# Scaling performance In Power-Limited HPC Systems

**Prof. Dr. Luca Benini**

**ERC Multitherman Lab
University of Bologna – Italy**

**D-ITET, Chair of Digital Dircuits
& Systems - Switzerland**

Eidgenössische Technische Hochschule Zürich
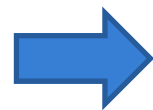Swiss Federal Institute of Technology Zurich

# Outline

- ❏ **Power and Thermal Walls in HPC**

- ❏ **Power and Thermal Management**

- ❏ **Energy-efficient Hardware**

- ❏ **Conclusion**

~170MW

Exascale computing in 2020

Sunway TaihuLight 93 PF, 15.3 MW **6 GF/W**

50 GFLOP/W!!

**We need almost 10x more energy efficiency**

30% energy budget of today's nuclear reactor

Feasible Exascale power budget ≤ 20MWatts

o500 ranks the new percomputers by FLOPS Linpack Benchmark

**The second, Tianhe-2 (ex 1st) consumes 17.8 MW for "only" 33.2 PetaFLOPs, but…**

24 MW

⚠ Cooling system matters!!!

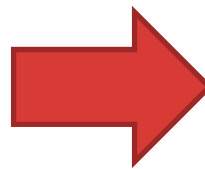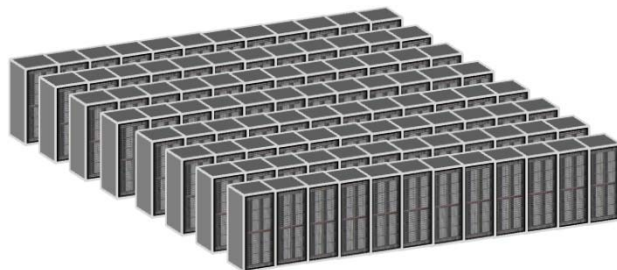**Dynamic Power management (DPM)**
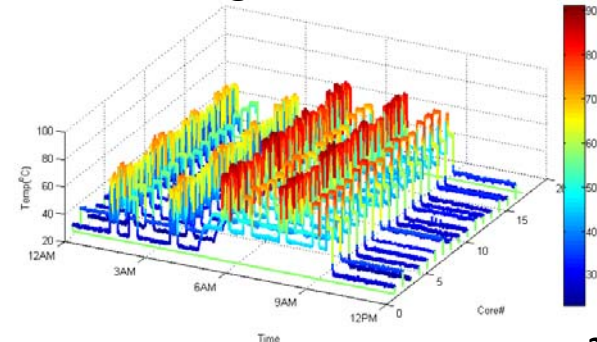
# Thermal Wall → Max+

**Intel Haswell – E5-2699 v3 (18 core)**

Up to **24°C** Temperature difference on DIE
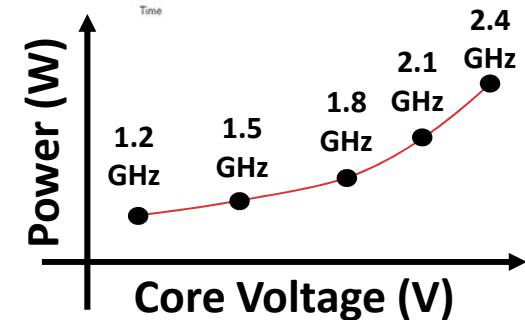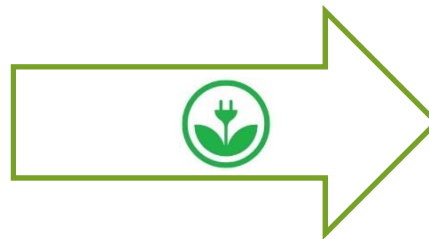More than **7°C** thermal heterogeneity under same workload

## HPC System

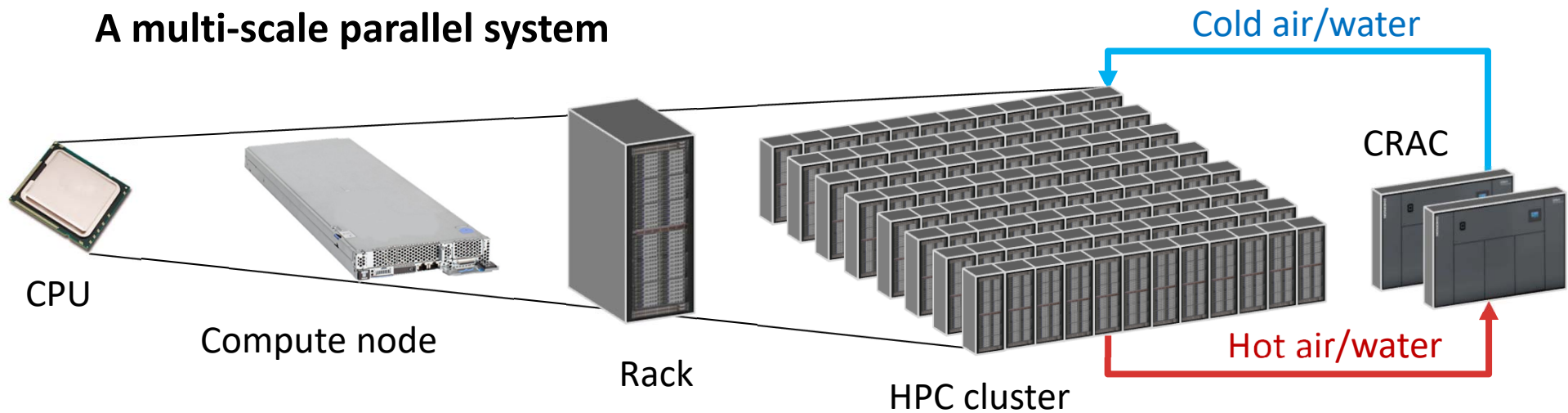**Thermal range: 69 C° – 101 C°**

**Per-core DVFS approach**

**Power consumption: 40% - 66%**

Power (W)

1.2 GHz
1.5 GHz
1.8 GHz
2.1 GHz
2.4 GHz

Core Voltage (V)

## Dynamic thermal management (DTM)

# HPC Architecture - Hardware



**A multi-scale parallel system**



Cold air/water

CPU

Compute node

Rack

HPC cluster

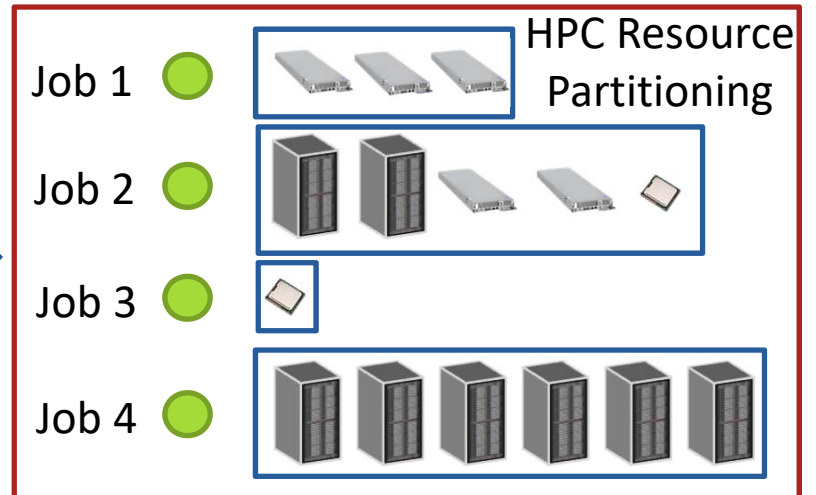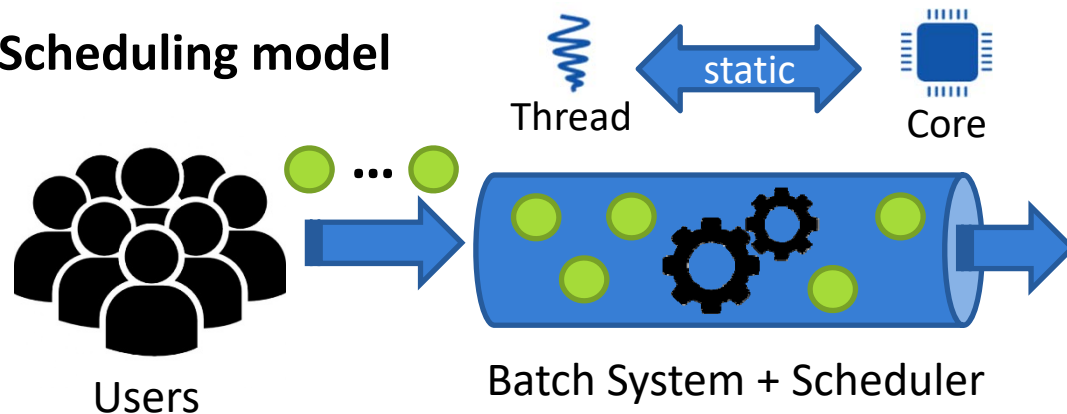CRAC

Hot air/water

**DPM, DTM are Multi-scale Problems!** ➡️ Multitherman erc

European Research Council
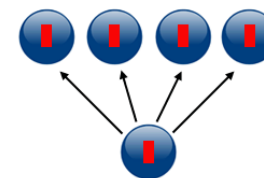
# HPC Architecture - Software

