



703308 VO High-Performance Computing WS2022/2023

Energy in HPC

Philipp Gschwandtner

Overview

- ▶ motivation / indications
 - ▶ Green500 & accelerators
- ▶ basics
 - ▶ power & energy
 - ▶ instrumentation, measurement and modeling
 - ▶ control mechanisms
- ▶ consequences & applications
 - ▶ multi-objective optimization, workload co-scheduling, MPI slack optimization, etc.

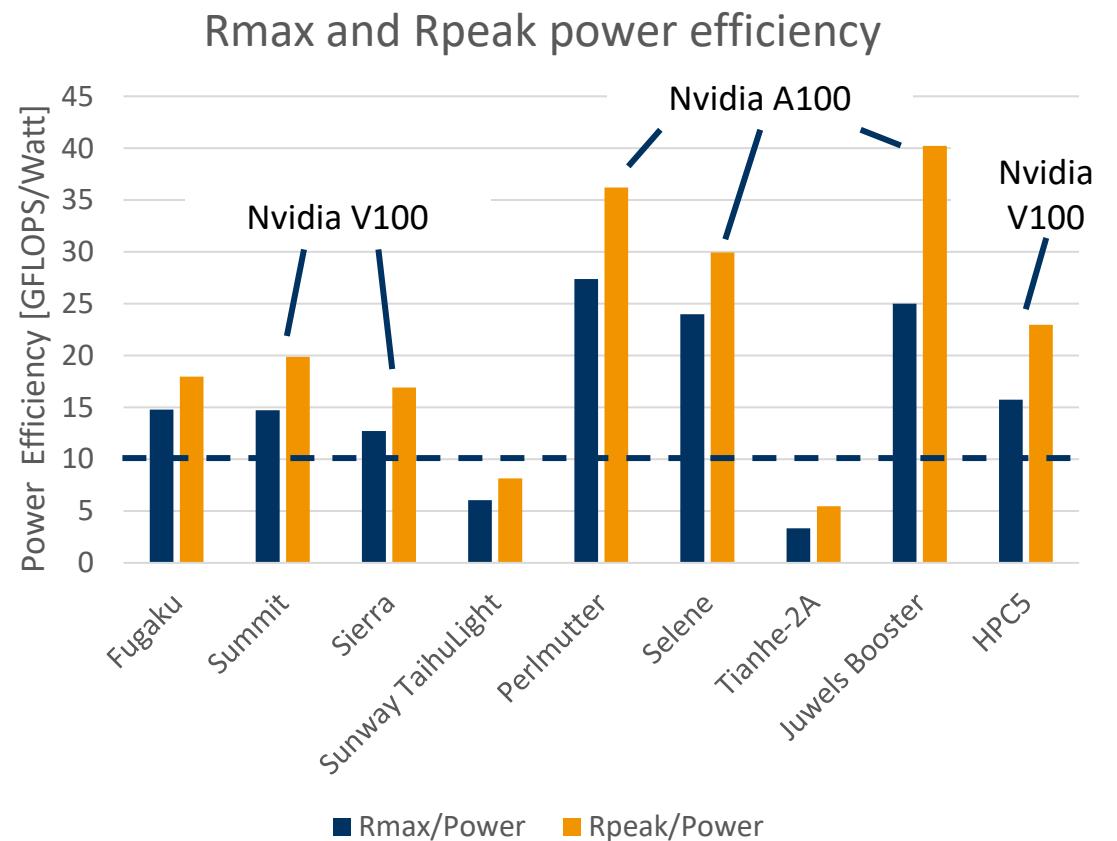
TOP500 & Green500

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <u>AMD Instinct MI250X</u> , Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	8,730,112	1,102.00	1,685.65	21,100
2	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
3	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <u>AMD Instinct MI250X</u> , Slingshot-11, HPE EuroHPC/CSC Finland	2,220,288	309.10	428.70	6,016
4	Leonardo - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, <u>NVIDIA A100 SXM4 40 GB</u> , Quad-rail NVIDIA HDR100 Infiniband, Atos EuroHPC/CINECA Italy	1,463,616	174.70	255.75	5,610
5	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, <u>NVIDIA Volta GV100</u> , Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148.60	200.79	10,096

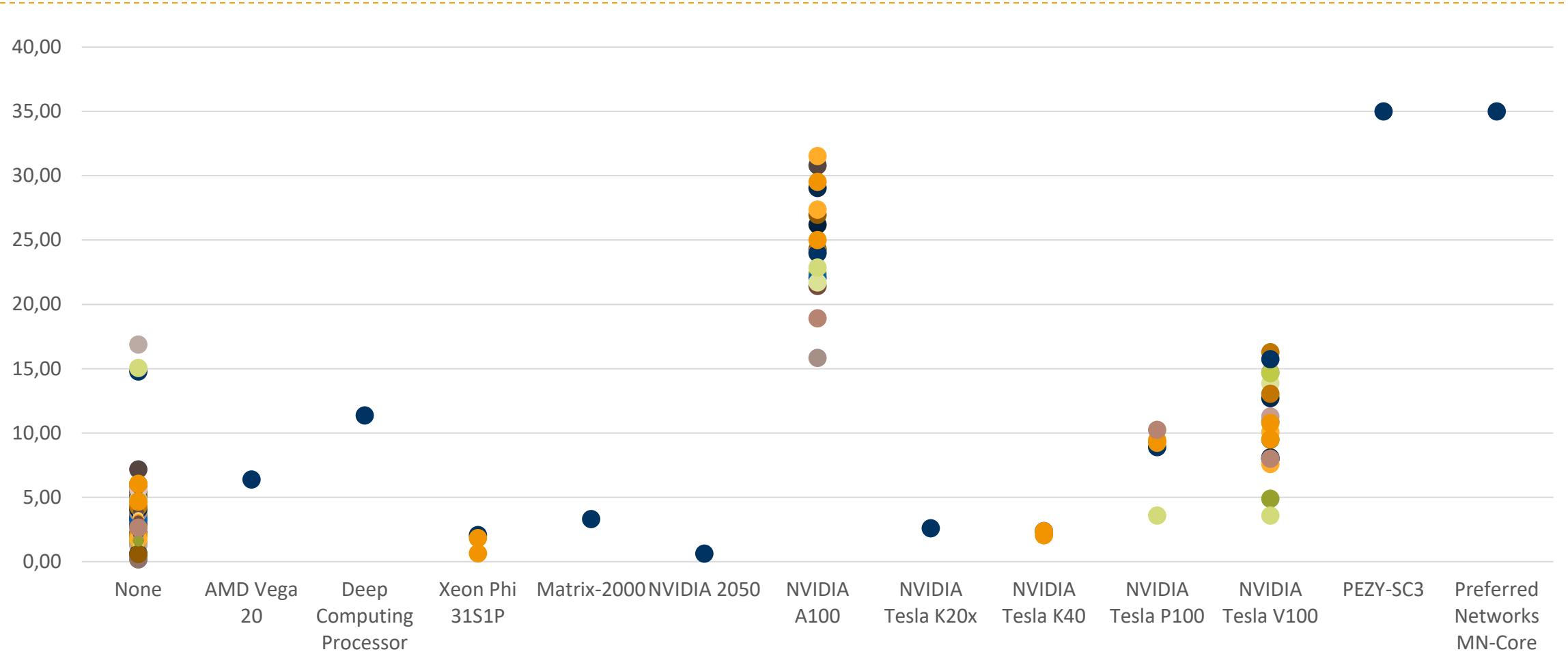
Rank	TOP500 Rank	System	Cores	Rmax (PFlop/s)	Power (kW)	Energy Efficiency (GFlops/watts)
1	405	Henri - Lenovo ThinkSystem SR670 V2, Intel Xeon Platinum 8362 2800Mhz (32C), <u>NVIDIA H100 80GB PCIe</u> , Infiniband HDR, Lenovo Flatiron Institute United States	5,920	2.04	31	65.091
2	32	Frontier TDS - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <u>AMD Instinct MI250X</u> , Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	120,832	19.20	309	62.684
3	11	Adastra - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, <u>AMD Instinct MI250X</u> , Slingshot-11, HPE Grand Equipement National de Calcul Intensif - Centre Informatique National de l'Enseignement Supérieur (GENCI-CINES) France	319,072	46.10	921	58.021

Why are we using accelerators?

- ▶ All top 10 systems above 10 GFLOPS/Watt use accelerators
 - ▶ Nov '21 data
- ▶ Exceptions:
 - ▶ Fugaku: ARM-based, no accelerators
 - ▶ Tianhe-2A: Matrix 2000 accelerators (128 core RISC CPUs)

