# A SOFTWARE APPROACH FOR ENERGY-EFFICIENT HPC WITH MERIC
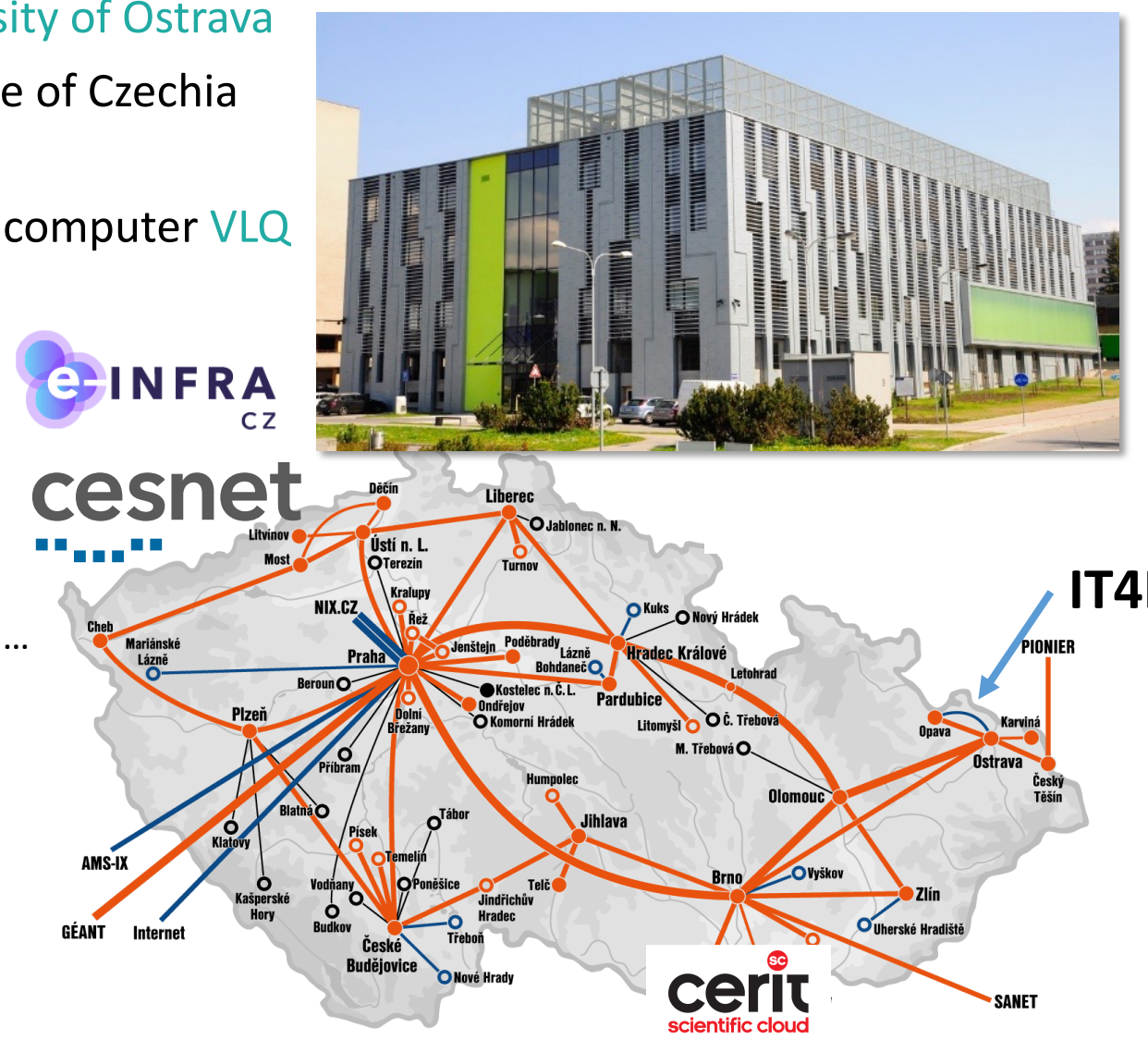
Ondrej Vysocky, Lubomir Riha, Tomas Kozubek

IT4Innovations National Supercomputing Centre, VSB-TU Ostrava, Ostrava, Czech Republic

https://code.it4i.cz/energy-efficiency/meric-suite/meric

## IT4INNOVATIONS NATIONAL SUPERCOMPUTING CENTER CZECH REPUBLIC

**About us:**
- Established in 2011 in Ostrava, at VSB – Technical University of Ostrava
- Member of e-INFRA CZ, a strategic research infrastructure of Czechia
  - co-operating LUMI supercomputer (TOP #9 in the world)
- operating Barbora & Karolina supercomputers, quantum computer VLQ
- co-operating LUMI supercomputer (TOP #9 in the world)