**Scheduling model**



Users

static

Thread ← static → Core

Batch System + Scheduler

HPC Resource Partitioning

Job 1
Job 2
Job 3
Job 4

**Programming Model**

MPI_COMM_WORLD



group1

group2

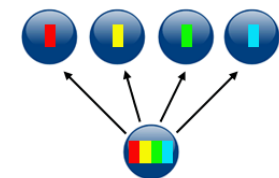comm1

comm2
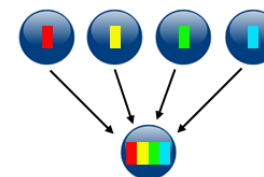
communications

**COMMUNICATIONS**

**one-to-one, one-to-many, many-to-one and many-to-many**
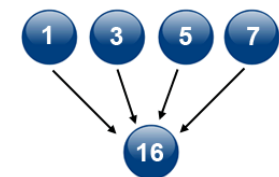


broadcast

scatter

gather

reduction

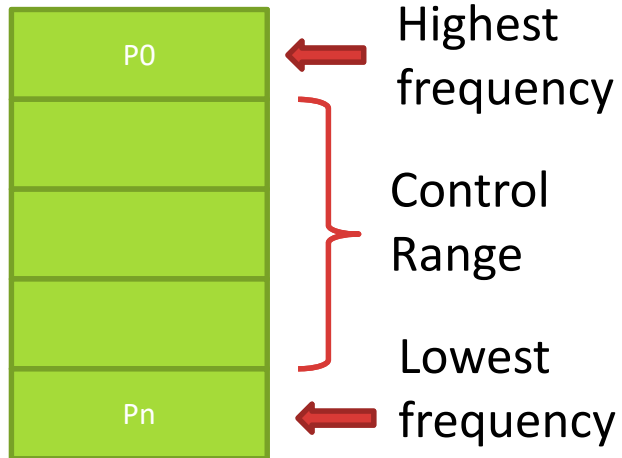**Programming & Scheduling model is essential!**

# Outline

❑ **Power and Thermal Walls in HPC**

❑ **Power and Thermal Management**

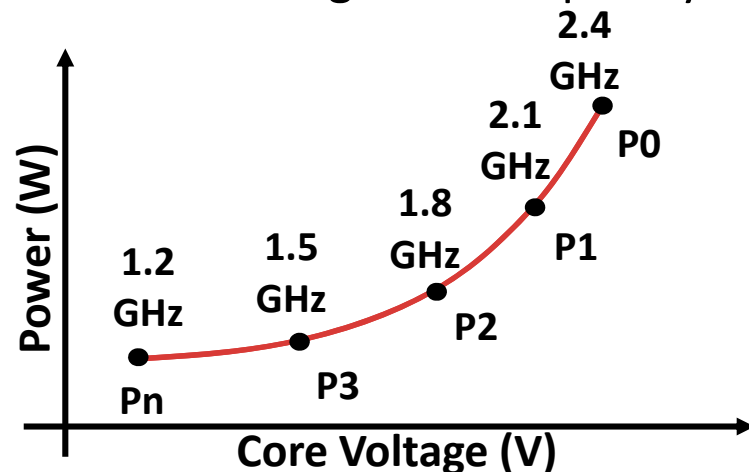❑ **Energy-efficient Hardware**

❑ **Conclusion**

# HW Support for DPM, DTM
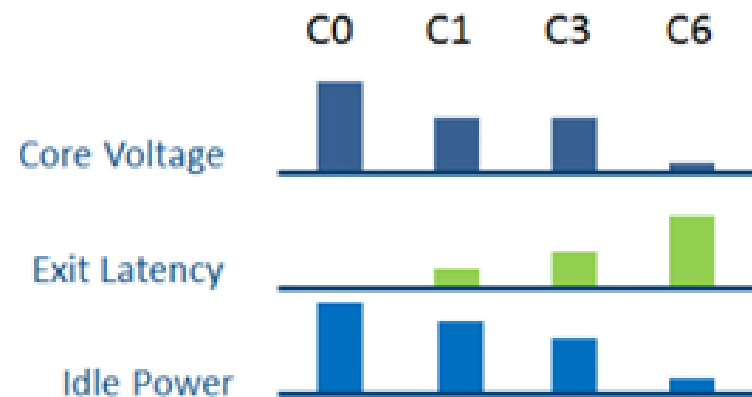
**ACTIVE STATES**
**DVFS (P-State)**



P0 → Highest frequency

} Control Range

Pn → Lowest frequency

P-State: both a voltage and frequency level



Power (W) vs Core Voltage (V):
- Pn — 1.2 GHz
- P3 — 1.5 GHz
- P2 — 1.8 GHz
- P1 — 2.1 GHz
- P0 — 2.4 GHz

**Intel provides a HW power controller called Running Average Power Limit (RAPL).**

**IDLE STATES**
**low power (C-State)**



| | | Core C-States | | | |
|---|---|---|---|---|---|
| | | C0 | C1 | C3 | C6 |
| Package C-States | C0 | | | | | Active State |
| | C1E | | | | | Lower P-State |
| | C2 | | | | | Only L3 Snoop |
| | C3 | | | | | Flush L3 - Off |
| | C6 | | | | | Low Voltage |
| | | Active State | Clock Gated | Flush L1,L2 Off | Power Gated | |

■ Possible combination of core/package states
■ Impossible combination of core/package states


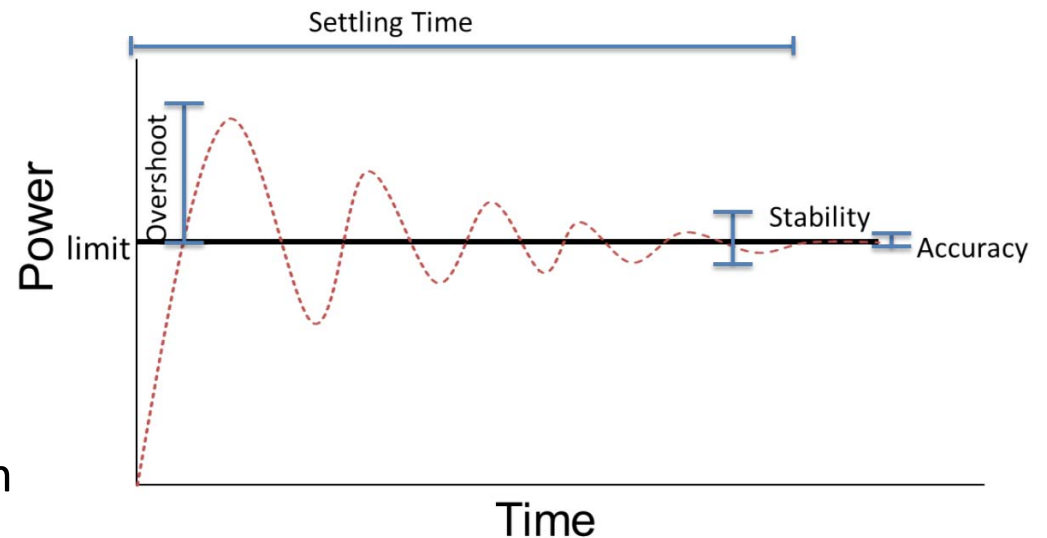
C0   C1   C3   C6

Core Voltage
Exit Latency
Idle Power

# Power Management → Reactive

A significant exploration work on RAPL control:

❖ *Zhang, H., & Hoffman, H. (2015). "A Quantitative Evaluation of the RAPL Power Control System". Feedback Computing.*

Quantify the behavior the control system in term of:
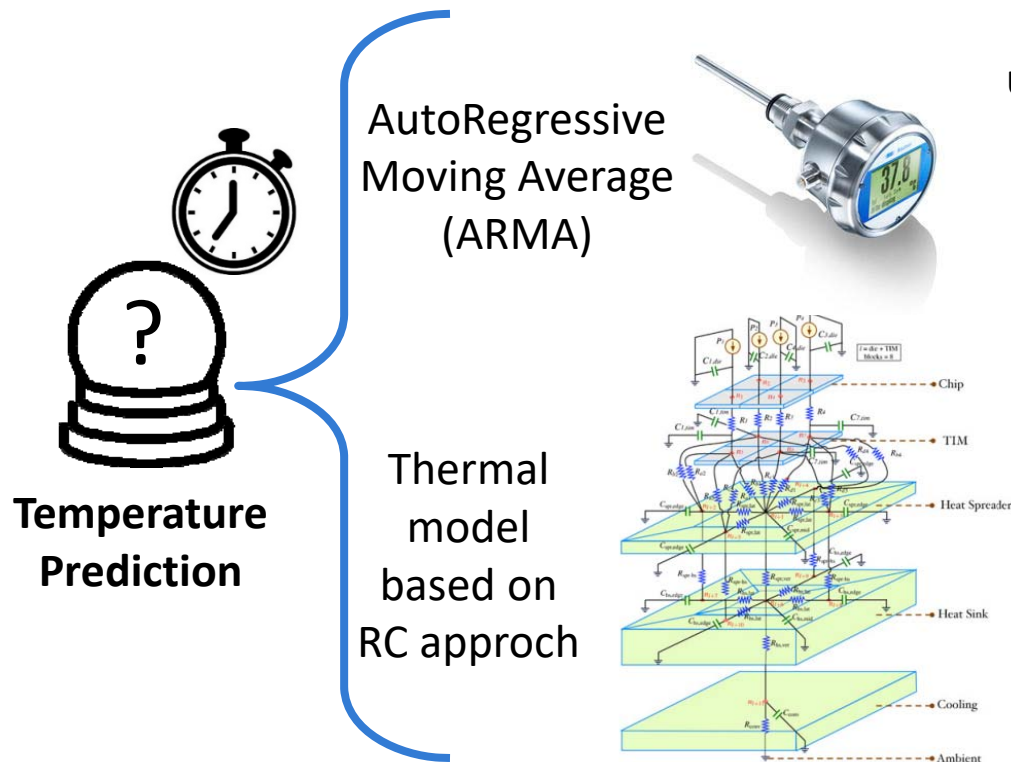
➢ **Stability**: freedom from oscillation