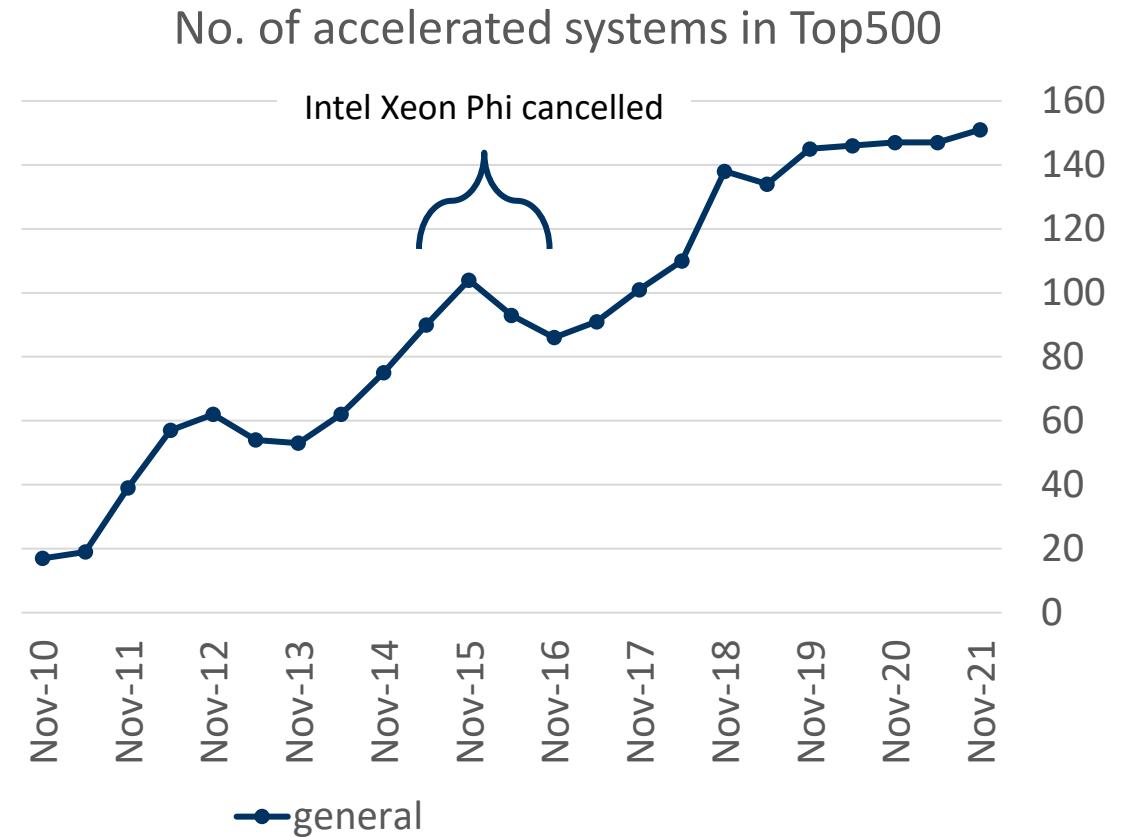


Power Efficiency of all Top 500 Systems



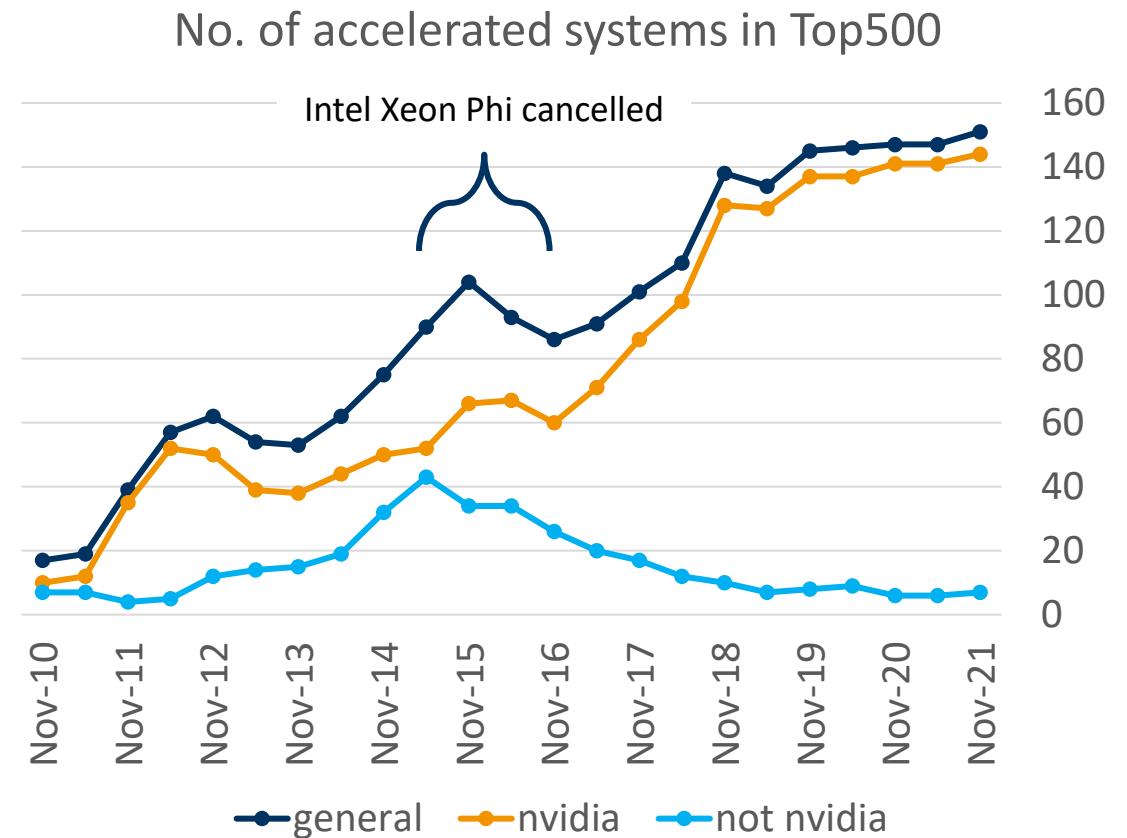
Market share

- ▶ Accelerator market share in HPC has been steadily increasing and will likely continue to do so
 - ▶ 10 out of top 10 on Green500 list (Nov. 2022)
- ▶ Application developers **need** to use accelerators to get high performance on modern systems



Market share cont'd

- ▶ Problem: No market competition
 - ▶ Nvidia is predominant
 - ▶ originally also Cell processor (Playstation 3!) and Intel Xeon Phi
 - ▶ disappeared since ~2015
 - ▶ Right now (2022): some AMD comeback, new Intel attempts
- ▶ Nvidia encourages using CUDA
 - ▶ vendor lock-in
- ▶ There are alternatives!
 - ▶ SYCL, ROCm/HIP, OpenMP, etc.



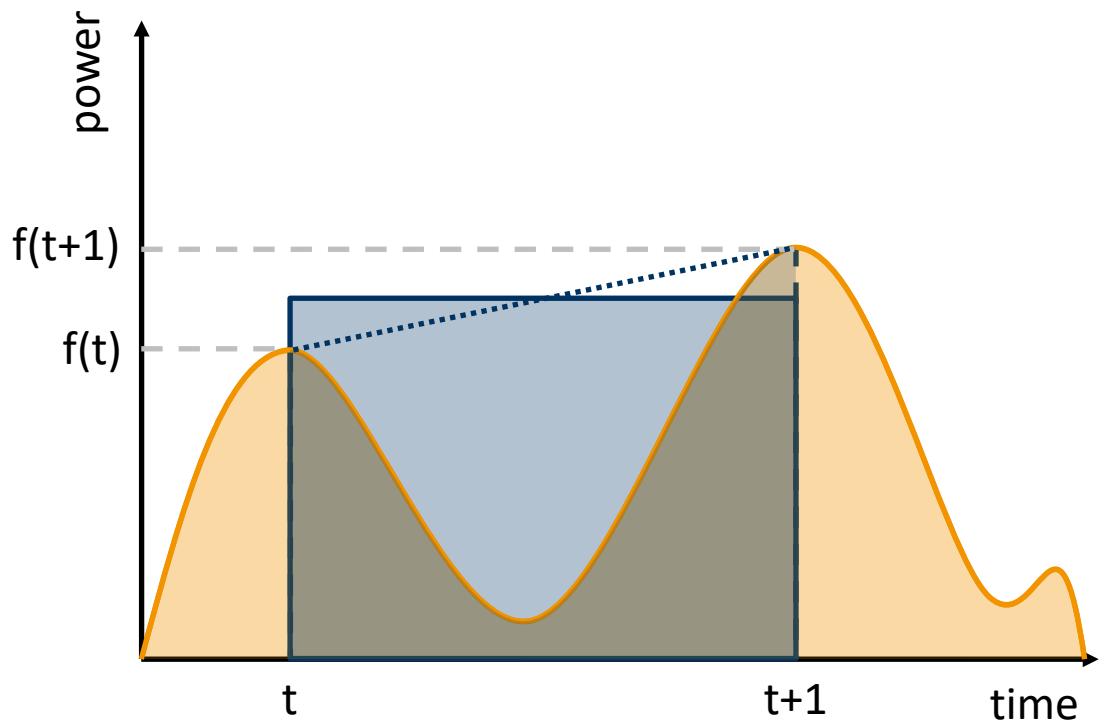
Basics

Disclaimer

- ▶ According to thermodynamics, there is no power or energy “consumption”
 - ▶ We are just efficiently converting electricity to heat
 - ▶ Computational result is a side-product
- ▶ We’re still going to call it “power consumption” and “energy consumption” for practical purpose

Power vs. energy

- ▶ Power is instantaneous, i.e. measured at a specific point in time
 - ▶ E.g. 10 W (watts)
 - ▶ does not have a time component
- ▶ energy is power over time ($E = \int P$)
 - ▶ E.g. 10 Wh (watt-hours) or 36 kJ – kiloJoules
 - ▶ Could be measured using analogue means
 - ▶ Often just power sampling ($E \cong P_{avg} \times T$)
 - ▶ Often requires very high temporal resolution for reasonable accuracy



