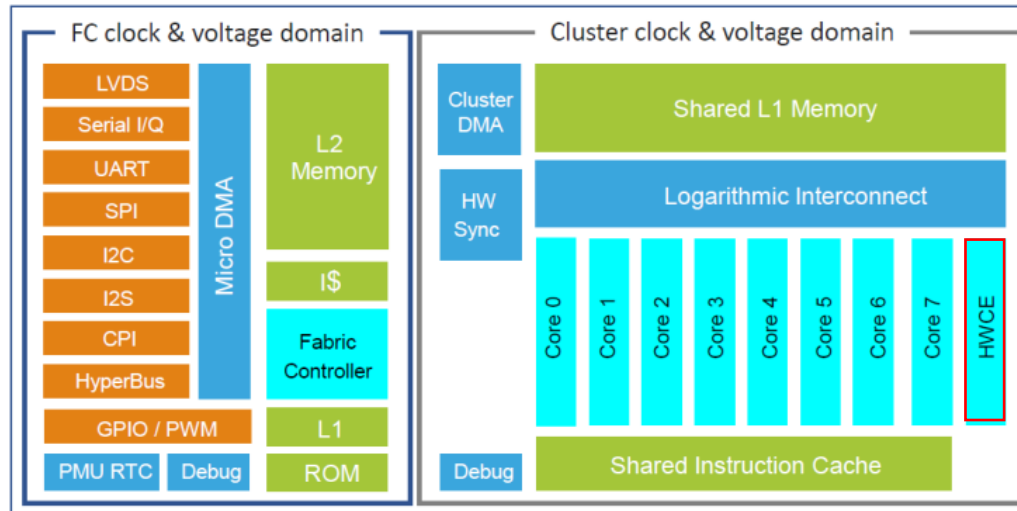


# Open source collaboration scheme explained



# PULP cluster+MCU+HWCE: GWT's GAP8

Two independent clock and voltage domains, from 0-133MHz/1V up to 0-250MHz/1.2V



## MCU Function

Extended RISC-V core  
Extensive I/O set  
Micro DMA  
Embedded DC/DC converter  
Secured execution

## Computation engine

8 extended RISC-V cores  
Fully programmable  
Efficient parallelization  
Shared instruction cache  
Multi channel DMA  
HW synchronization  
→ HW convolution Engine

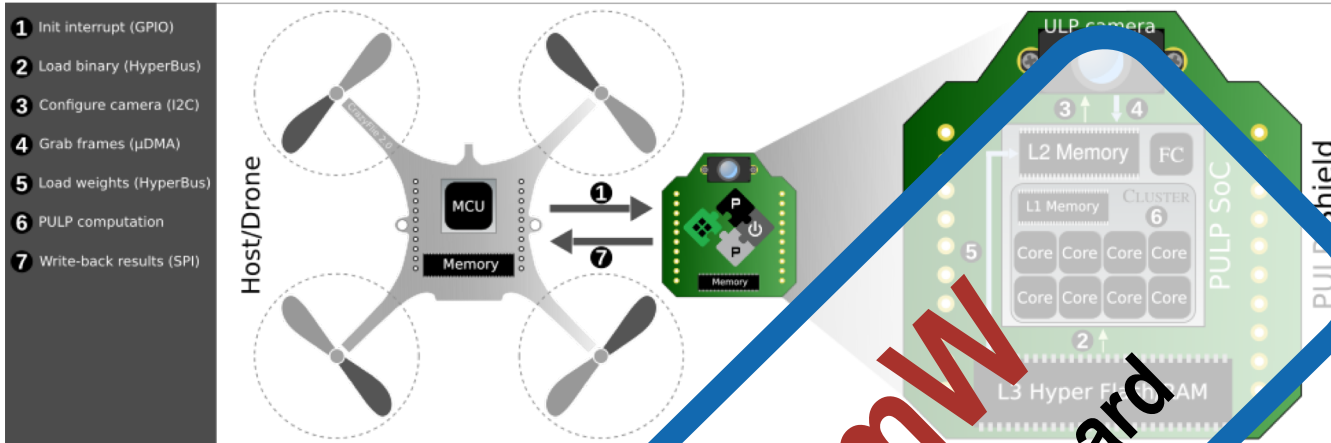


What	Freq MHz	Exec Time ms	Cycles	Power mW
40nm Dual Issue MCU	216	99.1	21 400 000	60
GAP8 @1.0V	15.4	99.1	1 500 000	3.7
GAP8 @1.2V	175	8.7	1 500 000	70
GAP8 @1.0V w HWCE	4.7	99.1	460 000	0.8