➢ **Accuracy**: convergence to the limit

➢ **Settling time**: duration until limit is reached

➢ **Maximum Overshoot**: the maximum difference between the power limit and the measured power
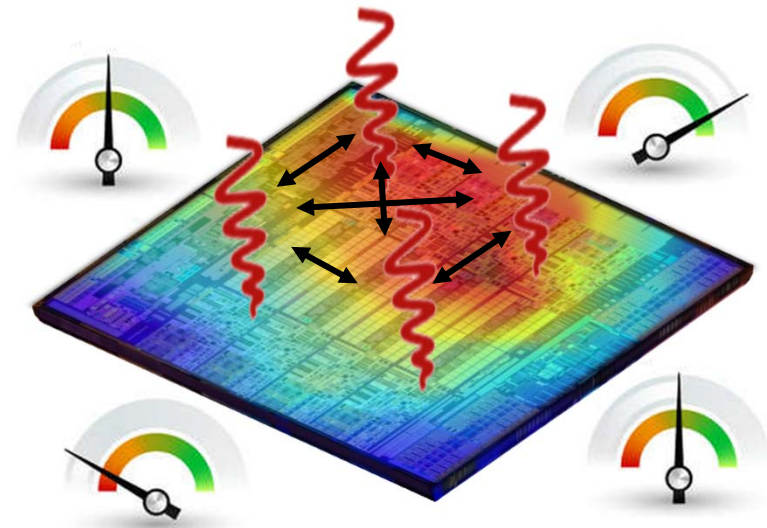
# Power Management → HW Predictive

❖ **on-line optimization policies**
  - *A. Bartolini et al. "Thermal and Energy Managementof High-Performance Multicores: Distributed and Self-Calibrating Model-Predictive Controller." TPDS'13*



AutoRegressive Moving Average (ARMA)

Implement proactive and reactive policies using DVFS selections and thread migrations

**Temperature Prediction**
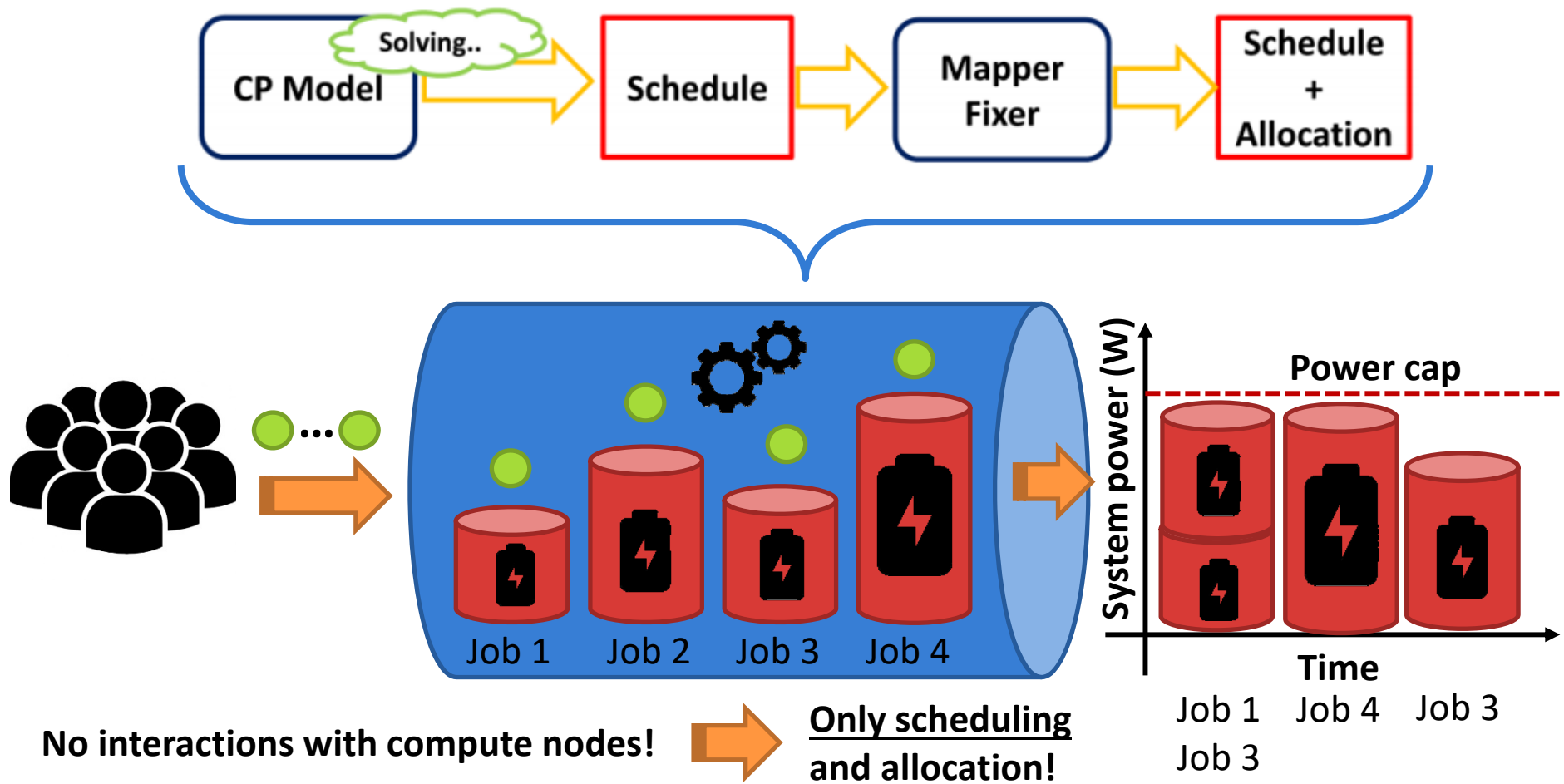
Thermal model based on RC approch

Scheduler based on convex optimization for DVFS selections and thread migrations

Online techniques are capable of sensing changes in the workload distribution and setting the processor controls accordingly.

# Power Management → SW predictive

❖ **Predictive models to estimate the power consumption**

- *Borghesi, A., Conficoni, C., Lombardi, M., & Bartolini, A. "MS3: a Mediterranean-Stile Job Scheduler for Supercomputers-do less when it's too hot!". HPCS 2015*
- *Sîrbu, A., & Babaoglu, O. "Predicting system-level power for a hybrid supercomputer". HPCS 2016*



No interactions with compute nodes! ➡ **Only scheduling and allocation!**

# Challenges

SW policies

HW mechanisms

Application aware

High overhead
Coarse granularity
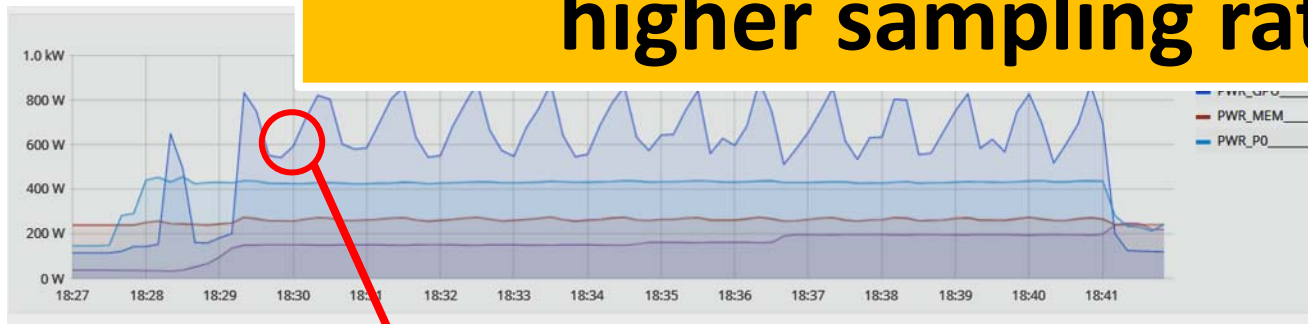(seconds)

Low overhead
Fine granularity
(milliseconds)

No application awareness

1) **Low-Overhead, accurate monitoring**
2) **Scalable data collections, analytics, decisions**
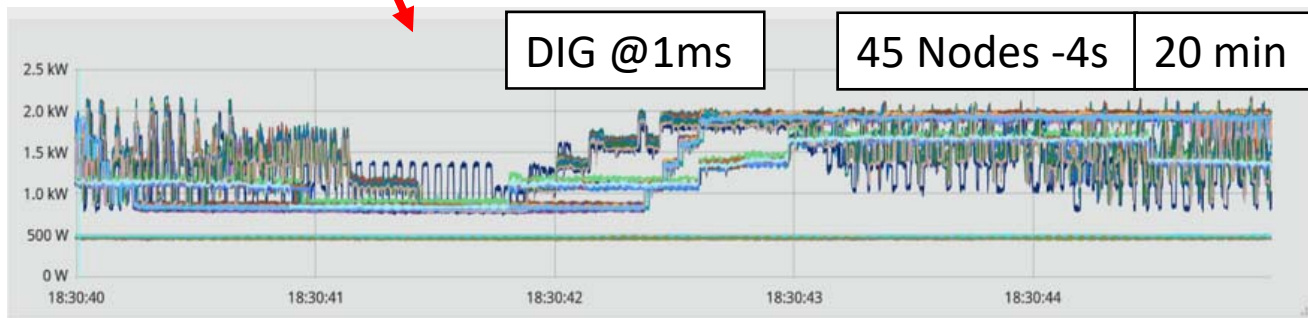3) **Application awareness**