**Research, collaboration & training**
- 5 research laboratories, 120 FTE in HPC, HPDA, AI, QC
- Participating in EU HPC initiatives:
  - EuroHPC JU, EUDAT, ETP4HPC, BDVA, EOSC, QUIC, VI-HPS, WHPC, …
- Strong international collaboration
  - 25+ HE/DEP ongoing projects
- Cooperation with industry and public institutions
  - NCC in HPC, EDIH OVA, LUMI and Czech AI Factories
- Training and educational activities

## RUNTIME SYSTEM

**MERIC runtime system provides dynamic tuning of parallel applications running in the HPC environment**
- Performance and power aware
- lightweight & easy to install & easy to use
- C/C++ API, Fortran module, Python module
- MPI, OpenMP, CUDA parallelization

**Goals:**
- Application energy consumption measurement
- Application dynamism & energy efficiency analysis
- Dynamic HW power knobs tuning for energy savings
- HW & SW power management co-design

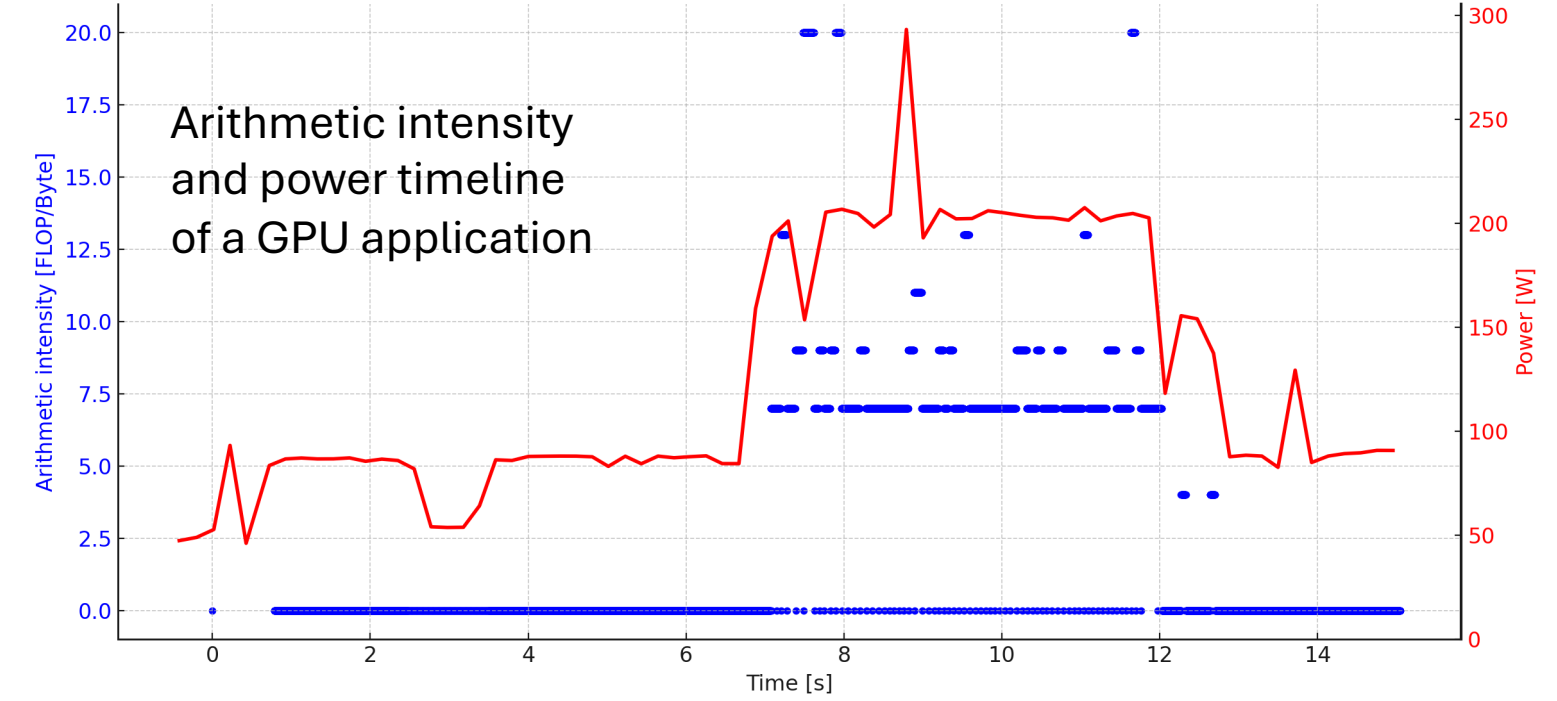**Support for a wide range of architectures**
- 86, IBM OpenPOWER, ARM
- Nvidia/AMD GPUs