4x More efficiency at less than 10% area cost

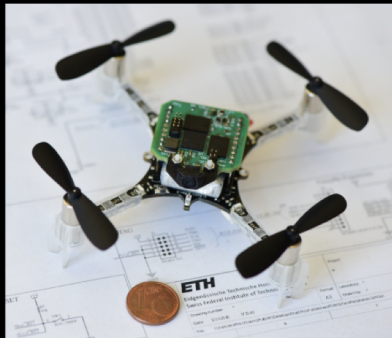


# Complete Application: DroNET on NanoDrone



Pluggable PCB:  
PULP-Shield

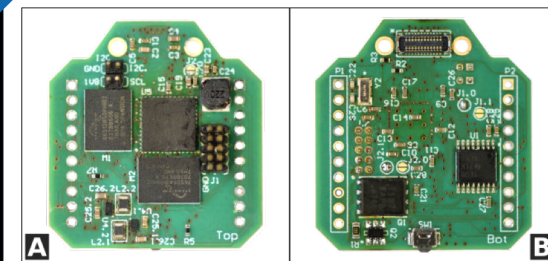
- ~5g, 30×28mm
- GAP8 SoC
- 8 MB HDRAM
- 16 MB HFlash
- QVGA ULP
- HiMax camera
- Crazyflie 2.0 nano-drone (27g)



Copyright 2019 © ETH zürich



Credit: F. K. Gürkaynak & Daniele Palossi



Only onboard computation for autonomous flight + obstacle avoidance  
no human operator, no ad-hoc external signals, and no remote base-station!



# PULP has released a large number of IPs

## RISC-V Cores

<b>RI5CY</b>	<b>Ibex</b>	<b>Snitch</b>	<b>Ariane + Ara</b>
32b	32b	32b	64b

## Peripherals

<b>JTAG</b>	<b>SPI</b>
<b>UART</b>	<b>I2S</b>
<b>DMA</b>	<b>GPIO</b>

## Interconnect

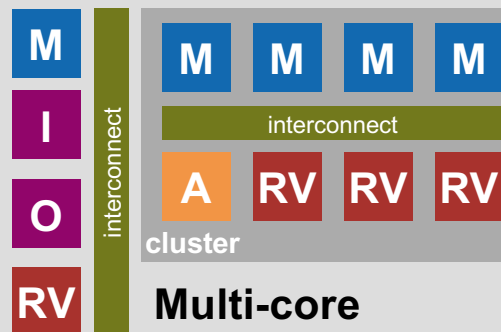
<b>Logarithmic interconnect</b>
<b>APB – Peripheral Bus</b>
<b>AXI4 – Interconnect</b>