## Low Overhead, accurate Monitoring

High-resolution monitoring → more information available

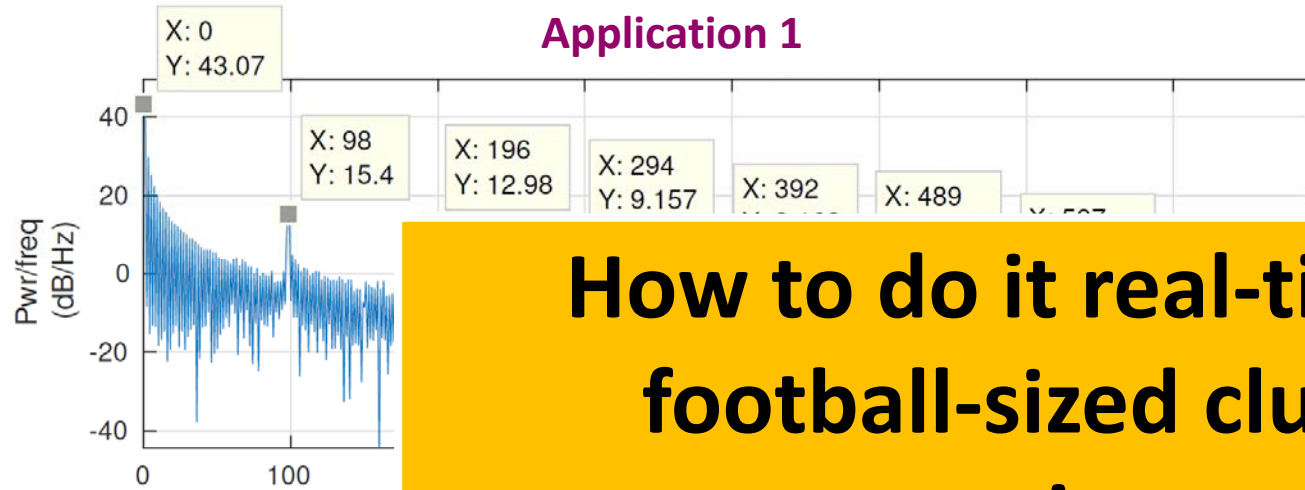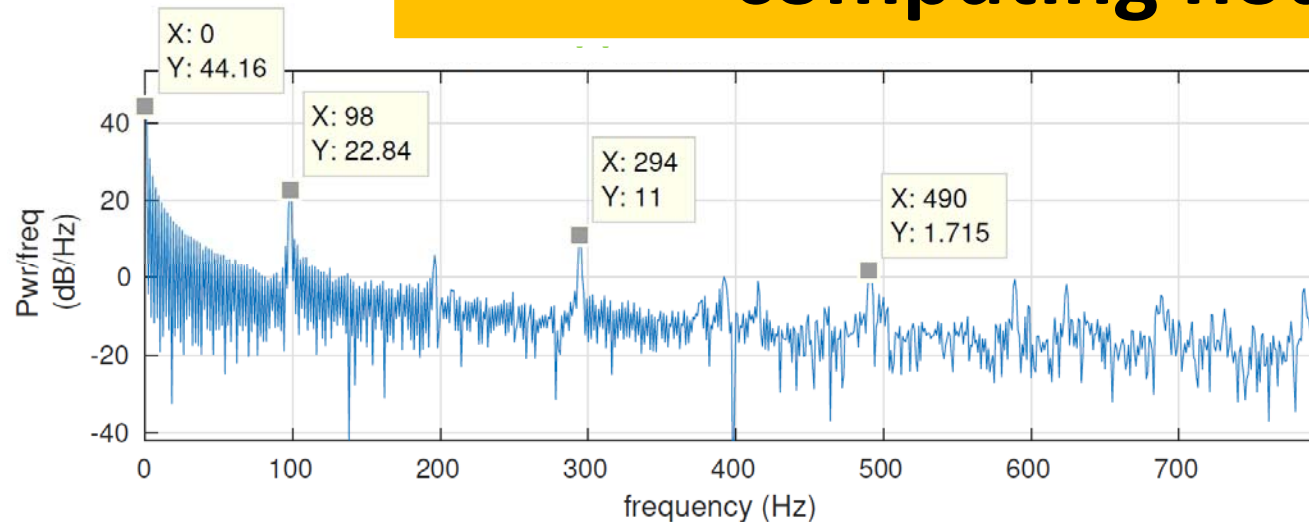**How to analyze real-time with higher sampling rates?**

Max. Ts = 1s

DIG @1ms     45 Nodes -4s     20 min

# Low Overhead, accurate Monitoring
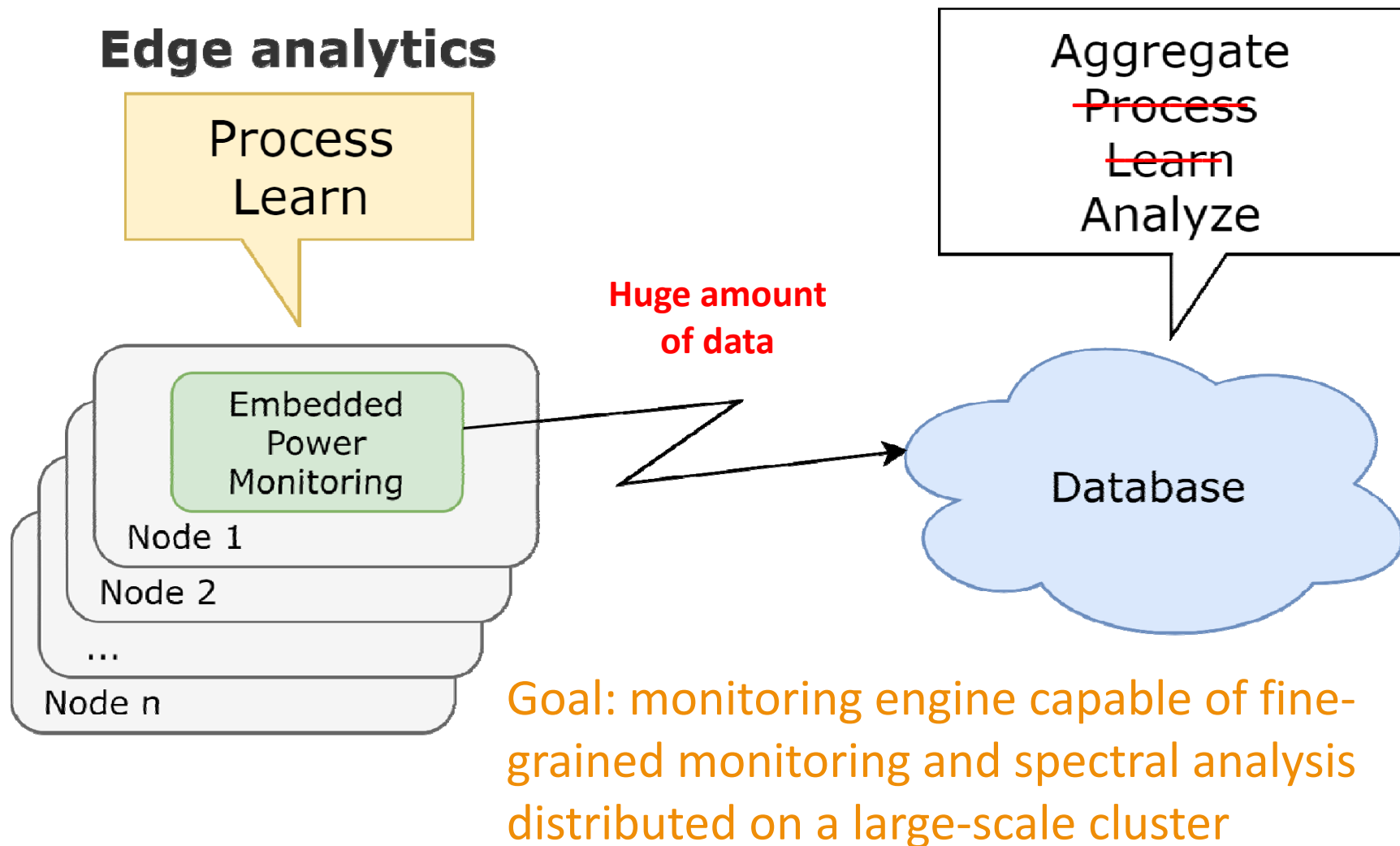
Real-time Frequency analysis on power supply and more…

**Application 1**

How to do it real-time for a football-sized cluster of computing nodes?