**Power monitoring systems**
- CPU: Intel/AMD RAPL, IBM OCC, A64FX, HWMON (Nvidia Grace)
- GPU: NVML, ROCm
- System: ATOS HDEEM

**Performance parameter tuning**
- CPU frequency, GPU SM frequency, memory frequency, power limit, number of active CPU cores

## SAMPLING-BASED GPU TUNING

**CUDA energy efficiency runtime system**
- Realtime monitoring of GPU utilization
  - CUPTI PM Sampling API is used to collect SM utilization & memory activity metrics
- Arithmetic intensity modeling
- Dynamic frequency tuning
  - On A100 SXM-4, 50ms between freq. configuration changes
- Special daemon tool for GPU frequency tuning


Arithmetic intensity and power timeline of a GPU application

| GPU | Switching latency range [ms] | Transition latency range [ms] |
|---|---|---|
| RTX Quadro 6000 | 0.55 - 350.4 | 0.09 - 335.8 |
| A100 SXM-4 | 4.43 - 22.7 | 0.11 - 11.5 |
| GH200 | 4.91 - 477.3 | 0.08 - 471.1 |

**LATEST tool**
- Evaluation of GPU frequency change latency
- Utilization of synthetic, frequency-sensitive workload
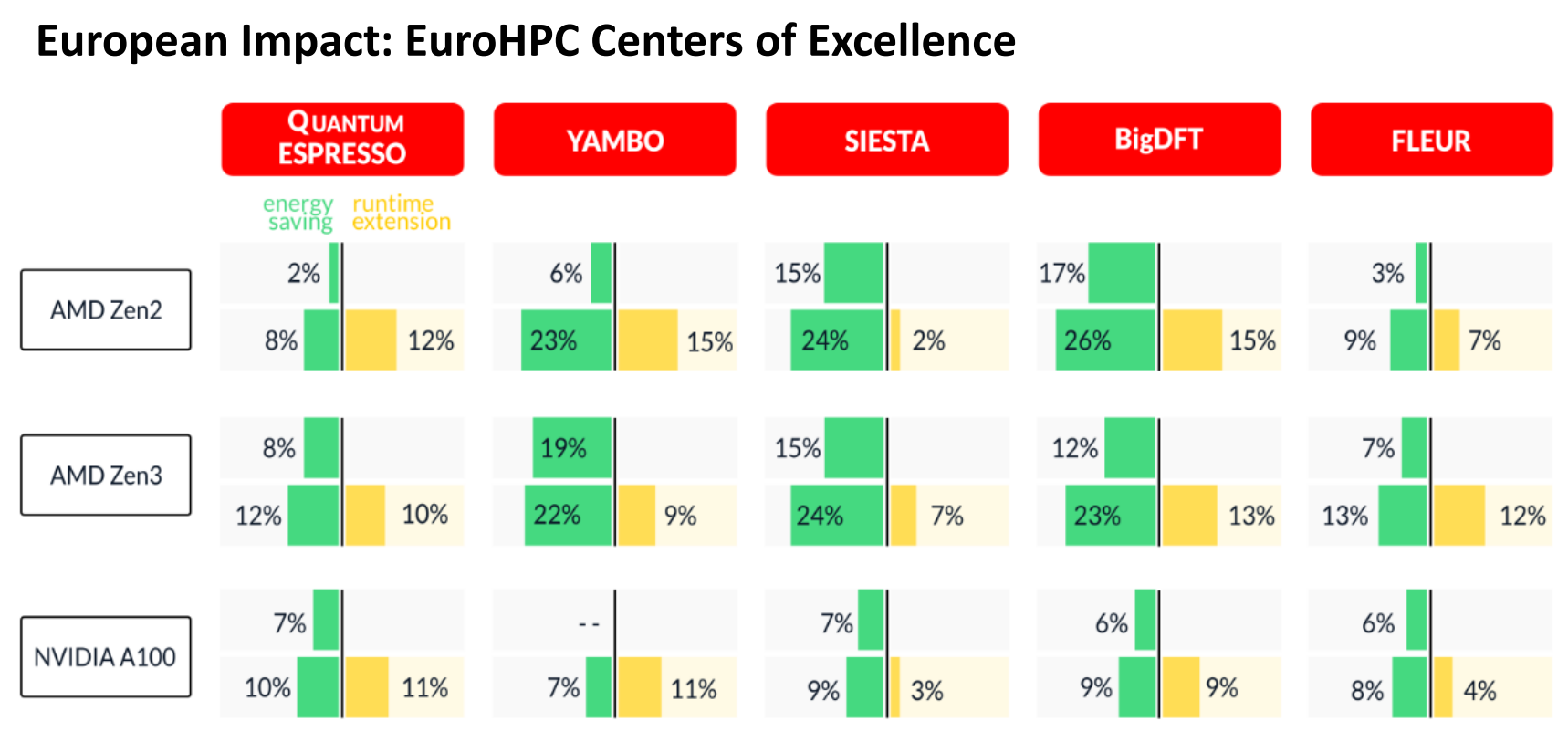- Analysis of the frequency transition – for each freq. pair