## Platforms



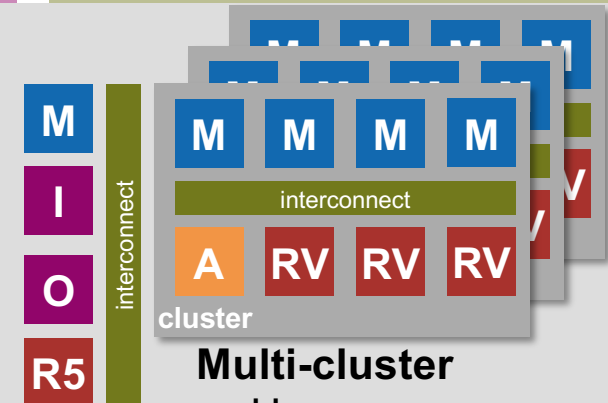
### Single Core

- PULPino
- PULPissimo



### Multi-core

- Fulmine
- Mr. Wolf



### Multi-cluster

- Hero
- Open Piton

## IOT

### Accelerators

**HWCE**  
(convolution)

**Neurostream**  
(ML)

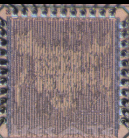
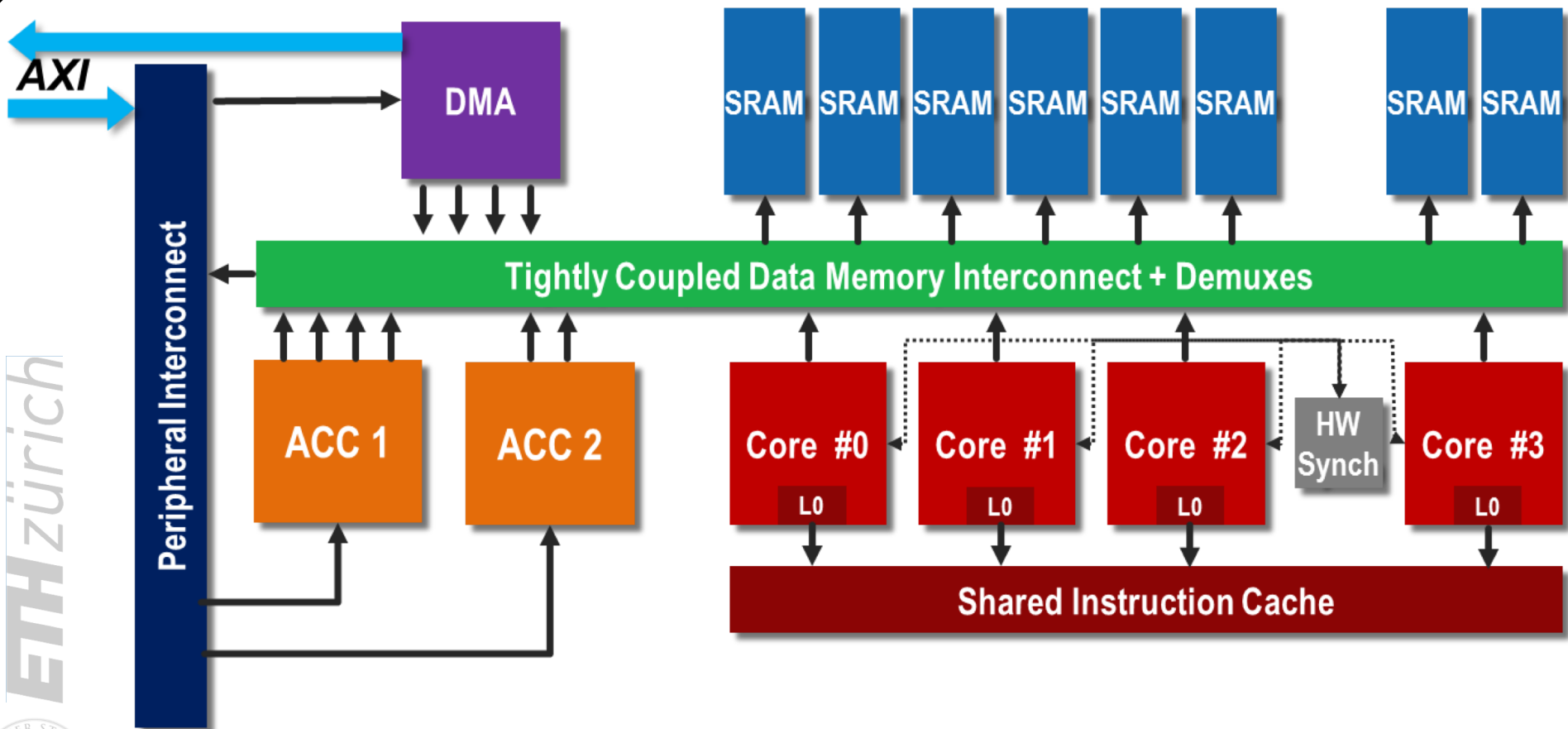
**HWCrypt**  
(crypto)