# Solution – Dwarf In a Giant (DIG)



**Edge analytics**

Process
Learn

Embedded
Power
Monitoring

Node 1
Node 2
...
Node n

**Huge amount
of data**

Aggregate
~~Process~~
~~Learn~~
Analyze

Database

Goal: monitoring engine capable of fine-grained monitoring and spectral analysis distributed on a large-scale cluster
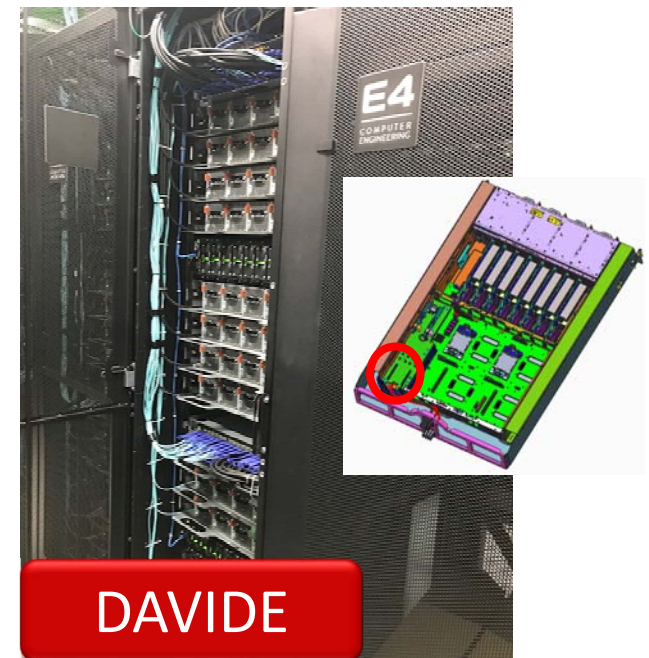
# DIG in Real Life

Developing hardware extensions for fine-grained power monitoring: DIG deployed in production machines



**"Galileo"**

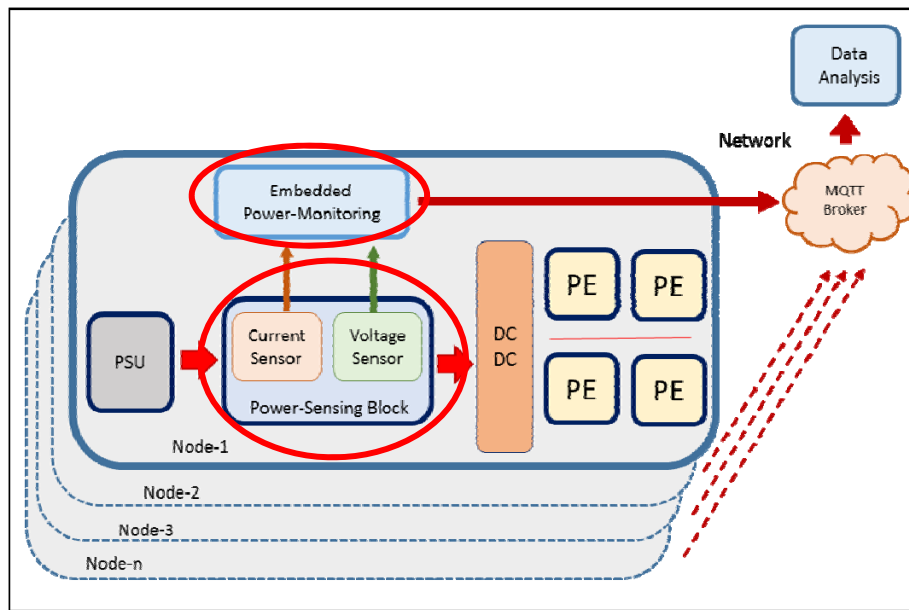- Intel Xeon E5 based
- Used for prototyping

**ARM64**

- ARM64 Cavium based
- Commercial system
- with E4 - PCP II

**DAVIDE**

- IBM Power8 based
- Commercial system
- with E4 - PCP III
- 18th in Green500

## High Resolution Out-of-band Power Monitoring



- Overall node power consumption
- Can support edge computing/learning
- Platform independent (Intel, IBM, ARM)
- Sub-Watt precision
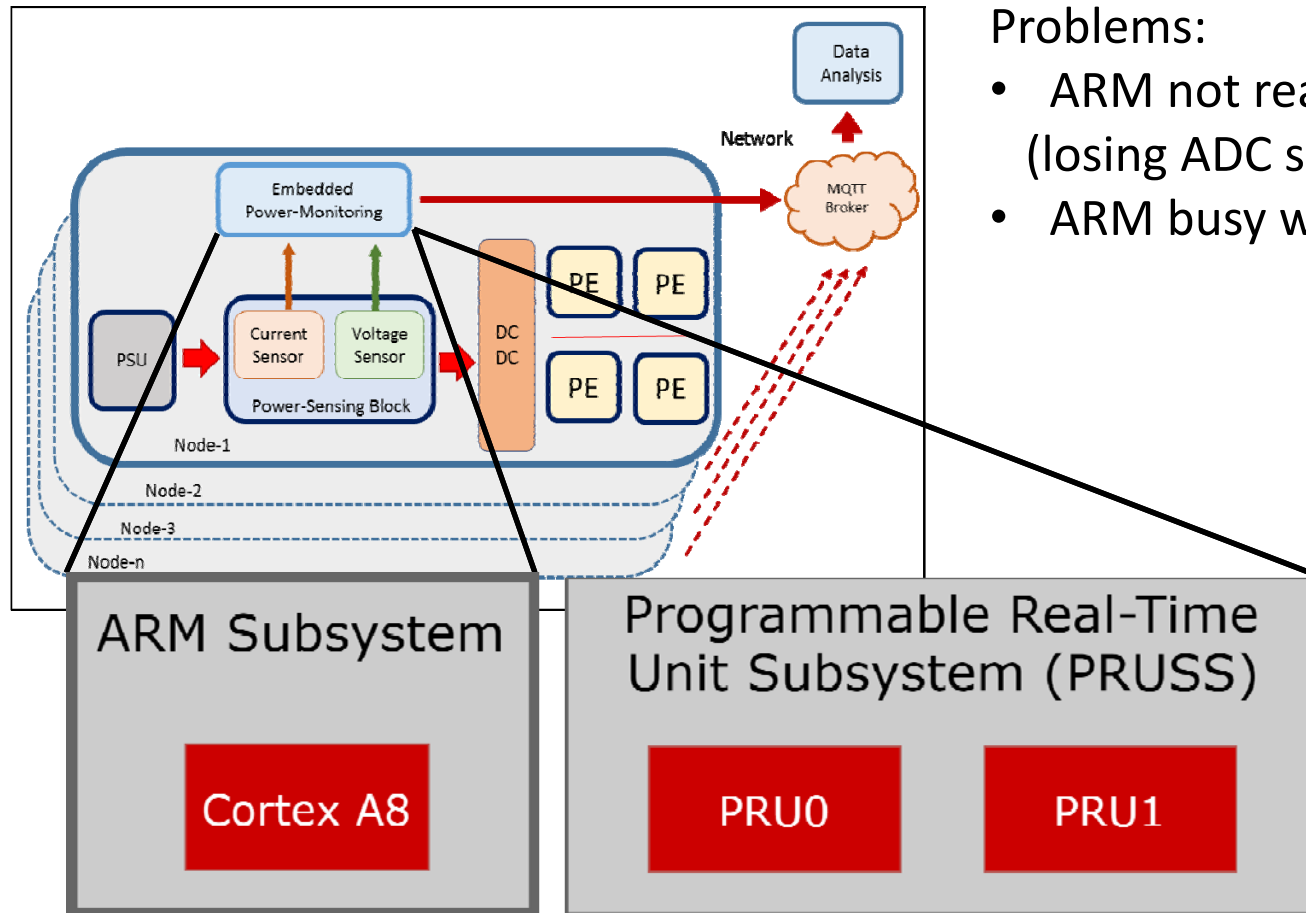- Sampling rate @50kS/s (T=20us)

State-of-the art systems (Bull-HDEEM and PowerInsight)
- Max. 1 ms sampling period
- Use data only offline

Hackenberg et al. "HDEEM: high definition energy efficiency monitoring"
Laros et al. "Powerinsight-a commodity power measurement capability."

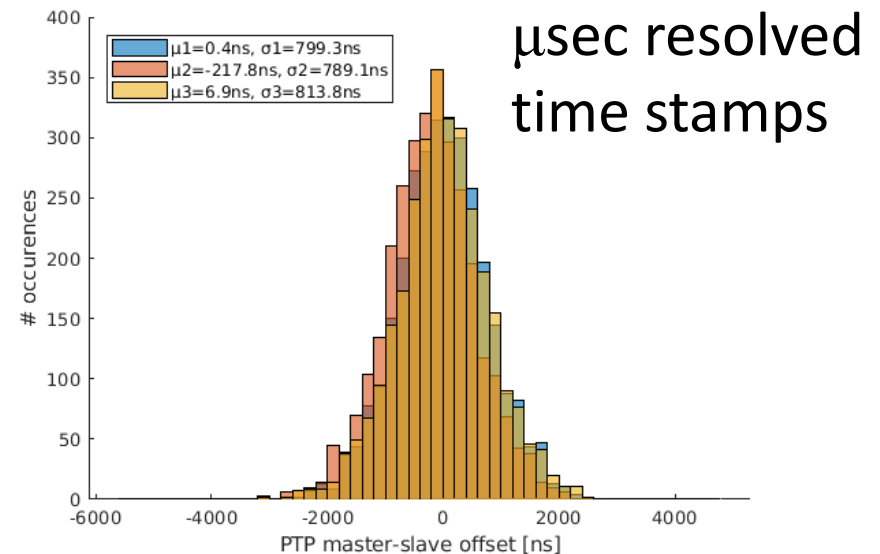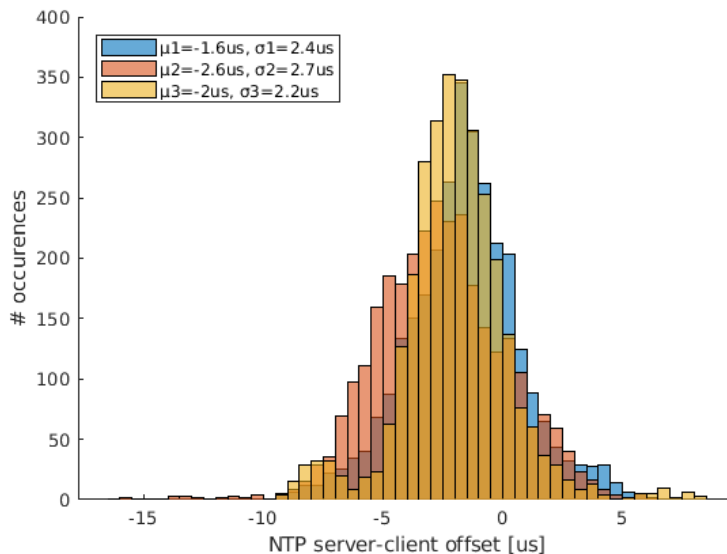# Real-time Capabilities