## MERIC: ENERGY EFFICIENCY SW SUITE FOR HPC

**1.) Parallel application behavior analysis & optimization**

AIRBUS — customer
SCALABLE — EuroHPC JU project
ProLB — code

| | no penalty | w. penalty |
|---|---|---|
| Runtime [%] | +0.5% | +4.1% |
| Energy [%] | -12.1% | -19.5% |

ES GROUP    RENAULT

**European Impact: EuroHPC Centers of Excellence**



| | QUANTUM ESPRESSO | YAMBO | SIESTA | BigDFT | FLEUR |
|---|---|---|---|---|---|
| AMD Zen2 | 2% / 8% / 12% | 6% / 23% / 15% | 15% / 24% / 7% | 17% / 26% / 9% | 3% / 15% / 7% |
| AMD Zen3 | 8% / 12% / 10% | 19% / 22% / 9% | 15% / 24% / 7% | 12% / 23% / 13% | 7% / 15% / 12% |
| NVIDIA A100 | 7% / 10% / 11% | -- / 7% / 11% | 9% / 3% / 9% | 6% / 9% / 9% | 6% / 8% / 4% |

**2.) HW&SW co-design for energy efficiency**

dare — European RISC-V processor & accelerators

EUPEX — European modular Exascale-ready pilot system

**3.) Datacenter monitoring & optimization**

# K A R 0 L 1 N A

- Karolina system – power consumption ±780 kW
- **103 kW power savings** equals to 883 MWh / year
- 1MWh ~ 6000 CZK => 5.4 M CZK ~ 250 000 USD
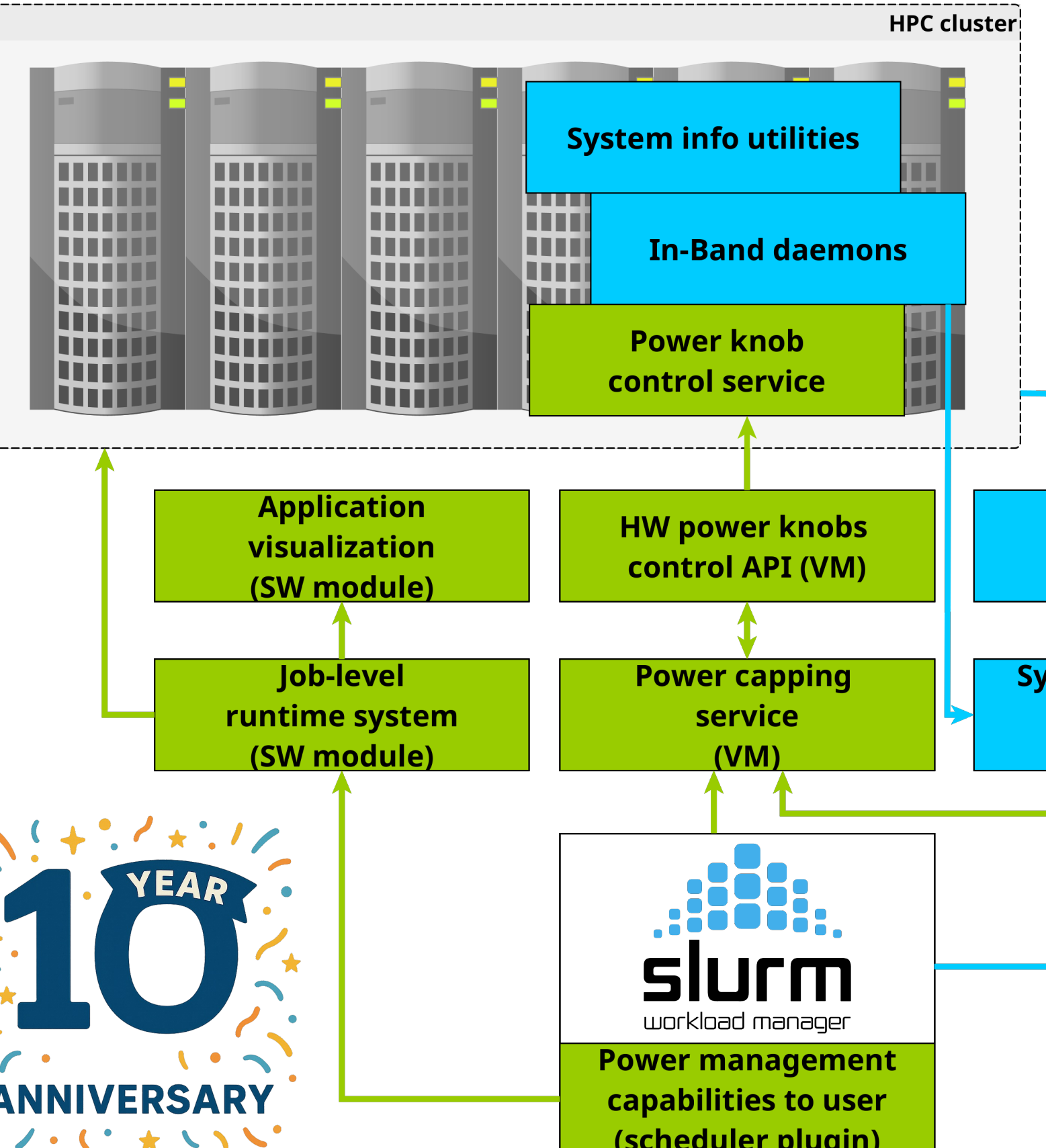- annual 315t CO2e reduction equals 12 600 trees
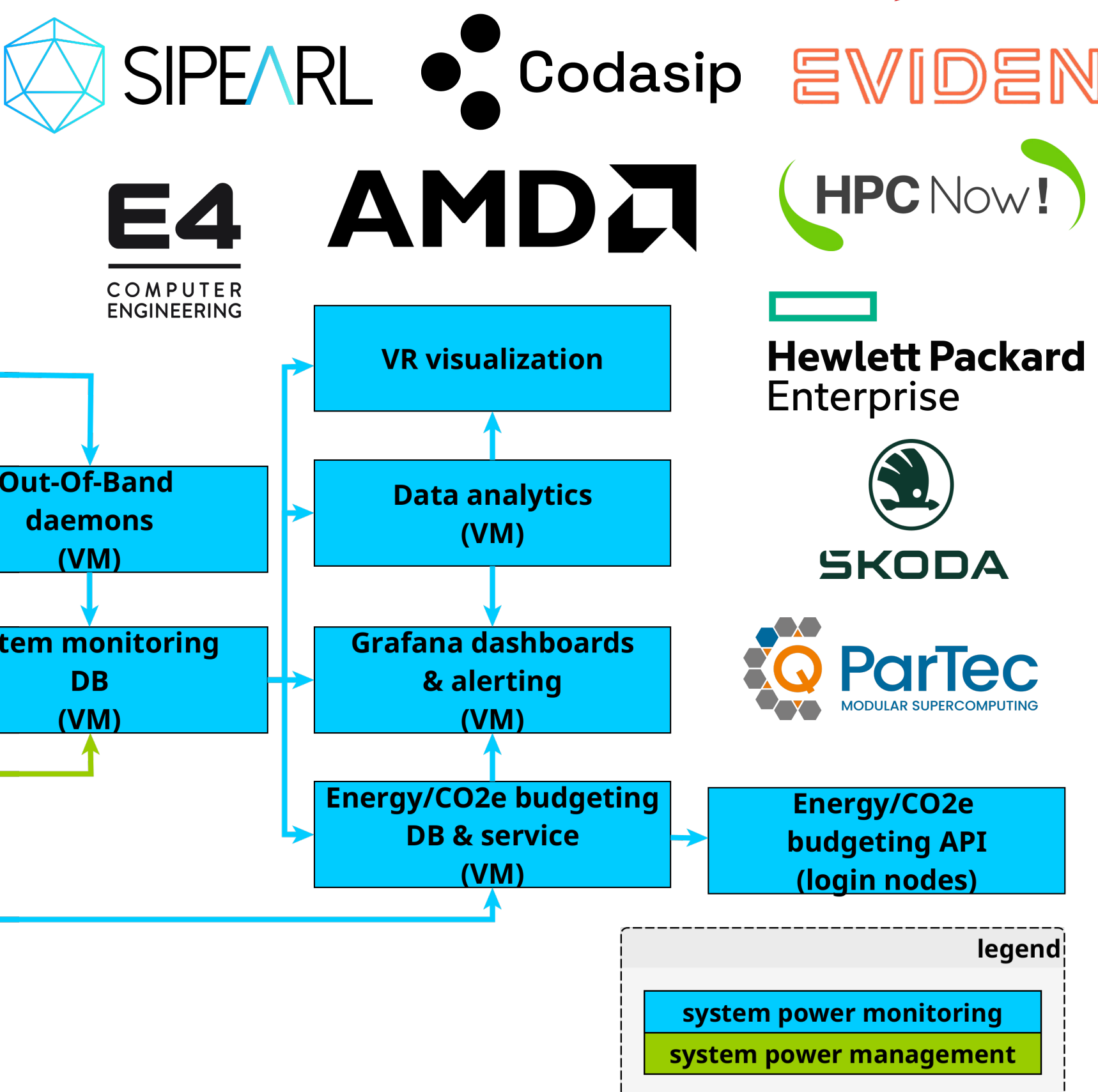
x 1 000