**PULPO**  
(1<sup>st</sup> ord. opt)

## HPC

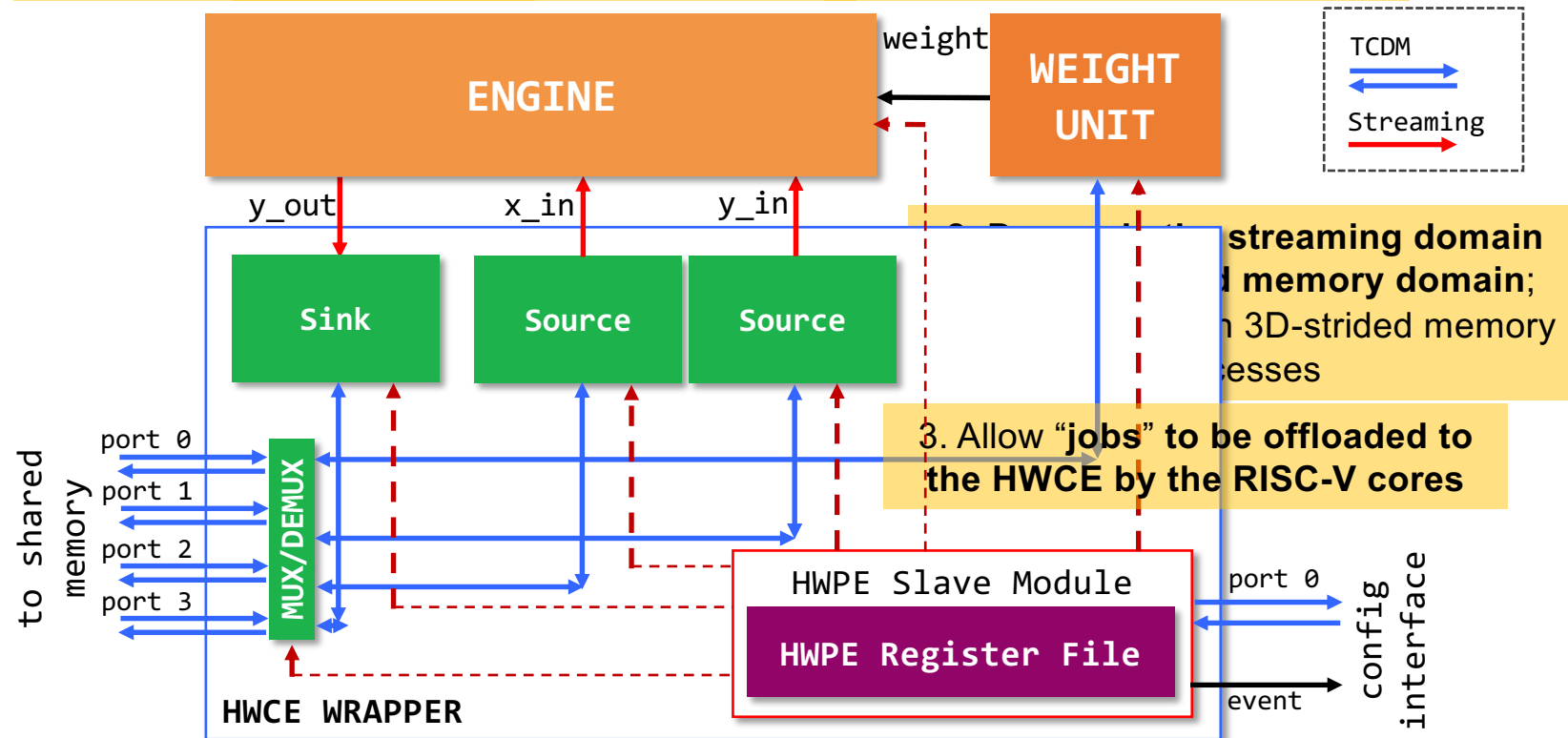


# Shared Memory Accelerators

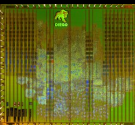


# Accelerator: HW Convolution Engine

1. **Block gating** involve-accumulate to minimize dynamic power in streaming fashion
2. **Weights** for each convolution filter are **stored privately**
3. Allow "jobs" to be offloaded to the HWCE by the RISC-V cores
4. **Weights** for each convolution filter are **stored privately**
5. **Fine-grain clock gating** to minimize dynamic power in streaming fashion



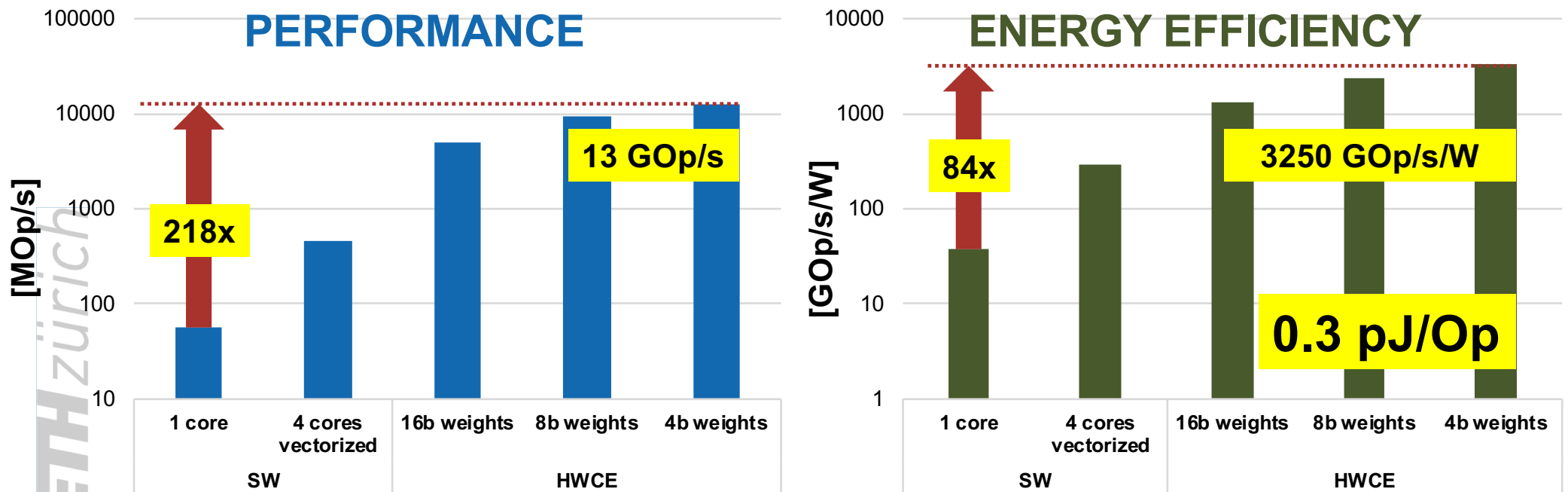
F. Conti and L. Benini, "A ultra-low-energy convolution engine for fast brain-inspired vision in multicore clusters," *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2015, pp. 683-688.