Problems:
- ARM not real-time (losing ADC samples )
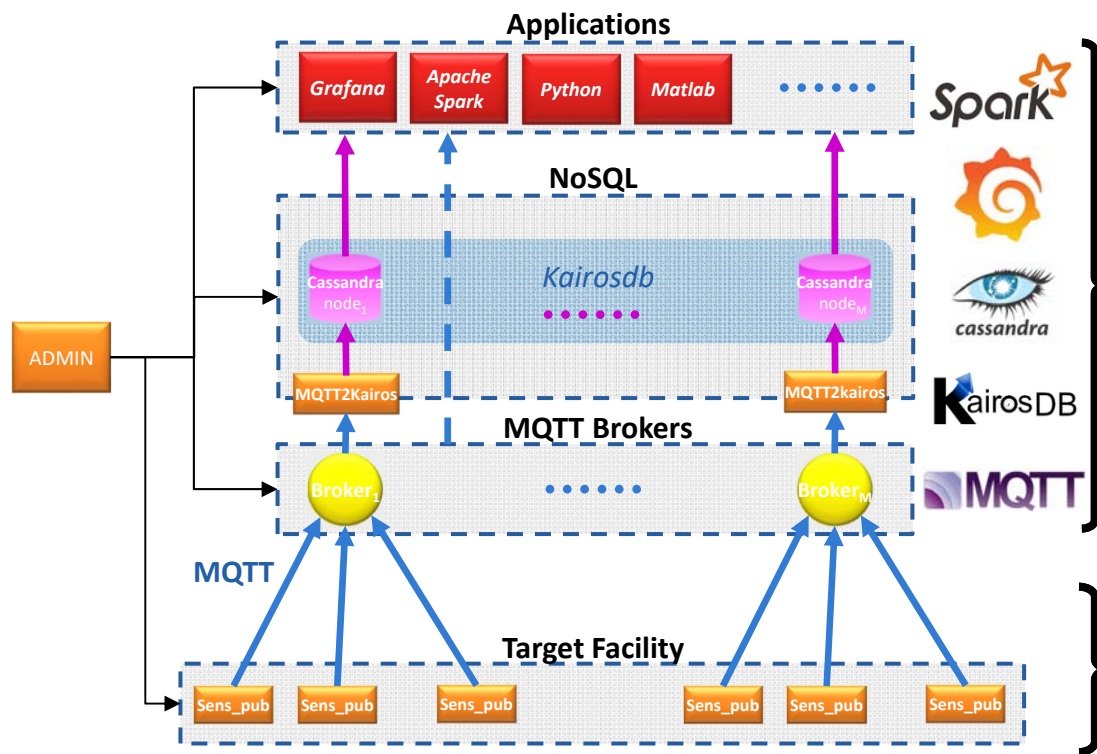- ARM busy with flushing ADC

Goal:
Offload the processing to the PRUSS

# DIG in production: E4's D.A.V.I.D.E.

Possible tasks of the PRUs: Averaging @ 1ms, 1s → offline Computing, FFT → edge analysis

| Framework | $Fs_{max}$ [kHz] | CPU Overhead |
|---|---|---|
| DIG | 50 | ~40% |
| DIG+PRU, edge analysis | 400 | <5% |
| DIG+PRU, offline | 800 | <5% |
| Bull-HDEEM | 1 | ? |
| PowerInsight | 1 | ? |

μsec resolved time stamps

# Scalable Data Collection, Analytics

# MQTT to NoSQL Storage: MQTT2Kairosdb

● = {Value;Timestamp}

**MQTT**

MQTT Broker

facility/sensors/# → MQTT2Kairosdb

facility/sensors/A

facility/sensors/B

facility/sensors/C
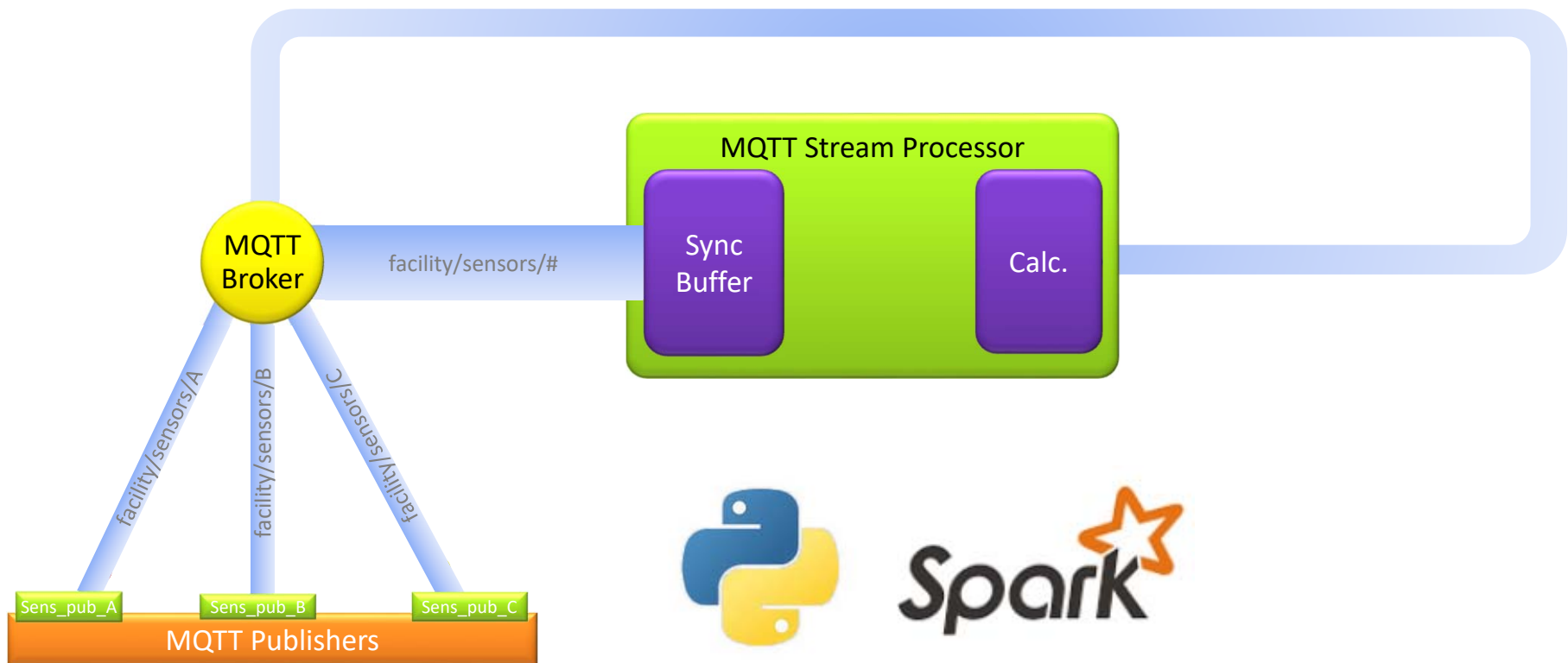
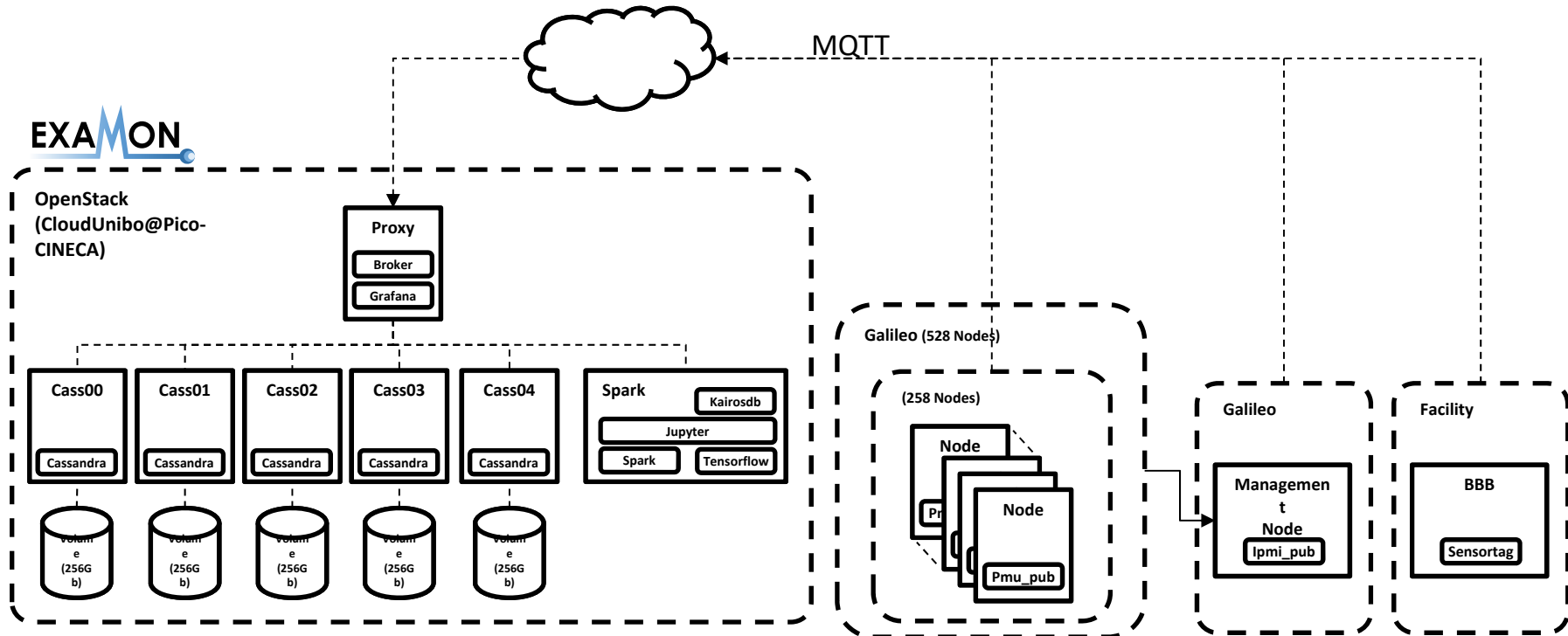Sens_pub_A  Sens_pub_B  Sens_pub_C