**Contractual research**

# DEUCALION

EuroHPC system in Portugal - deployment of MERIC suite

### SOFTWARE STACK



System info utilities
In-Band daemons
Power knob control service
Application visualization (SW module)
HW power knobs control API (VM)
Out-Of-Band daemons (VM)
Job-level runtime system (SW module)
Power capping service (VM)
System monitoring DB (VM)
VR visualization
Data analytics (VM)
Grafana dashboards & alerting (VM)
Energy/CO2e budgeting DB & service (VM)
Energy/CO2e budgeting API (login nodes)

slurm workload manager
Power management capabilities to user (scheduler plugin)

**10 YEAR ANNIVERSARY**

legend:
system power monitoring
system power management

### INDUSTRIAL COLABORATION

EXELIZ SOLUTIONS
FUJITSU
EVIDEN
SIPEARL
Codasip
HPC Now!
E4 COMPUTER ENGINEERING
AMD
Hewlett Packard Enterprise
ŠKODA
ParTec MODULAR SUPERCOMPUTING

## MERIC DEVELOPMENT FUNDING

READEX Runtime Exploitation of Application Dynamism for Energy-efficient eXascale computing
e-INFRA CZ
SEANERGYS Energy Efficient Exascale

## ENERGY-EFFICIENCY SERVICES

EUPEX European Pilot for Exascale
POP
dare
SCALABLE
MAX DRIVING THE EXASCALE TRANSITION

## ENERGY EFFICIENCY SERVICES

- How much energy does my application consume? What is its carbon footprint?
- Which parts of the code are power hungry? Does it activate power capping?
- How energy efficient the code is?
- Which hardware platform is the most energy efficient for my code?
- Which parts of the application may give opportunity for energy savings?
- How much energy can be saved by static versus dynamic tuning of power management knobs without impacting application performance? And if the performance penalty is 5%, 10%, … ?
- Does my hardware power/thermal management work as intended?
- When is the capping mechanism a performance-limiting factor?

VI-HPS

MAX — DRIVING THE EXASCALE TRANSITION

## MAX3 CENTER OF EXCELLENCE

**Hardware platform used for energy efficiency evaluation**

| Code | Instrumented for static tuning and compiled with MERIC | IT4I Barbora CPU partition Intel CascadeLake | EuroHPC Karolina CPU partition AMD Zen2 | EuroHPC Karolina GPU partition CPU AMD Zen3 | EuroHPC Karolina GPU partition Nvidia A100 GPU | Intel Sapphire Rapids CPU w. DDR / HBM | IBM Power 10 (S1022) | Fujitsu A64FX |
|---|---|---|---|---|---|---|---|---|
| Yambo | ✓ | ✓ | ✓ | ✓ | ✓ | ✓✓ | -- | -- |
| Quantum ESPRESSO | ✓ | ✓ | ✓ | ✓ | ✓ | ✓✓ | ✓ | -- |
| Siesta | ✓ | ✓ | ✓ | ✓ | ✓✓ | ✓✓ | -- | -- |
| BigDFT | ✓ | ✓ | ✓ | ✓ | ✓ | ✓✓ | -- | -- |
| Fleur | ✓ | ✓ | ✓ | ✓ | ✓ | ✓✓ | ✓ | ✓ |

### THE MOST ENERGY EFFICIENT PLATFORMS FOR MAX CODES

| SPR + HBM = 4.56 | SPR + HBM = 1.33 | SPR + HBM = 2.57 | SPR + HBM = 0.28 | SPR + HBM AMD Zen2 = 1.82 | Arithmetic intensity and power timeline of a GPU application |