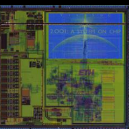
# HW Convolution Engine Performance

Cluster performance and energy efficiency on a 64x64 CNN layer (5x5 conv)

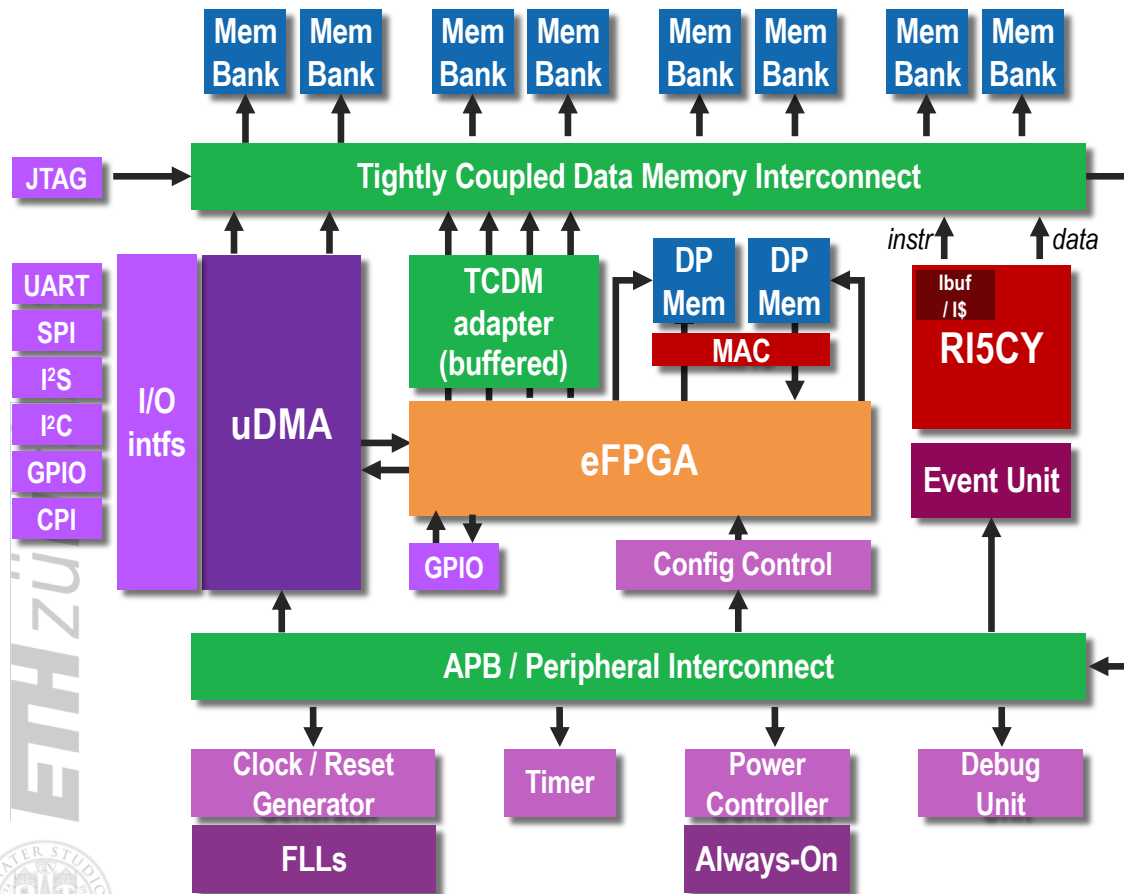
Scaled to ST FD-SOI 28nm @  $V_{dd}=0.6V$ ,  $f=115MHz$