MQTT Publishers

# Examon Analytics: Batch & Streaming

# Streaming Analytics: virtual sensors!

# Examon in production: CINECA's GALILEO



| | | | |
|---|---|---|---|
| Data Ingestion Rate | ~67K | Metrics/s | |
| DB Bandwidth | ~98 | Mbit/s | |
| DB Size | ~1000 | GB/week | |
| DB Write Latency | 20 | us | |
| DB Read Latency | 4800 | us | |

Tier1 system 0.5-1TB every week
Tier0 *estimated* 10TB per 3.5 Days

➡ **Stream analytics & distributed processing are a necessity**

# Application Aware En2Sol Minimization

## HARDWARE

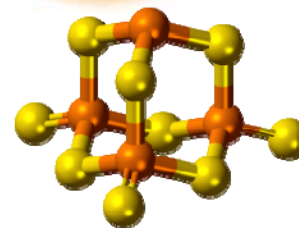**Galileo: Tier-1 HPC system based on an IBM NeXtScale cluster**

- **Cluster**: 516 nodes (14 rack)
- **Node**: Dual socket Intel Haswell E5-2630 v3 CPUs with 8 cores at 2.4 GHz (85W TDP), DDR3 RAM 128 GB
- **Power consumption**: 360 KW
- **OS**: SMP CentOS Linux version 7.0
- **Top500**: Ranked at 281th



**Compute node**

## SOFTWARE

Quantum ESPRESSO is an integrated suite of HPC codes for electronic-structure calculations and materials modelling at the nanoscale.
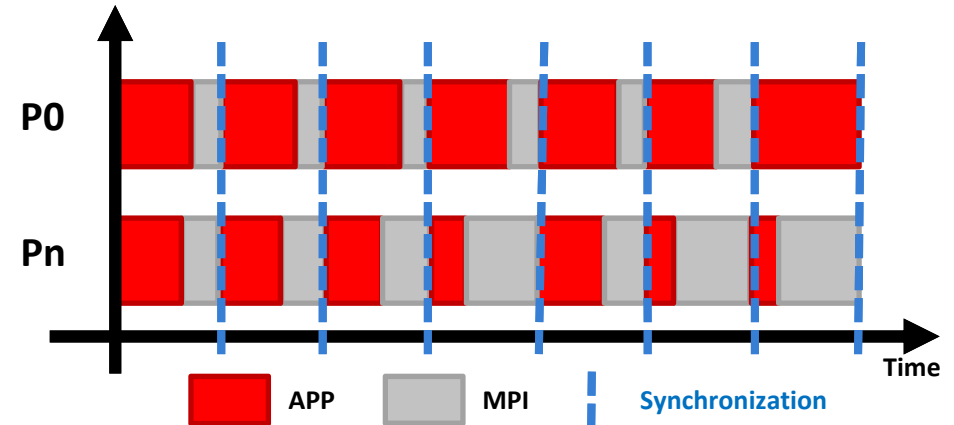


**Car-Parrinello**

**Kernels**

**MPI**

# PMPI

**MPI profiling interface** Augment each standard MPI function with profiling collection functionality



P0

Pn

Time

| APP | MPI | Synchronization |

```c
Include <mpi.h>                                    hello.c

main()
{
    int world_size, world_rank;
    char message[] = "Hello world to everyone from MPI root!"

    // Initialize the MPI environment
    MPI_Init(NULL, NULL);

    // Get the number of processes
    MPI_Comm_size(MPI_COMM_WORLD, &world_size);

    // Get the rank of the process
    MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);

    // Send a broadcast message from root MPI to everyone
    MPI_Bcast(message, strlen(message), MPI_CHAR, 0, MPI_COMM_WORLD);

    // Finalize the MPI environment
    MPI_Finalize();
}
```

```c
Include <mpi.h>                              pmpi_wrapper.c

int MPI_Bcast(void *buffer, int count, MPI_Datatype datatype,
                    int root, MPI_Comm comm)
{
    /* prologue profiling code */
    start_time = get_time();

    int err = PMPI_Bcast(buffer, count, datatype, root, comm);

    /* epilogue profiling code */
    end_time = get_time();
    int duration = end_time - start_time;
    printf("MPI_Bcast duration: %d sec\n", duration);

    return err;
}
```
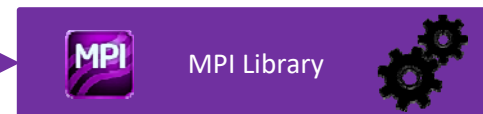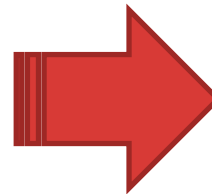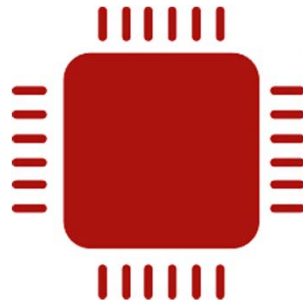
MPI Library

# PMPI Runtime

**Our PMPI implementation has the following features:**

- **Number MPI calls**: 50 MPI functions wrapped (all the QE's MPI calls)
- **Timing**: record TSC for timing (time clock accuracy)
- **Network data**: record all data sent and received from the process
- **Fixed perf counters**: monitor 3 fixed performance counters using low overhead *rdpmc()* instruction
    - Fixed 1: Number of instructions retired
    - Fixed 2: Clock at the nominal frequency at every active cycle
    - Fixed 3: Clock coordinated at frequency of the core at every active cycle
- **PMC perf counters**: monitor 8 configurable performance counters using low overhead *rdpmc()* instruction

**Time Overhead: 0,59%**    **Memory Overhead?**

It is related to:
- Number of MPI processes
- Application time
- Number of MPI calls

*Example*: 16 MPI processes, 7.40 min of application time and 3,5 Mln of MPI calls

**Memory overhead**: ≈250 MB

**Average timing error wrt Intel Trace Analyzer: 0.45%**
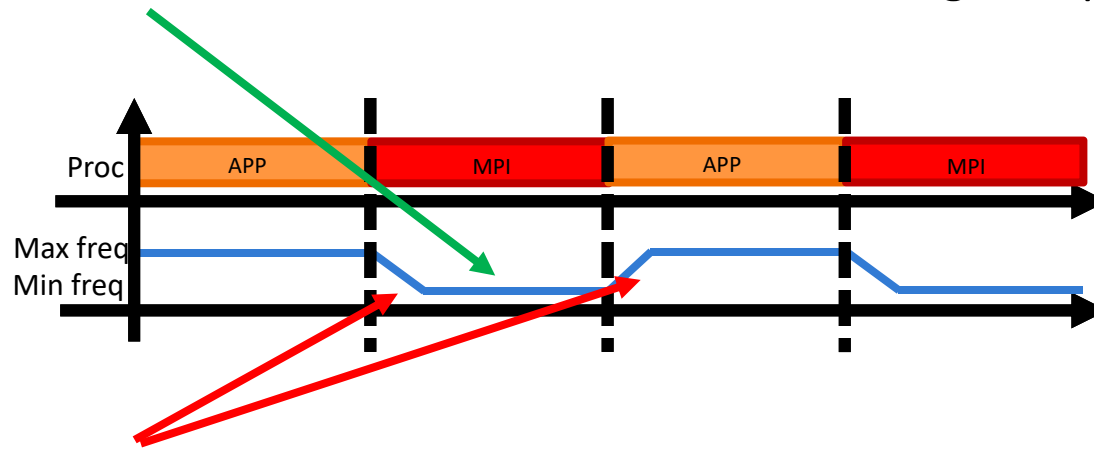
# APP time vs MPI Time

**Ndiag 1**  **Workload MPI root**: 10.25%  **AVG workload (no root)**: 5.98%
**Linear algebra is computed only by the root MPI → unbalanced workload**