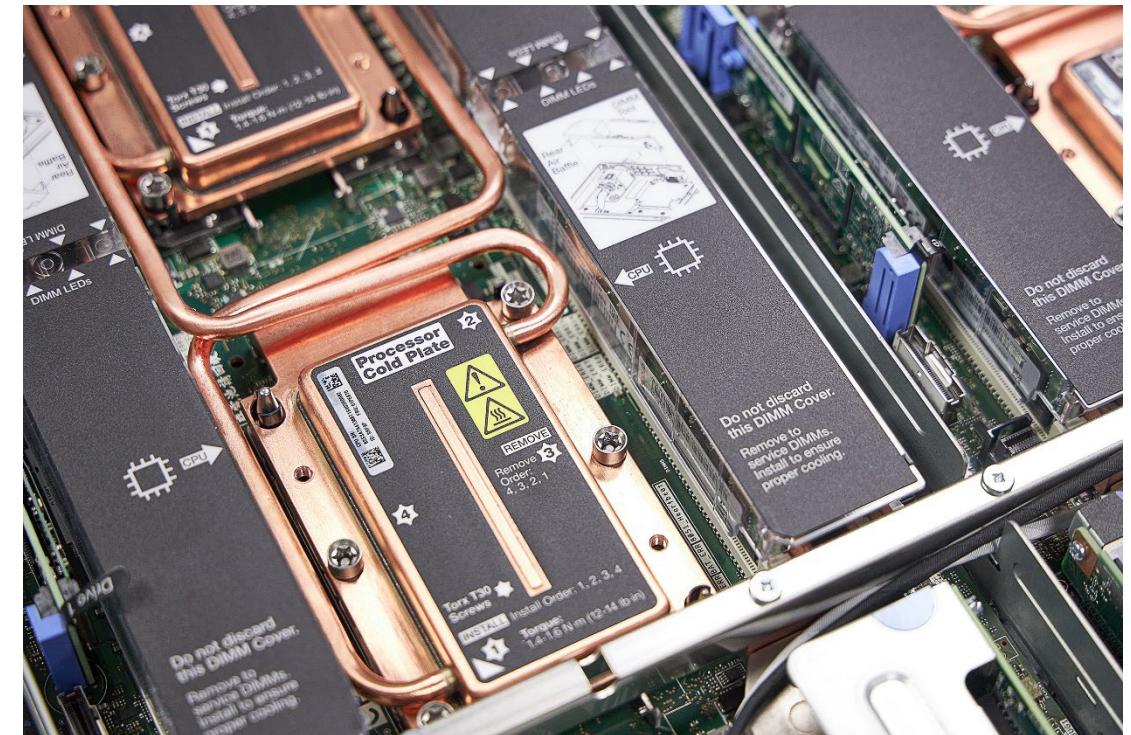
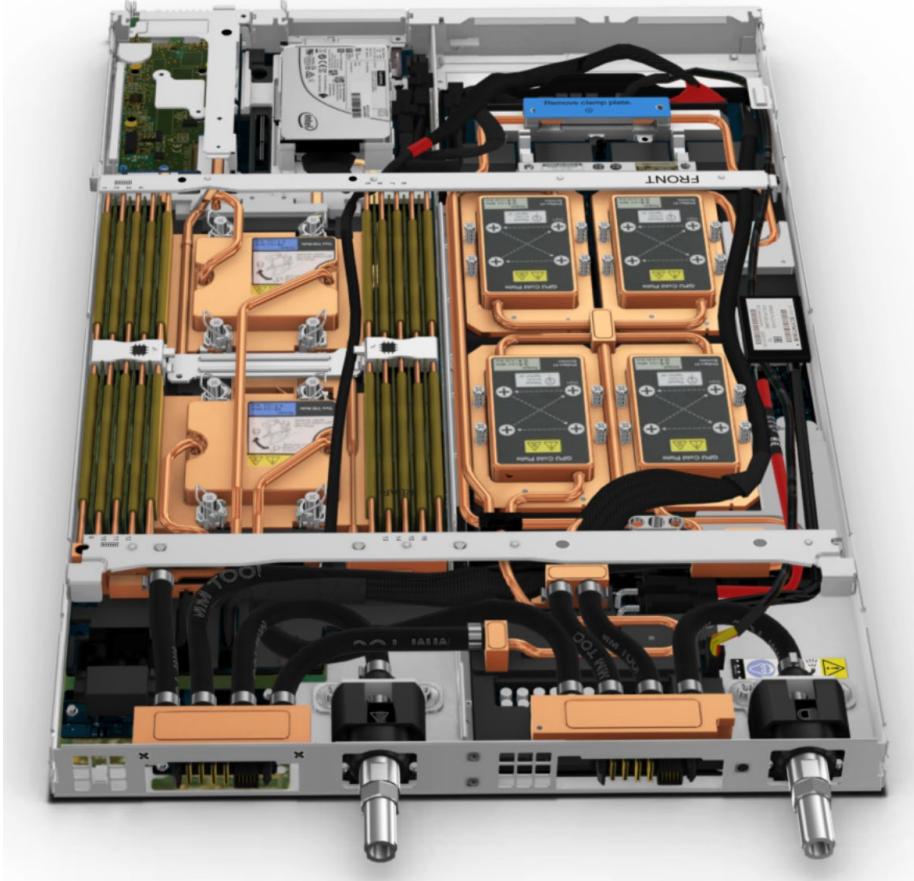
Power consumption of integrated circuits

- ▶ $P_{total} = P_{short-circuit} + P_{static} + P_{dynamic}$
 - ▶ $P_{short-circuit}$: Power due to short-circuit current during transistor switching
 - ▶ P_{static} : Power due to leakage current, increases with decreasing feature sizes
 - ▶ $P_{dynamic} = C \times F \times V^2 \times \alpha \approx F^3$
(C ... capacitance, F ... frequency, V ... voltage, α ... switching factor)
 - ▶ Frequency and voltage are tightly connected
 - ▶ Hence: Sometimes referred to as “cube rule”

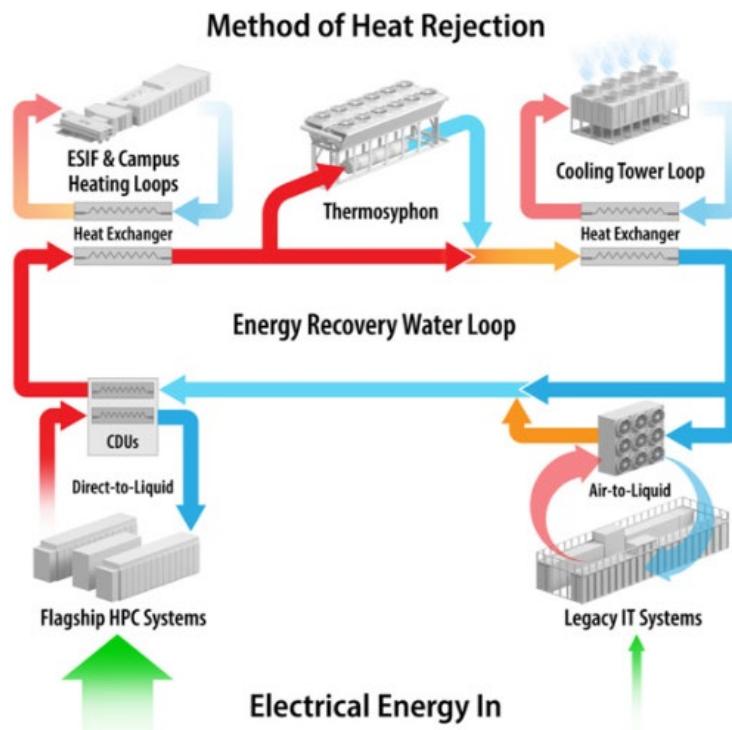
Power consumption of supercomputers

- ▶ Cores are not everything
 - ▶ Off-core entities: shared caches, memory controllers, system agents, etc.
 - ▶ Off-chip entities: RAM, mainboard, NIC, etc.
- ▶ Computing nodes are not everything
 - ▶ Network, storage, management, etc.
 - ▶ Cooling system
 - ▶ Lights, office equipment, etc.
- ▶ Efficiency measured via PUE – Power Usage Effectiveness
 - ▶ Ratio of power required by supercomputer vs. power for its entire facility
 - ▶ E.g. SuperMUC-NG @ LRZ: PUE of 1.08
 - ▶ 10-20 years ago: PUEs as high as 2-3
 - ▶ Mostly caused by cooling overhead

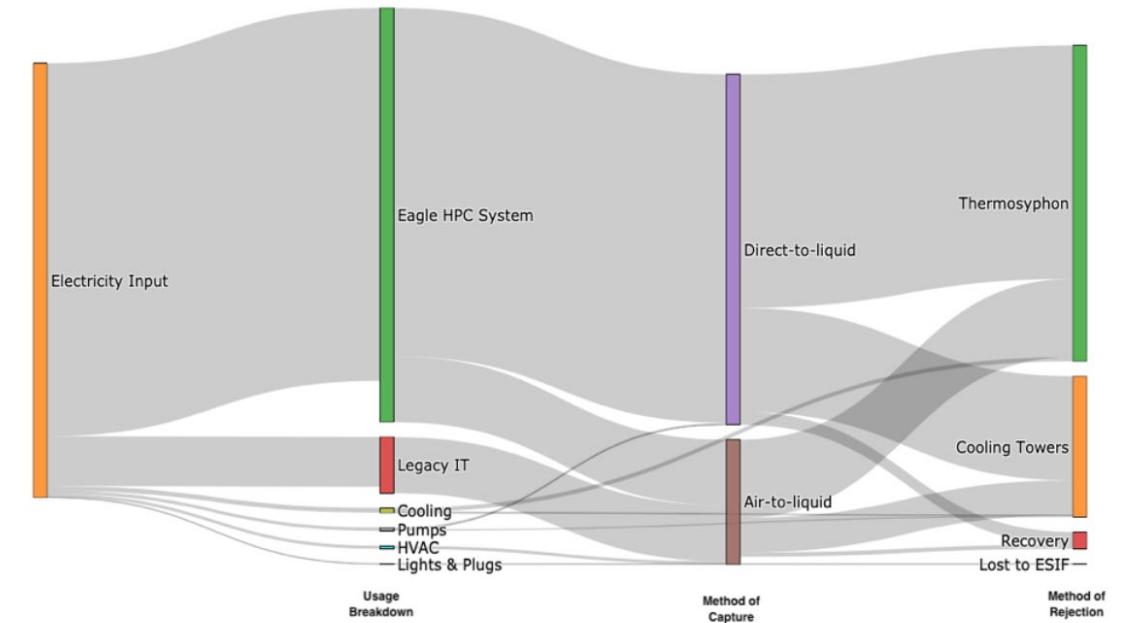
Supermuc-NG (Lenovo SD650 nodes, direct water cooling)



ESIF data center, NREL (PUE of 1.06)



NREL Data Center Energy Balance, From: 2020-11-27T06:00 To: 2020-11-29T23:59



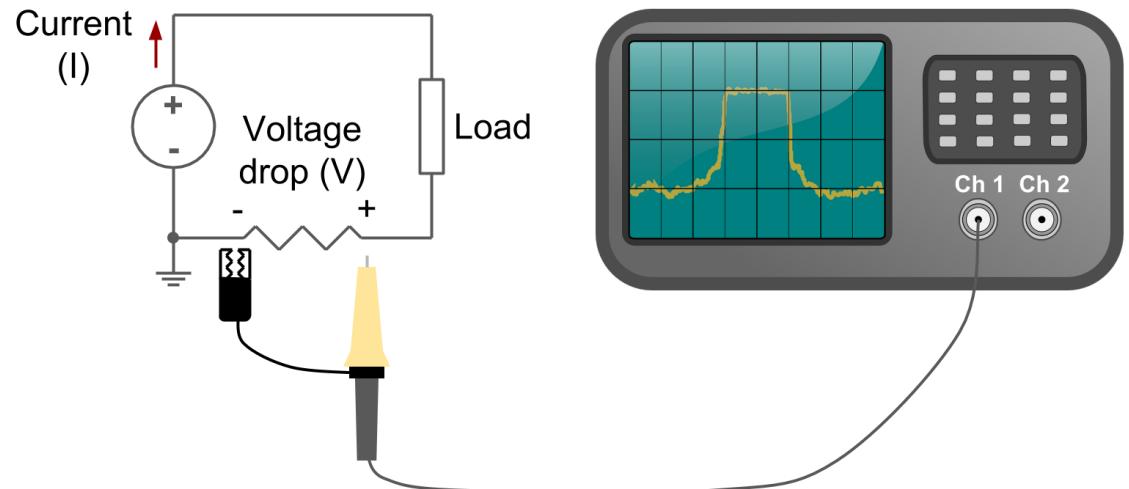


Instrumentation, measurement & modeling



Measurement methods

- ▶ **Hardware instrumentation**
 - ▶ in-bound: data available directly at compute resource (e.g. on- or off-core register)
 - ▶ out-of-bound: data available externally (e.g. via network interface)
 - ▶ Provided by vendor or self-made
 - ▶ Alternatively: wall socket measurements
- ▶ **Models and simulators**
 - ▶ Wattch, SimpleScalar, Sim-PowerCMP, CACTI, McPAT, GPUWATTCH, etc.
- ▶ **Additional considerations: scalability & deployment of instrumentation hardware**