Now coming: HWCE 2.0 – improves scalability & flexibility @ 3TOPS/W



# Arnold: a Collaboration with Quicklogic

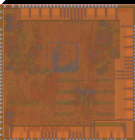
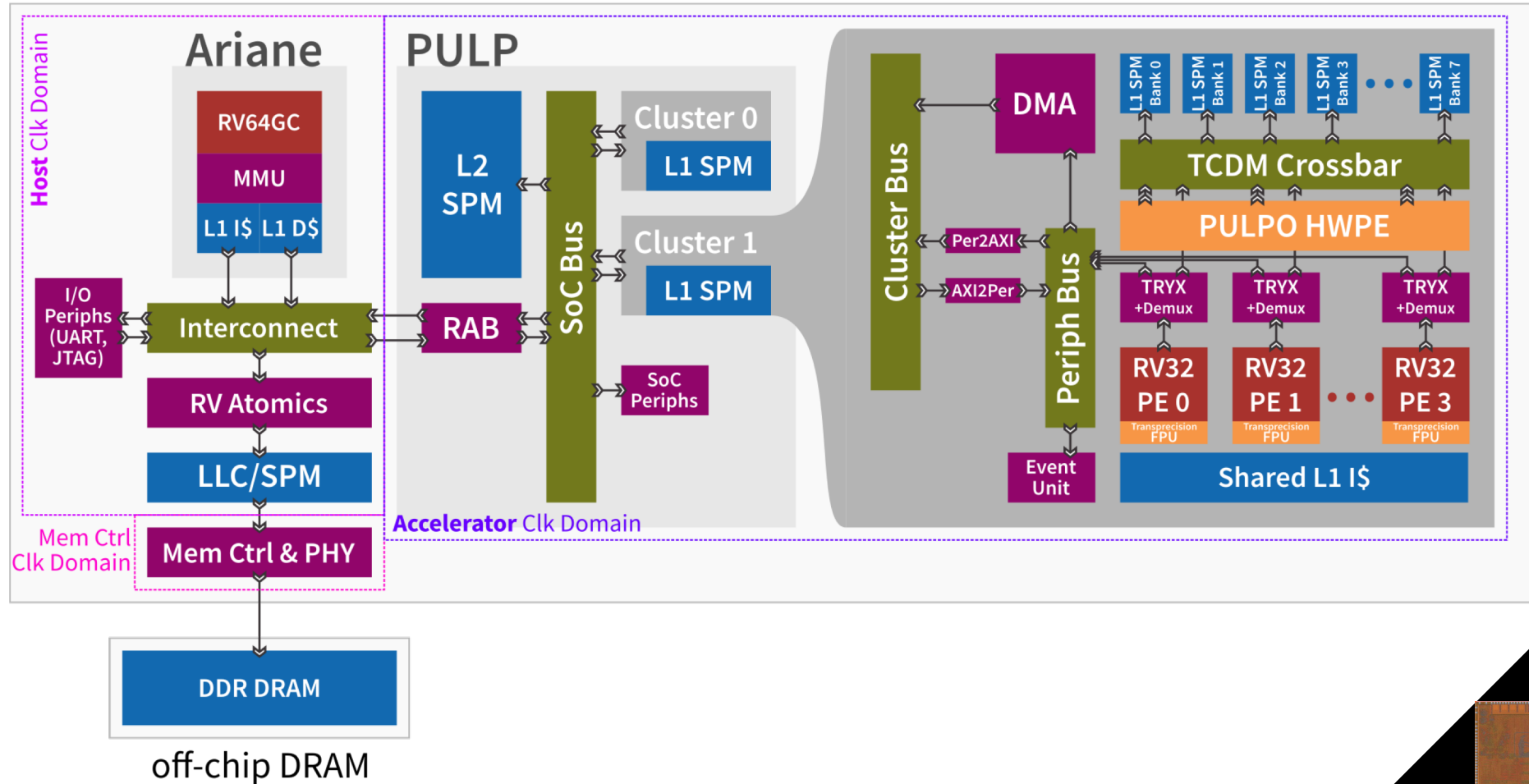


- **Chip in 22nm FDX**
  - Combines e-FPGA (Quicklogic)
  - with PULPissimo (single core uC)
- **Multiple operation modes**
  - Configurable peripheral
  - Accelerator for core
  - Accelerator for independent I/O



# Cluster as heterogenous accelerator (HERO)

HEROv3



# FPGA Prototyping Platforms

## Available:

### ■ Digilent Genesys2

- \$999 (\$600 academic)
- 1-2 cores at 66MHz



### ■ Xilinx VC707

- \$3500
- 1-4 cores at 60MHz

### ■ Digilent Nexys Video

- \$500 (\$250 academic)
- 1 core at 30MHz



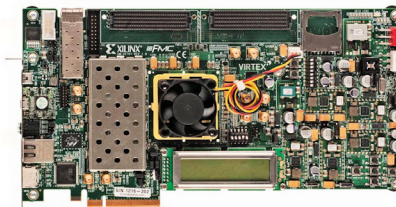
## In progress:

### ■ Xilinx VCU118, BittWare XUPP3R

- \$7000-8000
- >100MHz

### ■ Amazon AWS F1

- Rent by the hour

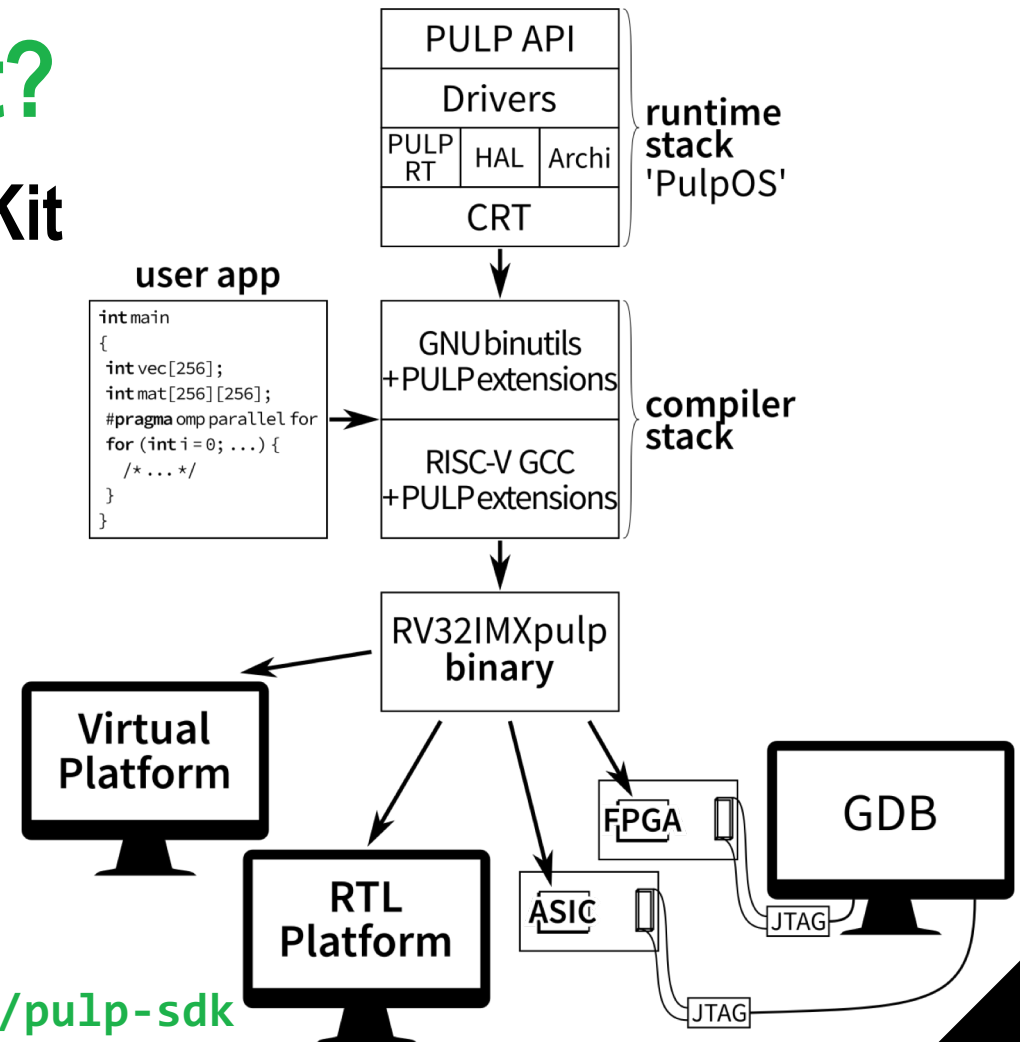




# How about SW support?

## PULP Software Development Kit

- Package for compiling, running, debugging and profiling applications on PULP platforms
- Supports all recent and upcoming PULP chips: Mr.Wolf, GAP, Vega, ...
- Supports all targets: Virtual Platform, RTL platform, FPGA
- RISC-V GCC with support for PULP extensions
- Basic OpenMP support
- <https://github.com/pulp-platform/pulp-sdk>



# What is PULP doing for maintain our cores?

- **We (ETHZ and University of Bologna) are research groups**
  - Motivated to develop new architectures and systems
  - We needed efficient RISC-V cores (and peripherals) for our work
  - Not so good (or interested) in providing industrial level support for these cores
- **We need help to**

ETH zürich



# Academic open source → Industrial open source



**OPENHW** GROUP  
— PROVEN PROCESSOR IP —

Rick O'Connor (OpenHW CEO, former RISC-V foundation director)

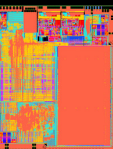
- **OpenHW Group** is a not-for-profit, global organization (EU,NA,Asia) driven by its members and individual contributors where HW and SW designers collaborate in the development of open-source cores, related IP, tools and SW such as the **Core-V** family of cores.
- **OpenHW Group** provides an infrastructure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



RI5CY, ARIANE



**CORE-V**™



# A vertical, application-focused open approach



- OpenTitan is **the first open source** silicon project building a transparent, high-quality reference design for silicon root of trust (RoT) chips.
- Founding Partners

**ETH** zürich

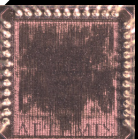
**GD** Giesecke+Devrient  
Creating Confidence

**Google**

**lowRISC**

**nuvoTon**

**Western Digital.**



# Feel the momentum!

Ibex RISC-V core, flash interface, communications ports, cryptography accelerators, and more.