Efficiency (GFLOPs/W) of the most efficient HW with no runtime extension

### ENERGY EFFICIENCY EVALUATION FOR FLEUR CODE ON SELECTED PLATFORM

| Hardware | Energy efficiency | Node energy consumption | Monitoring system | HW configuration | Runtime |
|---|---|---|---|---|---|
| AMD Zen2 (Rome) | 1.78 GFLOPs/W | 53.36 kJ | AMD RAPL + baseline | default | 109 s (100%) |
| | 1.82 GFLOPs/W | 52.00 kJ (-3%) | | CF 2.9 GHz | 101% |
| | 1.94 GFLOPs/W | 48.81 kJ (-9%) | | CF 2.1 GHz | 107% |
| AMD Zen3 (Milan) | 1.67 GFLOPs/W | 56.96 kJ | AMD RAPL + baseline | default | 93 s (100%) |
| | 1.79 GFLOPs/W | 53.05 kJ (-7%) | | CF 2.7 GHz | 101% |
| | **1.91 GFLOPs/W** | 49.73 kJ (-13%) | | CF 2.0 GHz | 112% |
| Intel Cascade lake | 1.00 GFLOPs/W | 94.94 kJ | HDEEM | default | 217 s (100%) |
| | 1.04 GFLOPs/W | 91.26 kJ (-4%) | | CF 2.8 GHz, UCF 2.2 GHz | 101% |
| | 1.13 GFLOPs/W | 84.51 kJ (-11%) | | CF 1.9 GHz, UCF 1.8 GHz | 123% |
| Intel Sapphire Rapids w. HBM | 1.77 GFLOPs/W | 73,31 kJ | RAPL + baseline | default | 82 s (100%) |
| | **1.82 GFLOPs/W** | 71,83 kJ (-2%) | | CF 3.1 GHz, UCF 1.8 GHz | 101% |
| | **1.82 GFLOPs/W** | 71.83 kJ (-2%) | | CF 3.1 GHz, UCF 1.8 GHz | 101% |
| Intel Sapphire Rapids w. DDR memory | 1.43 GFLOPs/W | 90.22 kJ | RAPL + baseline | default | 100 s (100%) |
| | 1.47 GFLOPs/W | 88.48 J (-2%) | | CF 2.9 GHz, UCF 2.0 GHz | 101% |
| | 1.54 GFLOPs/W | 86.50 kJ (-4%) | | CF 2.3 GHz, UCF 1.8 GHz | 110% |
| Nvidia A100 | -- | 180.6 kJ | AMD RAPL + NVML + baseline | default | 111 s (100%) |
| | -- | 169.26 kJ (-6%) | | 1230 MHz | 101% |
| | -- | 166.3 kJ (-8%) | | 990 MHz | 104% |
| IBM Power10 | 0.459 GFLOPs/W | 198.6 kJ | PDU | default | 199 s |
| A64FX | 0.321 GFLOPs/W | 282.5 kJ | perf. counters + baseline | default | 812 s |

## SPACE CENTER OF EXCELLENCE

| Code energy efficiency | NVidia GRACE CPU [MFLOPS/W] | Intel Sapphire Rapids with DDR [MFLOPS/W] | Intel Sapphire Rapids with HBM [MFLOPS/W] |
|---|---|---|---|
| Pluto | 805.4 | 264.0 | 309.0 |
| OpenGADGET | 716.2 | 138.2 | 149.2 |
| iPIC3D | 791.3 | 238.2 | 321.4 |
| RAMSES | 854.9 | 399.7 | 417.8 |
| BHAC | 292.9 | 121.9 | 125.5 |
| FIL | 522.2 | 223.9 | 248.9 |
| ChaNGa | 1478.9 | 779.6 | 1018.3 |

Results are shown for hardware configurations **generating very small runtime extensions (less than 5% of default).**

| Code energy efficiency | NVidia GRACE CPU [MFLOPS/W] | Intel Sapphire Rapids with DDR [MFLOPS/W] | Intel Sapphire Rapids with HBM [MFLOPS/W] |
|---|---|---|---|
| Pluto | 964.4 | 266.2 | 314.1 |
| OpenGADGET | 940.7 | 138.2 | 151.7 |
| iPIC3D | 964.4 | 242.6 | 321.4 |
| RAMSES | 968.6 | 401.6 | 418.5 |
| BHAC | 340.7 | 127.8 | 125.5 |
| FIL | 600.0 | 228.1 | 253.0 |
| ChaNGa | 1774.7 | 825.5 | 1026.5 |

Results are shown for hardware configurations **generating maximum energy savings and noticeable runtime extension.**