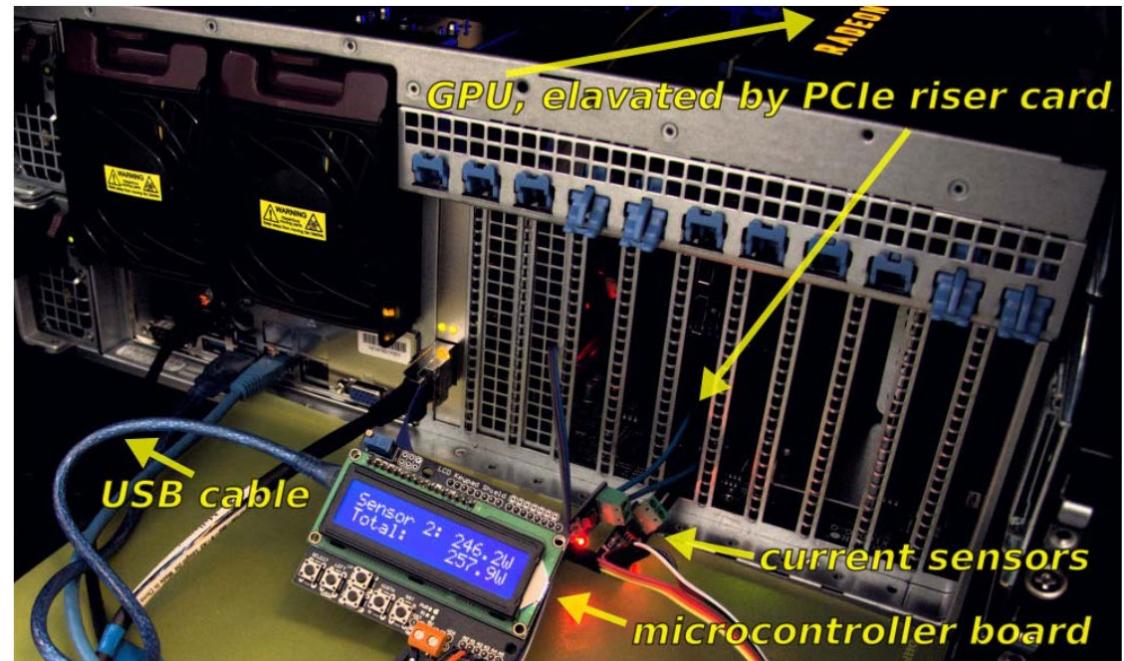
Out-of-bound examples



Voltech PM1000+



PowerMon2



PowerSensor 2

Intel RAPL – Running Average Power Limit

- ▶ In-bound hardware method of monitoring and controlling power and energy
 - ▶ Available since SandyBridge
 - ▶ Model or measurement method depends on CPU microarchitecture
- ▶ Provides low-overhead measurements via model-specific registers (MSRs)
 - ▶ Directly provide energy consumption data in “energy units”, model-dependent
 - ▶ E.g. SandyBridge: 15.3 μ J; Haswell: 61 μ J
 - ▶ Update once every 976 μ s
 - ▶ Increase monotonically (similar to TSC register)
 - ▶ But only 32 bit! Overflows after some time, dependent on CPU stress
- ▶ Also supports controlling the hardware by power capping
 - ▶ E.g. “never exceed 43 watts”
 - ▶ Much more fine-grained than any OS/software mechanism

How to read Intel RAPL data

- ▶ Often requires root due to security implications
 - ▶ Reading MSRs requires raw register access to basically the entire CPU
 - ▶ Reading fine-grained energy/power data might enable side channel attacks
- ▶ Using the Linux kernel's perf_event interface
 - ▶ Using perf and sudo
 - ▶ `sudo perf stat -a -e "power/energy-cores/" /bin/ls`
 - ▶ Alternative without root: requires a `/proc/sys/kernel/perf_event_paranoid` setting of less than 1
- ▶ Manually reading the MSRs (e.g. <https://github.com/kentcz/rapl-tools>)
 - ▶ `sudo modprobe msr && sudo chmod o+rwx /dev/cpu/0/msr`
 - ▶ `sudo setcap cap_sys_rawio+ep <measurement_program>`

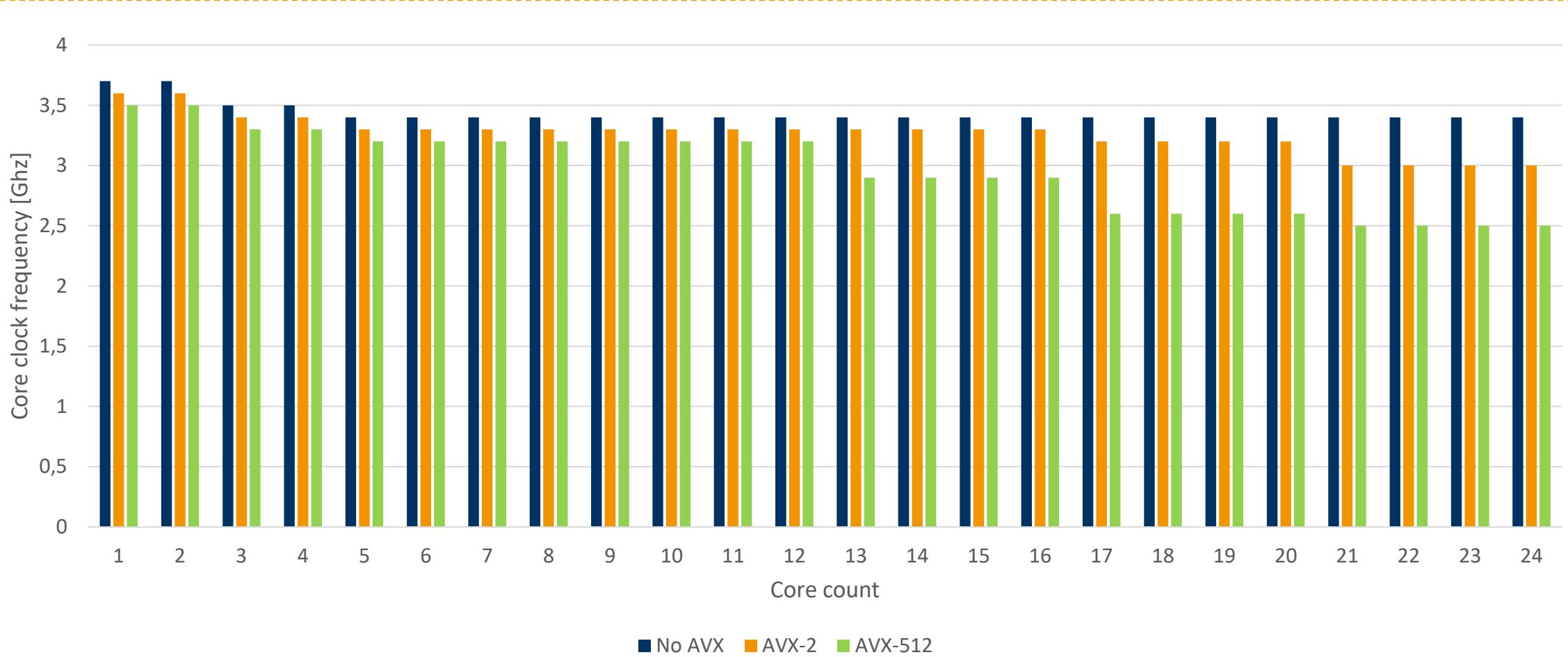
Other vendors

- ▶ A lot of “RAPL”-compatible systems from other vendors
 - ▶ Doesn’t really mean “compatible”, RAPL is used as a deonym here
 - ▶ Just means there is a user-accessible interface for getting power or energy data
- ▶ In-bound: AMD CPUs (RAPL/APM), Apple M1 (?), etc.
- ▶ Out-of-bound: Nvidia (?), AMD GPUs (?), IBM (Amester), Sun (?), etc
- ▶ Third-party: PowerMon, PowerPack, PowerInsight, etc.
 - ▶ Always out-of-bound

DVFS – Dynamic Frequency and Voltage Scaling

- ▶ Hardware can operate at multiple clock frequency points
 - ▶ Usually requires to scale voltage along with frequency
 - ▶ Originally one frequency for all components, these days individual frequencies for cores, last-level cache, etc.
- ▶ Originally controlled in software by selecting next DVFS state – 30ms latency
 - ▶ OS directly requests a so-called P-state by writing into a register
 - ▶ PCU halts CPU e.g. every few milliseconds, reads the current P-state, acts accordingly
- ▶ Modern CPUs have much more autonomy in this
 - ▶ OS requests a certain range of frequencies, minimum QoS or maximum performance
 - ▶ CPU controls the actual setting in hardware (e.g. Intel SpeedShift, ~1 ms latency)
 - ▶ Required by AVX, Turbo and alike to work properly

AVX Turbo Frequencies (Intel Xeon 8174)

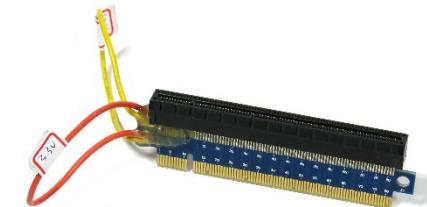


Green500 measurement methodology

- ▶ 33 pages of definitions: measurement devices, topology, workload requirements, averaging, etc.
- ▶ **Level 1 requires to measure**
 - ▶ The entire “core” phase ≥ 1 minute, compute-nodes + measure or estimate network interconnect
 - ▶ Power and take the average
 - ▶ At least $\text{std}::\max(\{2 \text{ kW}, 10\% \text{ of the system}, 15 \text{ nodes}\})$
- ▶ **Level 2**
 - ▶ Level 1 + average power of full run, intermediate measurements (at least 10 averages in core phase)
 - ▶ Compute-node subsystem + measure or estimate all other subsystems
 - ▶ At least $\text{std}::\max(\{10 \text{ kW}, 12\% \text{ of the system}, 15 \text{ nodes}\})$
- ▶ **Level 3**
 - ▶ Level 2 but measure energy and compute average power consumption
 - ▶ Energy measurement resolution: 120 Hz for DC, 5 KHz for AC
 - ▶ Entire system (all components, all nodes, no extrapolations!)