**Ndiag 16**  **Workload MPI root**: 6.59% **AVG workload (no root)**: 6.23%
**Linear algebra is computed by all MPI processes → balanced workload**



MPI time is dominated by short phases

MPI time is dominated by long phases

**Idea**: use DVFS to slow down cores during MPI-phases

**Challenge**: Account for DVFS inertia, and appl. slowdown

# PMPI-based E2Sol minimization

If QE has significant percentage of MPI time with MPI phases longer than 500us

➡️ **PMPI needed to gauge and exploit (PMPI + PM) power saving opportunity**



Unbalanced benchmark on a single node (negligible MPI communication time)

**Up to 11% of energy** and **12% of power saved** with no impact on performance

# Outline

❑ **Power and Thermal Walls in HPC**

❑ **Power and Thermal Management**

❑ **Energy-efficient Hardware**

❑ **Conclusion**

# The era of Eterogeneous Architecture



Massive presence of accelerators in TOP500

Absolute dominance in GREEN500

# Recipe for Energy-efficient Acceleration

- Many (thousands) "simple" cores, managing FP units and special-function units for key workload patterns (stencil, tensor units) → maximize FP/mm2

- Non-coherent caches and lots of "non-cache" memory (registers for multithreading, scratchpad) → maximize "useful" Bit/mm2 for on-chip

- Large memory bandwidth based on tightly coupled-memory (HBM) → maximize GBps/mm2 for off-chip

- Low Operating voltage and moderate operating frequency → keep W/mm2 under control

- From 2D to 3D (now 2.5D)

**Is there room for differentiation, or are GP-GPUs the only answer?**

# Pezy-SC2 (top 1-2-3 GREEN500 Nov17)

Pezy-SC highlights:

- Technology (16nm TSMC) - 54% power reduction
- Advance and integrated power delivery – 30% power reduction
- Low voltage operation (0.7v) – 16% power reduction
- Low performance host processor – 15% power reduction



**Combines low-power design, simple (no legacy!) instruction set, advanced power management**

# Opportunity for (EU) HPC: open ISA

**RISC-V** **open** RISC ISA developed by UCB and supported now by the RISC-V foundation (riscv.org), with 70+ members (including, NVIDIA, IBM, QUALCOMM, MICRON, SAMSUNG, GOOGLE…)

- Reasonable, streamlined ISA → distills many years of research, conceived for efficiency not for legacy support
- Safe-to-use free ISA → freedom to operate (see RISC-V genealogy project), freedom to change/evolve/specialize, no licensing costs
- Wide community effort already on-going on tools, verification… → leverage this to jumpstart and compensate for our initial inertia
- Rapidly gaining traction in many application domains (IoT, big data) → large "dual-use" markets opportunity
- **Spec covers 64bit, vector ISA (on-going), 128bit (planned)**
- **HPC-profile RISC-V startups already active (esperanto.ai)**

# PULP: An Open Source Parallel Computing Platform

**Started in 2013 (UNIBO, ETHZ)**

Programming Model



Virtualization Layer

Compiler Infrastructure

Processor & Hardware IPs

Low-Power Silicon Technology

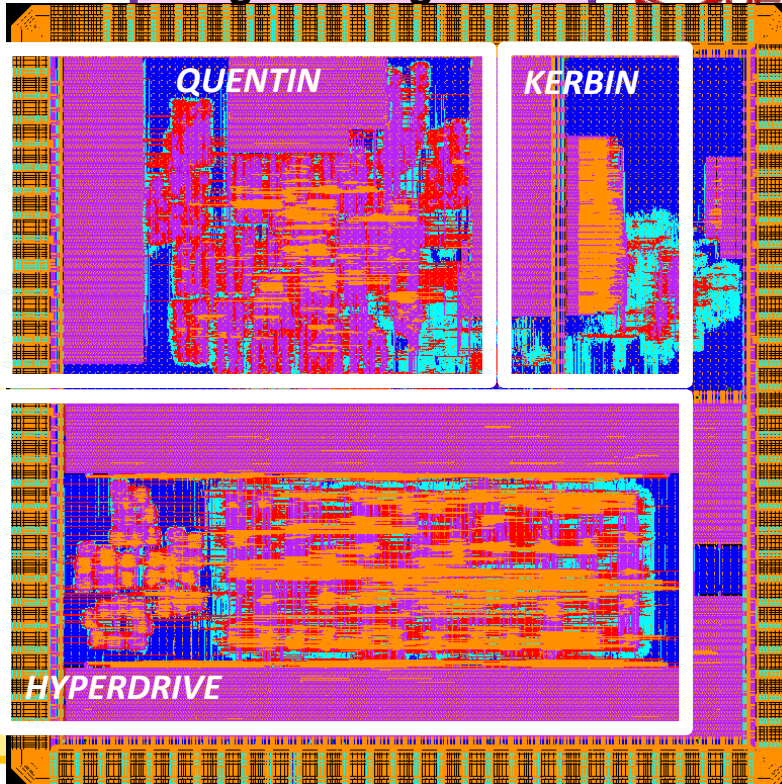PULP Hardware and Software released under Solderpad License

Used by tens of companies and universities, taped out in 14nm FINFET, 22FDX,...
64bit core "Ariane" + Platform to be launched in Q1 2018 (taped out in 22FDX)

# PULP: An Open Source Parallel Computing Platform

**Started in 2013 (UNIBO, ETHZ)**

Programming Model

QUENTIN

KERBIN

HYPERDRIVE

OpenVX

OpenMP

freeRTOS

Linux

GCC

LLVM COMPILER INFRASTRUCTURE
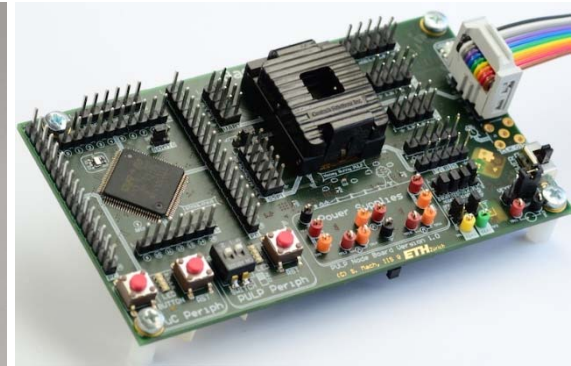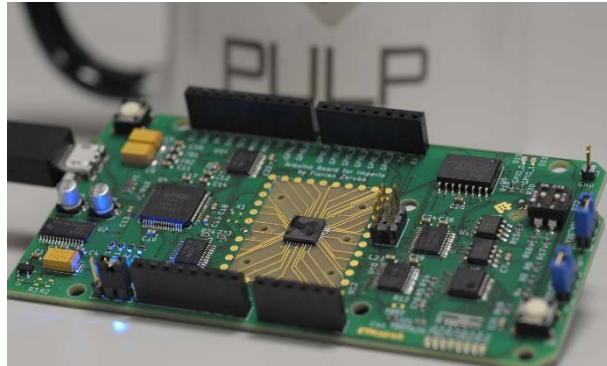
RISC-V

ST life.augmented

GLOBAL FOUNDRIES

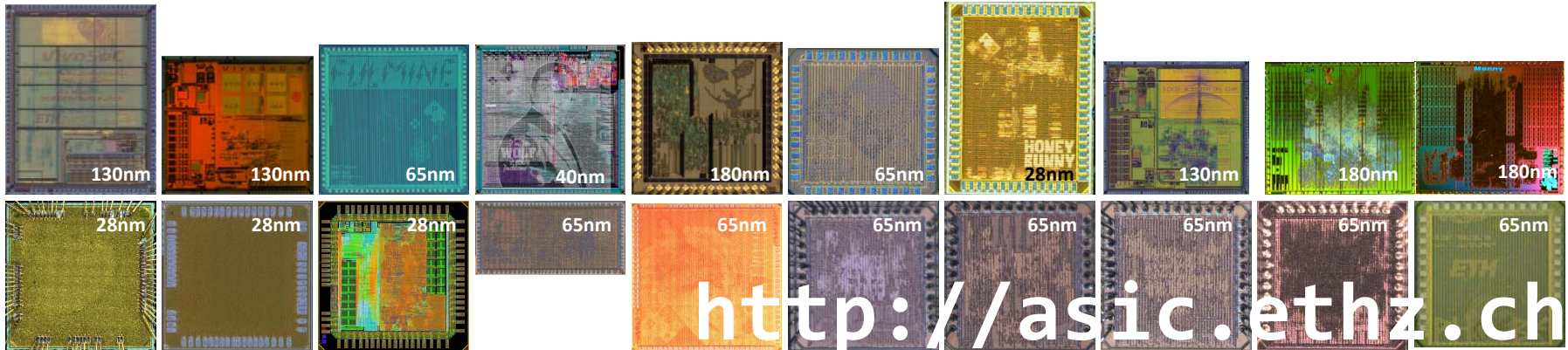PULP Hardware and Software released under Solderpad License

Used by tens of companies and universities, taped out in 14nm FINFET, 22FDX,…
64bit core "Ariane" + Platform to be launched in Q1 2018 (taped out in 22FDX)

# Thanks for your attention!



**www.pulp-platform.org**



http://asic.ethz.ch

## The fun is just beginning...