**Energy consumption reduction [kJ] / Runtime [s];**
- 1st row - results with no runtime extension (runtime close to 100% of default);
- 2nd row - results for maximum energy savings.

| Code / System | Pluto | OpenGADGET | iPIC3D | RAMSES | BHAC | FIL | ChaNGa |
|---|---|---|---|---|---|---|---|
| Nvidia A100 | -6% / 103% | -7% / 102% | -3% / 104% | – | – | -6% / 102% | -14% / 103% |
| | -9% / 113% | -7% / 102% | -5% / 111% | – | – | -7% / 103% | -20% / 107% |
| SPR w. DDR | -9% / 102% | -7% / 102% | -7% / 102% | -6% / 102% | -10% / 103% | -6% / 100% | -12% / 100% |
| | -10% / 106% | -7% / 102% | -9% / 108% | -7% / 104% | -14% / 110% | -6% / 104% | -13% / 103% |
| SPR w. HBM | -4% / 101% | -9% / 94% | -7% / 101% | -7% / 102% | -4% / 99% | -8% / 102% | -16% / 102% |
| | -6% / 105% | -11% / 98% | -7% / 101% | -8% / 104% | -4% / 99% | -6% / 104% | -30% / 135% |
| Grace CPU | -22% / 101% | -13% / 103% | -9% / 103% | -19% / 101% | -26% / 103% | -8% / 102% | -16% / 103% |
| | -35% / 122% | -33% / 128% | -29% / 126% | -28% / 137% | -36% / 109% | -20% / 117% | -30% / 135% |
| Cascade Lake | -6% / 102% | -9% / 103% | -6% / 101% | -7% / 102% | -5% / 102% | -3% / 102% | -30% / 102% |
| | -12% / 126% | -13% / 115% | -13% / 115% | -11% / 123% | -11% / 118% | -13% / 127% | -36% / 110% |

### PLUTO: DETAILED A100 GPU ANALYSIS INCLUDING STRONG SCALING

| SM freq. [MHz] | 1 node 8 GPUs | 2 nodes 16 GPUs | 4 nodes 32 GPUs | 8 nodes 64 GPUs | 16 nodes 128 GPUs | 32 nodes 256 GPUs |
|---|---|---|---|---|---|---|
| 1410 | 129 s | 68,2 s | 38,1 s | 21,6 s | 15,5 s | 11,2 s |
| default | 100% | 100% | 100% | 100% | 100% | 100,0% |
| 1350 | 102,8% | 102,3% | 102,8% | 102,6% | 100,5% | 100,8% |
| 1290 | 105,8% | 105,4% | 105,3% | 105,2% | 102,1% | 101,1% |
| 1230 | 108,9% | 108,4% | 108,4% | 108,0% | 104,8% | 103,7% |
| 1170 | 112,6% | 112,1% | 111,4% | 113,1% | 108,8% | 107,5% |
| 1110 | 116,2% | 115,3% | 115,5% | 115,9% | 112,9% | 110,0% |
| 1050 | 120,7% | 119,9% | 119,7% | 121,2% | 116,3% | 115,2% |
| 990 | 125,7% | 125,0% | 124,1% | 123,9% | 120,5% | 117,6% |

| SM freq. [MHz] | 1 node 8 GPUs | 2 nodes 16 GPUs | 4 nodes 32 GPUs | 8 nodes 64 GPUs | 16 nodes 128 GPUs | 32 nodes 256 GPUs |
|---|---|---|---|---|---|---|
| 1410 | 365 kJ | 368 kJ | 387 kJ | 423 kJ | 543 kJ | 719 kJ |
| default | 100% | 100% | 100% | 100% | 100% | 100% |
| 1350 | 95,7% | 94,9% | 96,2% | 96,4% | 95,7% | 96,4% |
| 1290 | 92,0% | 92,2% | 93,5% | 93,7% | 93,7% | 93,5% |
| 1230 | 90,1% | 90,4% | 91,7% | 91,8% | 91,9% | 93,7% |
| 1170 | 88,8% | 89,4% | 90,5% | 91,8% | 91,8% | 93,0% |
| 1110 | 87,6% | 88,4% | 90,2% | 90,9% | 92,1% | 92,8% |
| 1050 | 87,4% | 88,5% | 90,4% | 91,7% | 92,4% | 94,5% |
| 990 | 88,9% | 90,5% | 92,1% | 92,5% | 94,3% | 95,5% |

**Impact of the static tuning of the GPU SM frequency on the runtime (left) and energy consumption (rigth) of the 3D Orszag-Tang vortex benchmark.**
- the left panel shows runtime variations with respect to the default execution time shown in the first line.
- the panel shows relative energy consumption with respect to the energy consumption of the default execution