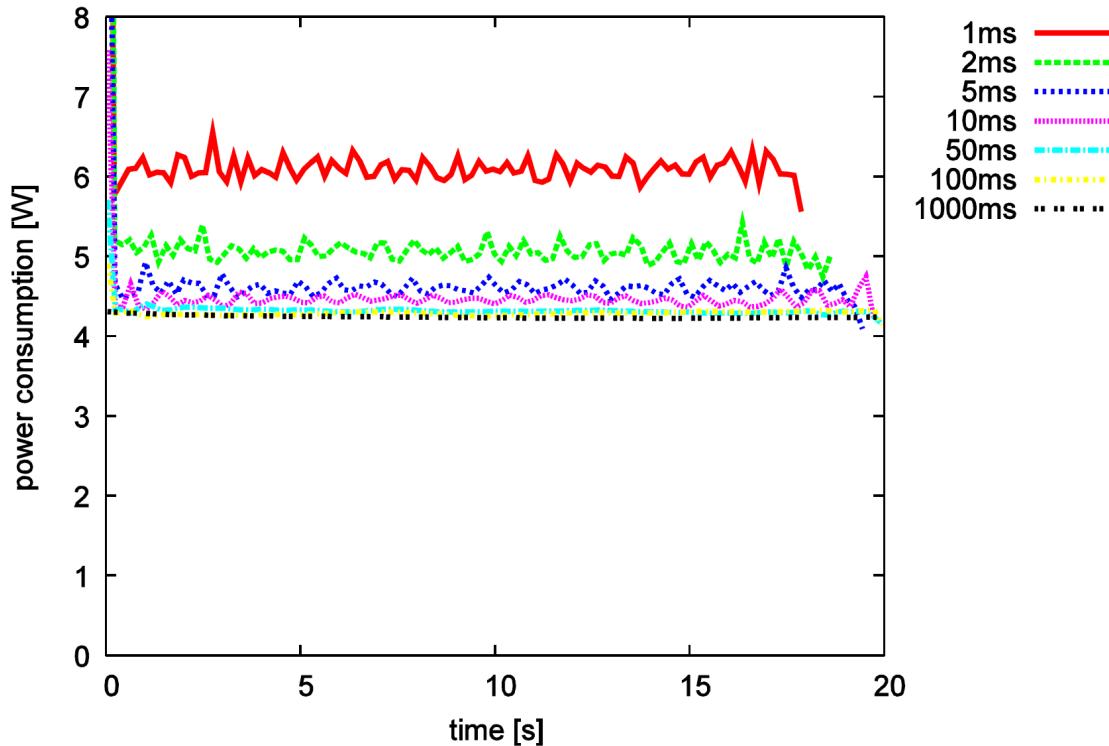
Issues

- ▶ Spatial resolution / topology
 - ▶ E.g. Intel RAPL offers separate readings for entire package, cores and off-core entities (e.g. RAM controller or iGPU)
 - ▶ What about RAM? Mainboard? GPUs? Storage? Network?
 - ▶ GPUs without in-bound measurements: up to 75 watts via PCIe, rest via ATX power connectors, requires PCIe riser cards
- ▶ Temporal resolution and accuracy
 - ▶ High-frequency CPU loads are hardly visible at the wall socket
 - ▶ Resolution != accuracy
 - ▶ E.g. RAPL can have 15.3 μJ resolution but a mJ accuracy or worse
- ▶ External conditions
 - ▶ Temperature
 - ▶ Measurement perturbation (often caused by high-frequency in-bound measurements)

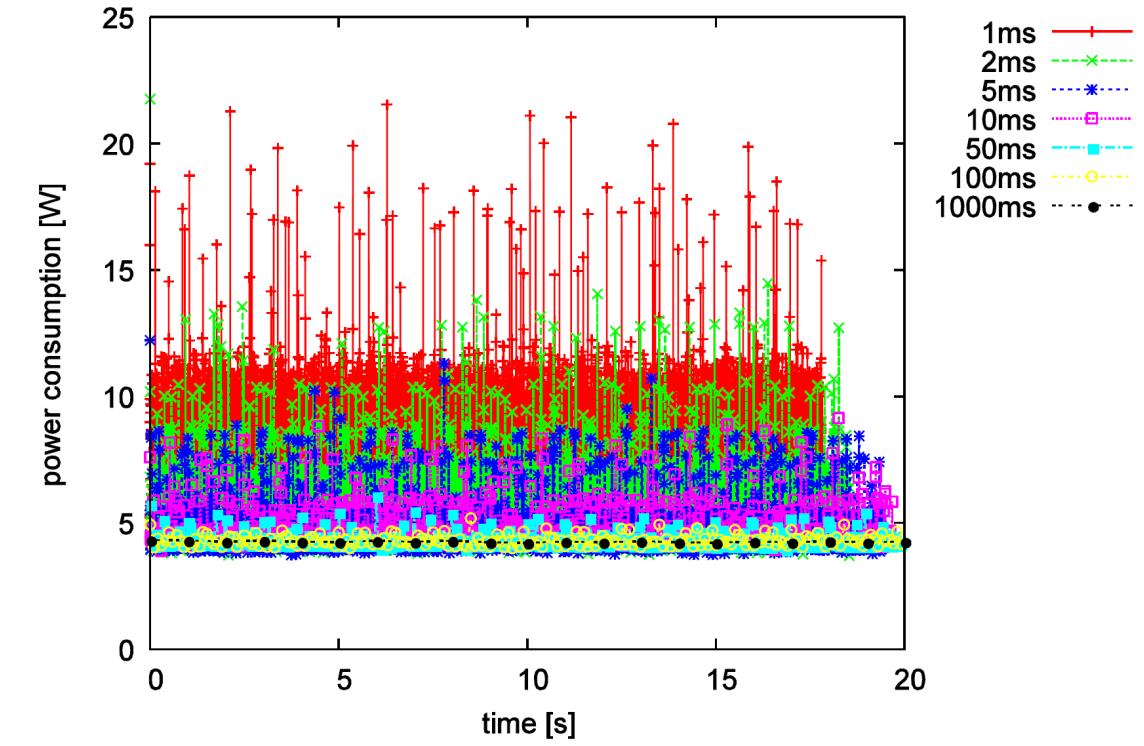


In-bound measurement perturbation on a i7-2600k (RAPL)

Power for SandyBridge idle, accuracy vs. time resolution, smoothed

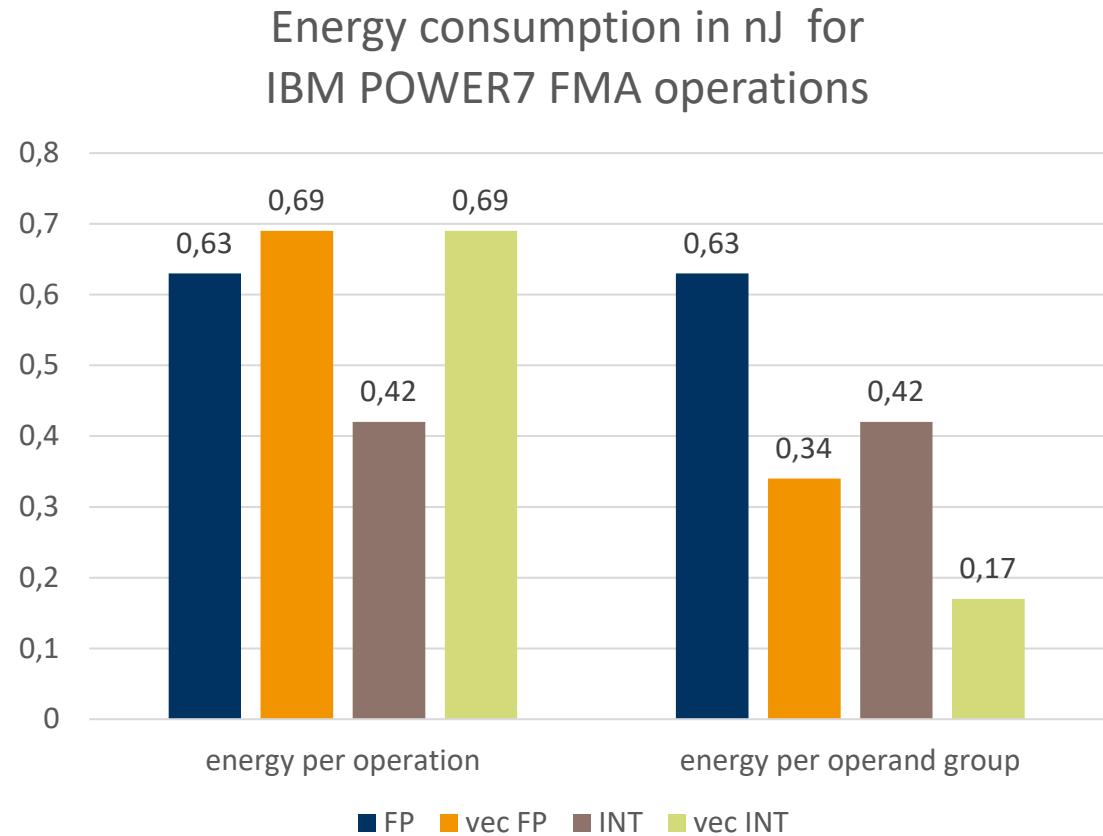


Power for SandyBridge idle, accuracy vs. time resolution



Modeling Example

- ▶ Lots of models available
 - ▶ Some as simple as linear combination of all performance counters
 - ▶ $E_{dynamic} = \sum \alpha_i c_i$
 - ▶ All the way up to support vector machines (SVEs) or Deep Learning





Consequences and applications



Developments in cooling technologies

- ▶ “Hot” water cooling
 - ▶ around 40-50° C inlet temperature
 - ▶ free air cooling (no active chilling)

- ▶ Use waste energy to heat office buildings
 - ▶ Can even use waste heat as energy to operate active cooling (e.g. SuperMUC-NG)