## Vibrant repository

**35+** Contributors  
**1300+** Contributions  
**470** GitHub Issues



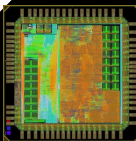
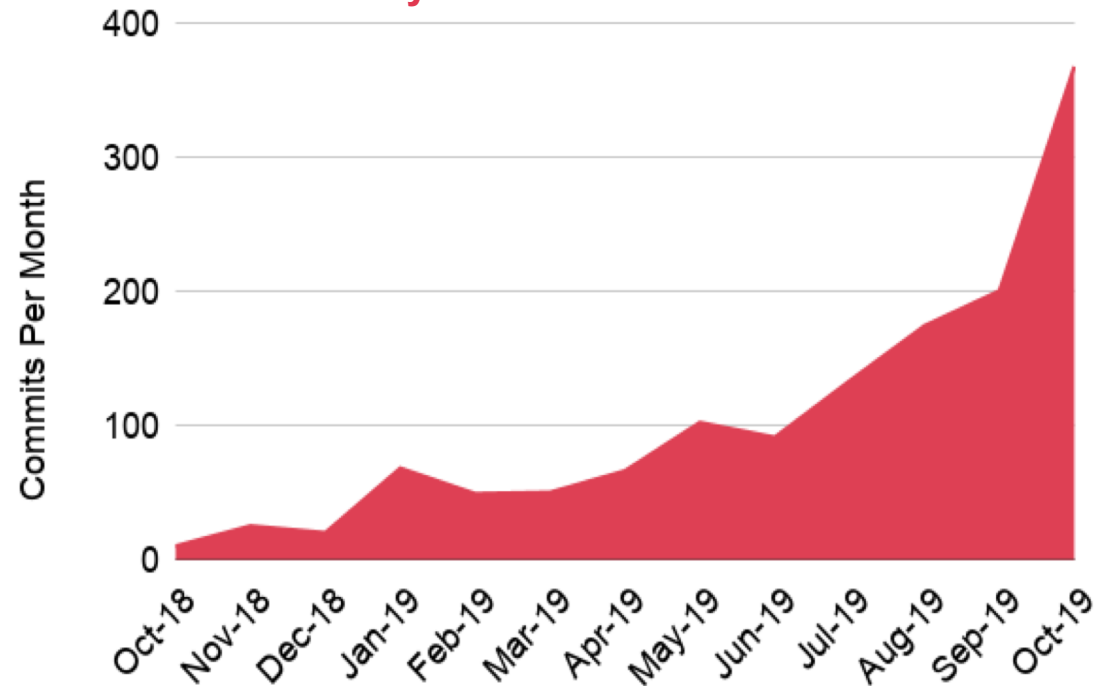
**PULP**  
 Parallel Ultra Low Power



**lowRISC**

**Zero-Riscy**

**Ibex**







# PULP

Parallel Ultra Low Power

Luca Benini, Davide Rossi, Andrea Borghesi, Michele Magno, Simone Benatti, Francesco Conti, Francesco Beneventi, Daniele Palossi, Giuseppe Tagliavini, Antonio Pullini, Germain Haugou, Lukas Cavigelli, Manuele Rusci, Florian Glaser, Renzo Andri, Fabio Montagna, Bjoern Forsberg, Pasquale Davide Schiavone, Alfio Di Mauro, Victor Javier Kartsch Morinigo, Tommaso Polonelli, Fabian Schuiki, Stefan Mach, Andreas Kurth, Florian Zaruba, Manuel Eggimann, Philipp Mayer, Marco Guermandi, Xiaying Wang, Michael Hersche, Robert Balas, Antonio Mastrandrea, Matheus Cavalcante, Angelo Garofalo, Alessio Burrello, Gianna Paulin, Georg Rutishauser, Andrea Cossettini, Luca Bertaccini, Maxim Mattheeuws, Samuel Riedel, Sergei Vostrikov, Vlad Niculescu, Frank K. Gurkaynak, *and many more that we forgot to mention*



<http://pulp-platform.org>



@pulp\_platform



# FOSSistanbul, March 13-15, Istanbul

**F**OSS**i**stanbul will bring together, enthusiasts, members of industry and academia that are working on open source hardware design, in a lively and attractive city.

With keynotes by: Luca Benini, Nele Mentens, Onur Mutlu

Register for **FREE**

<https://fossi-foundation.org/fossistanbul/>