- ▶ Immersive cooling

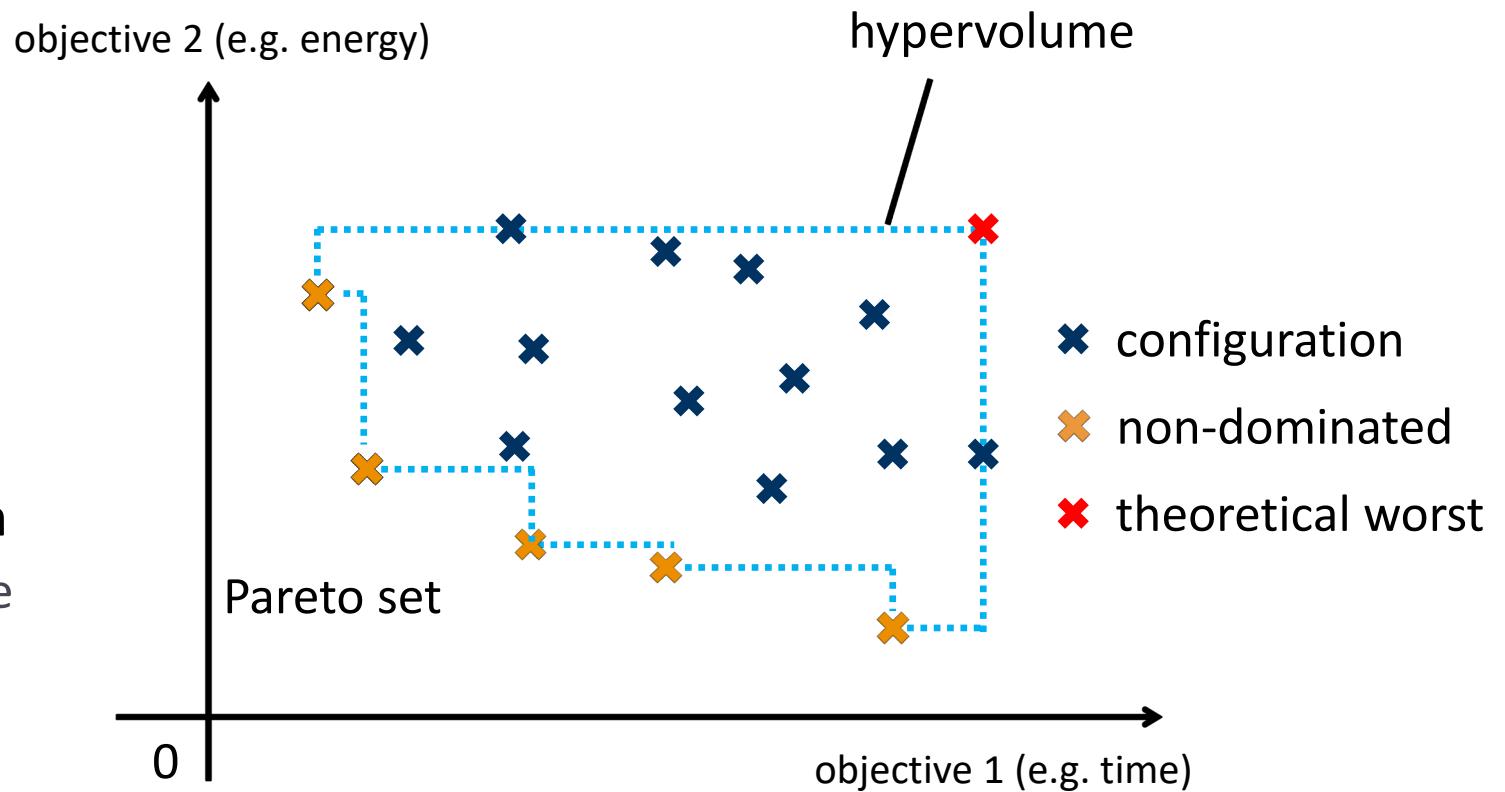


Multi-objective optimization

- ▶ We now have at least two objectives which are (partially) in conflict
 - ▶ Faster program execution (requires more power and energy)
 - ▶ Lower power or energy consumption (requires longer execution times)
- ▶ Two possible optimization approaches
 - ▶ Weight function to combine all objectives into a single one
 - ▶ E.g. EDP – Energy Delay Product, also: ED^2P , ED^3P , PDP, ...
 - ▶ Alternatively use true multi-objective optimization, keeps flexibility
 - ▶ e.g. Pareto optimality

Multi-objective optimization and Pareto optimality

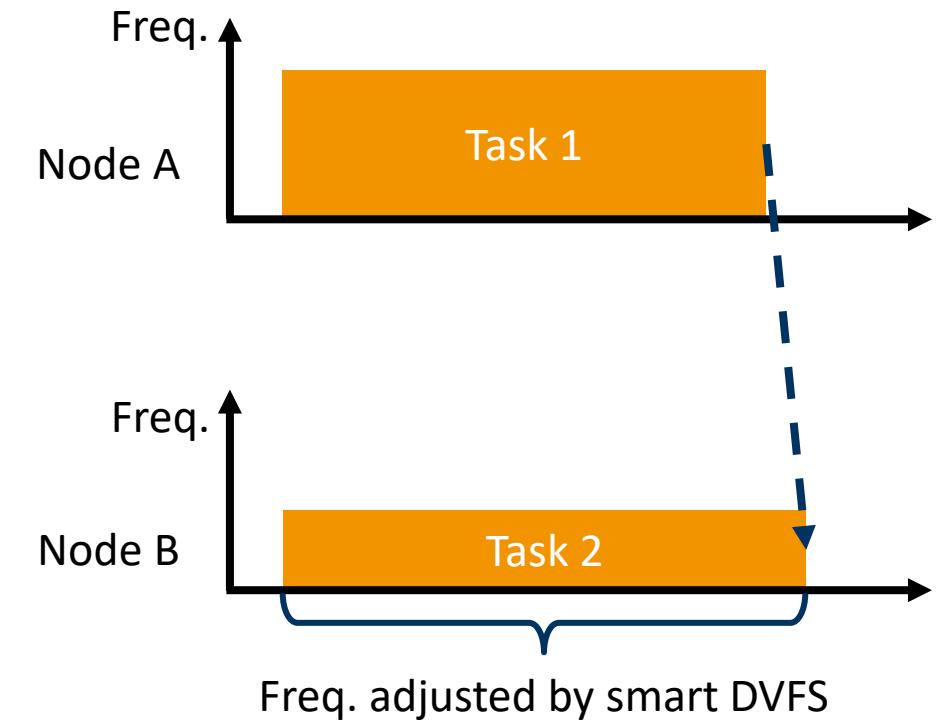
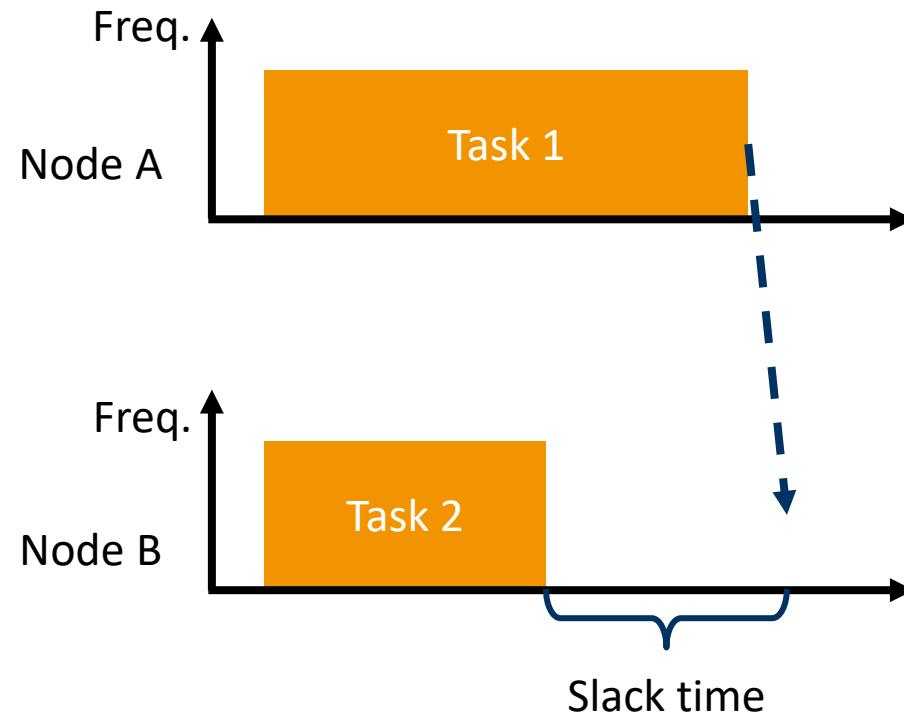
- ▶ Definition of Pareto set:
 - ▶ For each point in the set, no other point “dominates” it (=is better in all objectives)
 - ▶ Entails that it’s impossible to improve one objective without worsening another objective
- ▶ Superior to weighted approach
 - ▶ Provide all configurations in the Pareto set to the user, choose dynamically according to current preferences



MPI slack time optimization

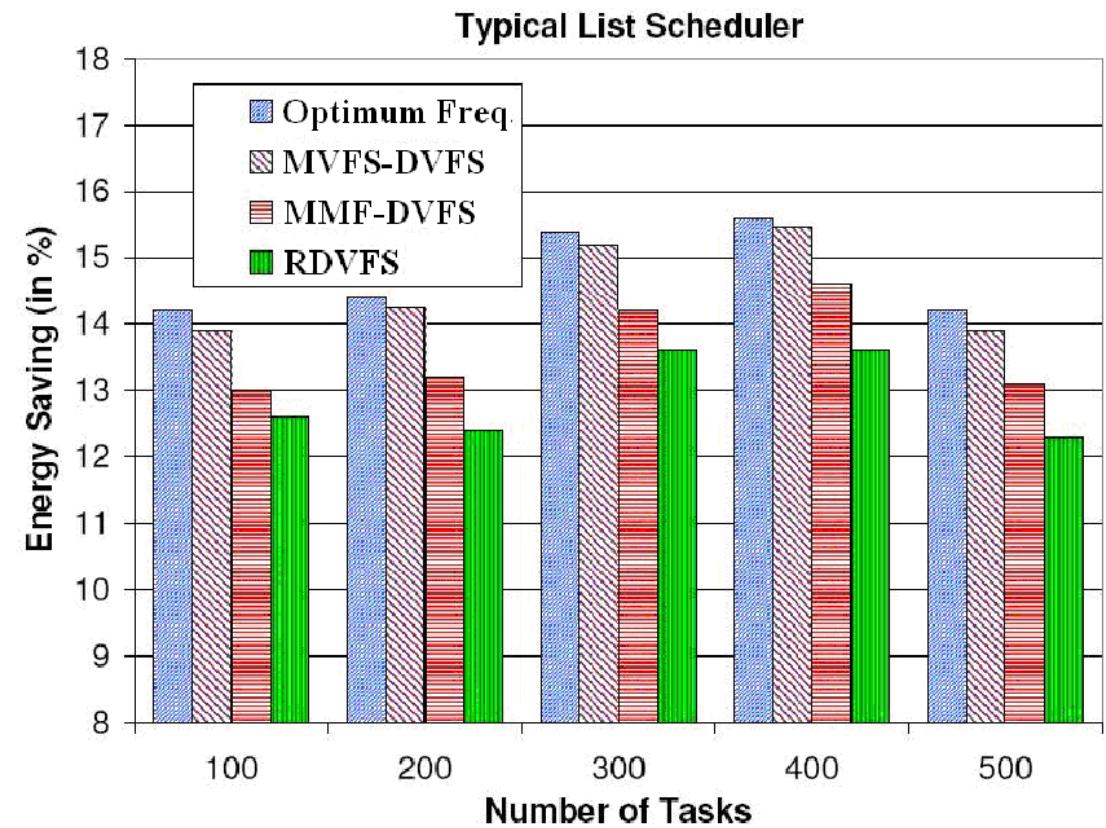
- ▶ Recognize slack time in parallel applications
 - ▶ Wait states
 - ▶ Periods of extended memcpy operations or I/O
 - ▶ Even computation if not on the critical path
 - ▶ etc.
- ▶ Use DVFS to reduce energy footprint with minimal impact on wall time
 - ▶ Lots of work on that from 5-15 years ago

Slack time optimization example



Slack time optimization results

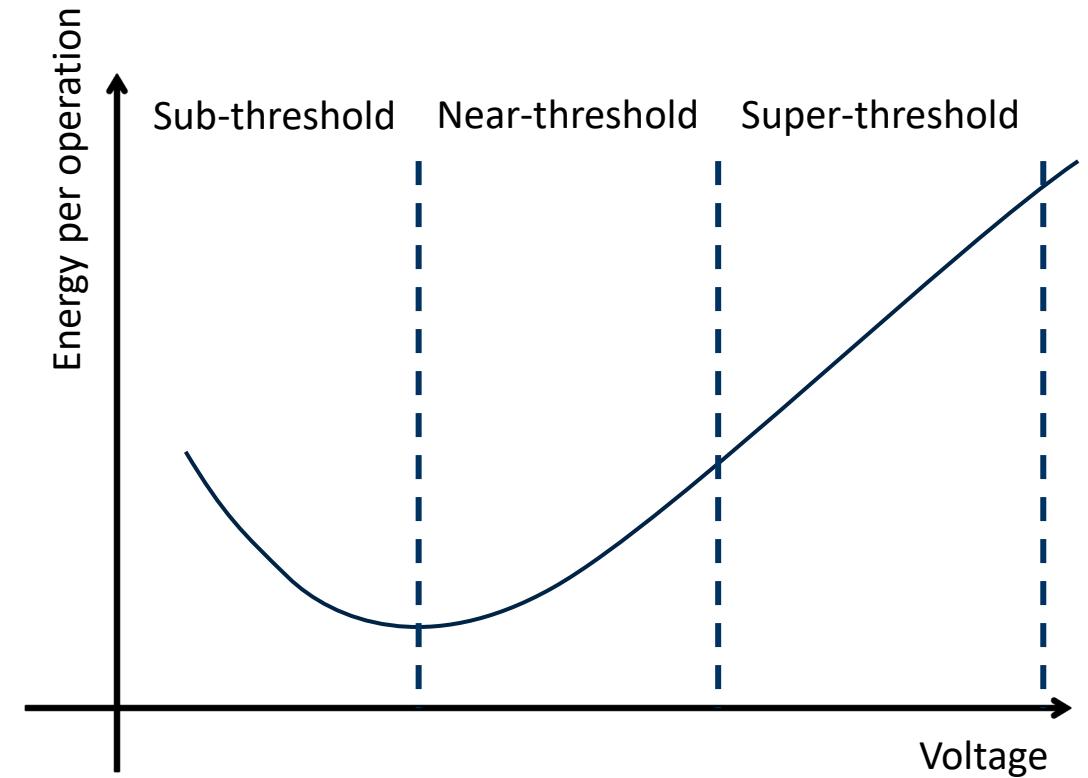
- ▶ Rizvandi et al. „Some Observations on Optimal Frequency Selection in DVFS-based Energy Consumption Minimization”
 - ▶ Simulations with 3000 randomly generated task graphs
 - ▶ Energy savings 10-20%



<https://arxiv.org/ftp/arxiv/papers/1201/1201.1695.pdf>

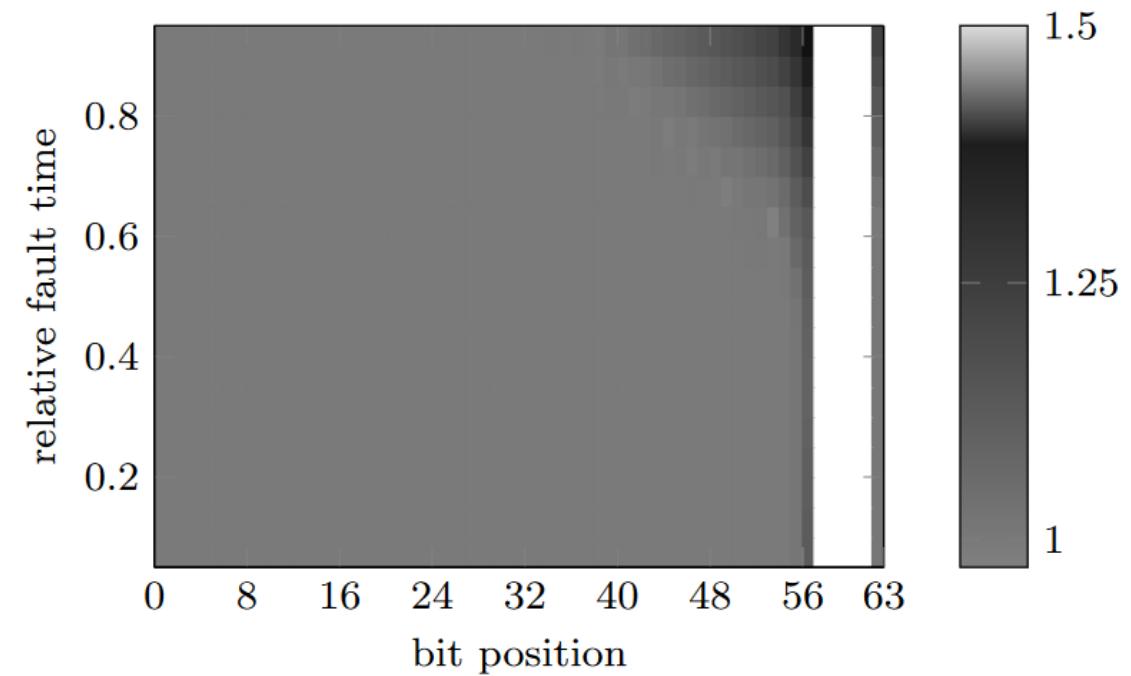
Near-Threshold Voltage

- ▶ Idea: reduce voltage below safe operating levels
 - ▶ Near-threshold voltage computing (NTV/NTC)
 - ▶ Can still operate transistors, but at large clock frequency reductions
 - ▶ Slowdown of 5x-10x



Approximate Computing

- ▶ Idea: use NTV and mitigate slowdown with parallelism
 - ▶ Single, reliable core at super-threshold voltage
 - ▶ Several unreliable ones at near-threshold voltage (under the same power envelope!)
 - ▶ Switch cores depending on state of computation
- ▶ Investigate effects of bit flips in floating point data for converging algorithms (e.g. jacobi)



Approximate Computing cont'd

- ▶ Idea: reduce voltage below safe operating levels
 - ▶ Near-threshold voltage computing (NTV/NTC)
 - ▶ Replace single, reliable core with several unreliable ones under the same power envelope
 - ▶ Speedup through parallelism results in energy savings
- ▶ Investigate effects of bit flips in floating point data for converging algorithms (e.g. jacobi)

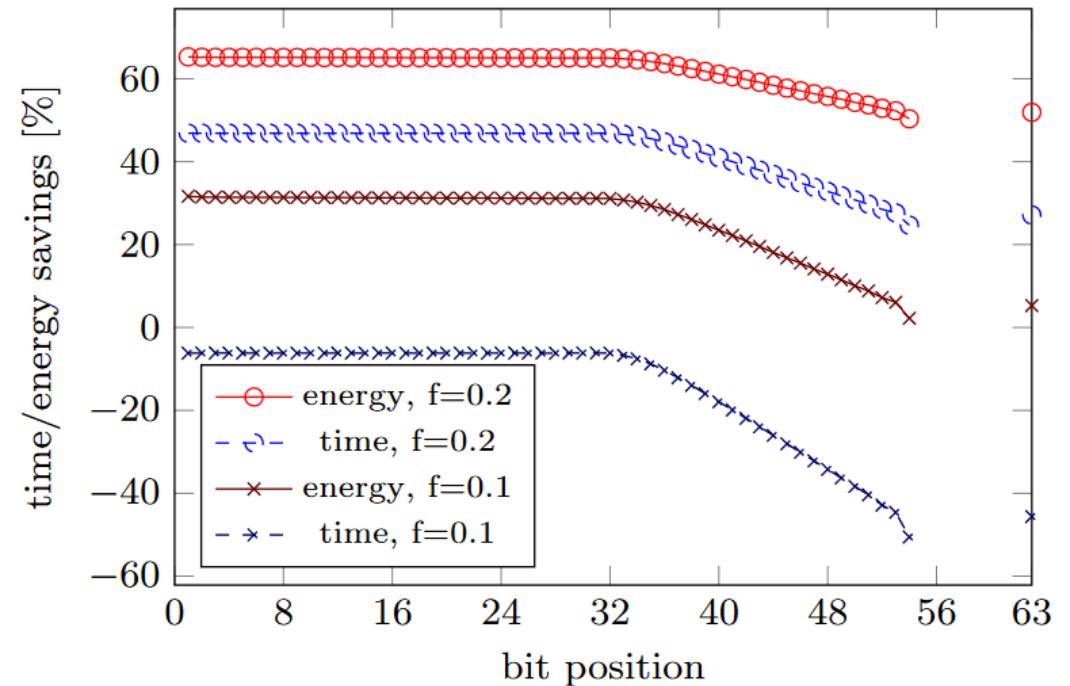


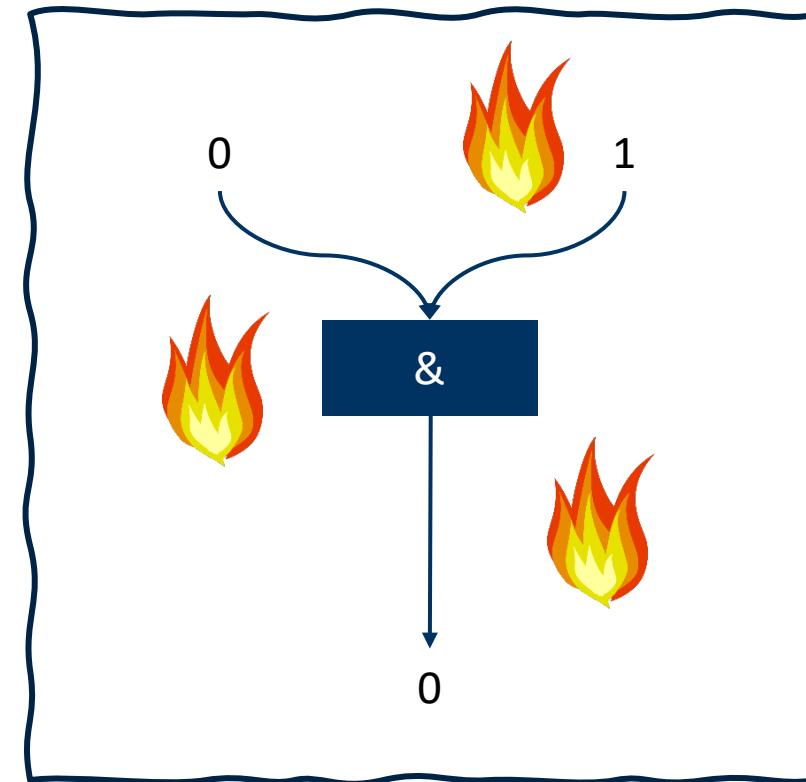
Fig. 4: Relative energy and time savings of an unreliable, parallel run of Jacobi on 16 cores compared to a reliable, sequential one. The missing data at bits 55–62 denotes divergence.

Additional Consequences

- ▶ Job Scheduling: Profile applications, compute roofline model, set optimal DVFS setting for consecutive runs
 - ▶ E.g. EAR – Energy Aware Runtime (SLURM tool @ SuperMUC, LRZ, Germany)
- ▶ Influences load balancing & scheduling decisions
 - ▶ Data movement is expensive
 - ▶ Move tasks to location of data instead of data to location of tasks
- ▶ Stacked memory, processing-in-memory, etc.

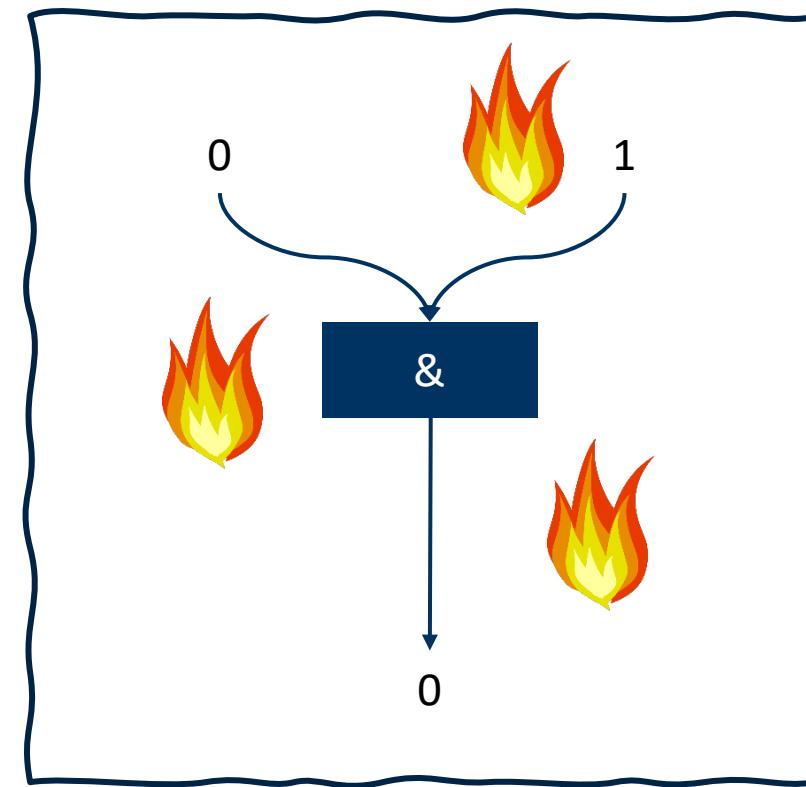
Reversible computing and Landauer principle: the future?

- ▶ There's a lower theoretical limit ("Landauer limit") to energy consumption of computation
 - ▶ Irreversible computation (e.g. logical AND) erases information, hence must be accompanied by corresponding entropy increase (=heat) in a closed system
 - ▶ because thermodynamics $\neg \backslash (\tau) / \neg$
 - ▶ Landauer limit is approx. 0.0175 eV or $2.805 * 10^{-21}$ J at room temperature
 - ▶ We're currently still several orders of magnitude away from that...



Reversible computing and Landauer principle: the future? cont'd

- ▶ Koomey's Law: The number of computations per joule doubles every 1.57 years
 - ▶ Coupled with Landauer limit: no more energy efficiency increase after 2080...
 - ▶ Also applies to quantum computing
- ▶ Solution: reversible computing
 - ▶ In theory, computing without losing information doesn't need to increase entropy, hence no heat



Summary

- ▶ Energy is a hot topic in HPC
- ▶ Instrumentation and measurements got a lot easier over the past decade
 - ▶ Availability not an issue these days
 - ▶ Accuracy and granularity however is
- ▶ Discussed several research perspectives on the topic
 - ▶ Multi-objective optimization, slack time optimization, approximate computing, ...

Image Sources

- ▶ Voltech PM1000: <https://www.voltech.com/media/8d8595acca6d01f/pm1000-user-manual-v14.pdf>
- ▶ Shunt transistor: <https://articles.saleae.com/oscilloscopes/how-to-measure-current-with-an-oscilloscope>
- ▶ PowerMon:
https://ieeexplore.ieee.org/abstract/document/5453824?casa_token=Ax4_mcUUFOIAAAAAA:YFY5X2H6aCU2pMDs5gwvMqbvA28huJePfLkRDveibf6d1TKkKmvqXfgCVtyVz1nZCp_z8-mIT9U
- ▶ PowerSensor 2:
https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8366941&casa_token=T70QZiDS9F4AAAAAA:O02oYwOTJXaLIRj8as2ZAGQSmYUDeigrjd5Mt-sAIGfegoz0NIH25rfXL1gsEM8mmM6WYtr6DdE
- ▶ PCIe riser: <https://www.igorslab.de/en/power-recording-graphics-card-power-supply-interaction-measurement/3/>
- ▶ Supermuc-NGn nodes: <https://www.lenovo.com/us/en/p/servers-storage/servers/high-density/thinksystem-sd650-n-v2/77xx7dsd672?orgRef=https%253A%252F%252Fwww.google.com%252F>
- ▶ Intel Architecture Day Slide: <https://download.intel.com/newsroom/2021/client-computing/intel-architecture-day-2021-presentation.pdf>
- ▶ Alder Lake Die Shots: https://www.reddit.com/r/intel/comments/qhbbow/10nm_esf_intel_7_alder_lake_die_shot/

Sources

- ▶ Green500 measurement methodology:
<https://www.top500.org/static/media/uploads/methodology-2.0rc1.pdf>

Consequences

- ▶ AVX instructions running at lower clock rate
 - ▶ show in some experiment on modern hardware?

Multi-objective optimization

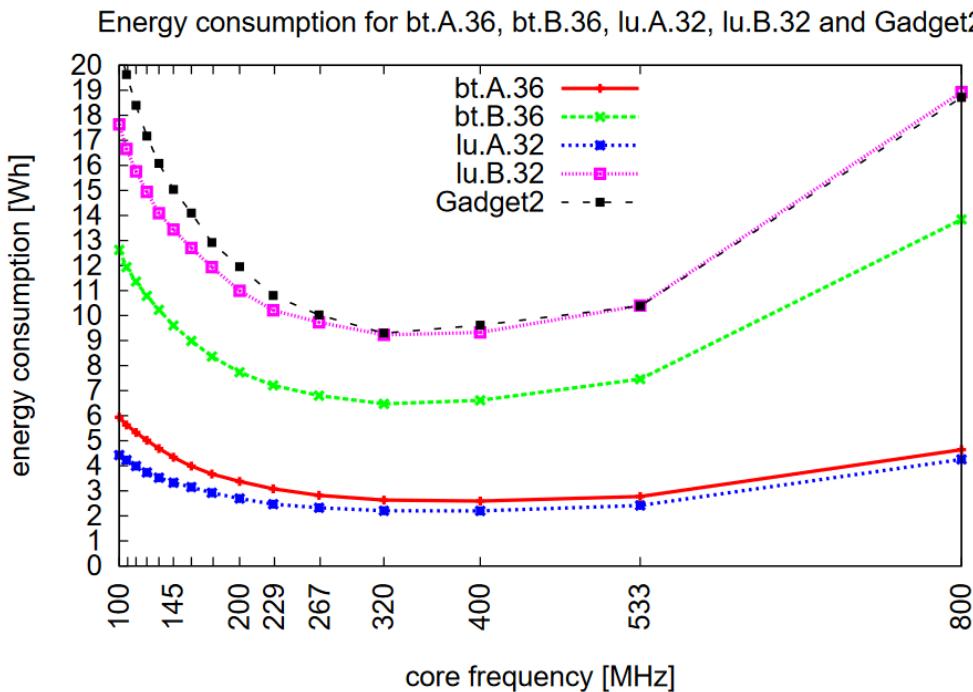


Fig. 6. Energy consumption of bt.A.36, bt.B.36, lu.A.32, lu.B.32 and Gadget2 for all possible core clock frequencies.

big.LITTLE

Thread-core Mappings (Beyond HPC?)

P core
(Golden Cove)
up to 5 GHz,
HT, private L2

E core
(Gracemont)
up to 4 GHz, no
HT, shared L2



Beginning Fall 2021

Alder Lake

Reinventing Multi Core Architecture

Single, Scalable SoC Architecture

All Client Segments – 9W to 125W – built on Intel 7 process

All-New Core Design

Performance Hybrid with Intel Thread Director

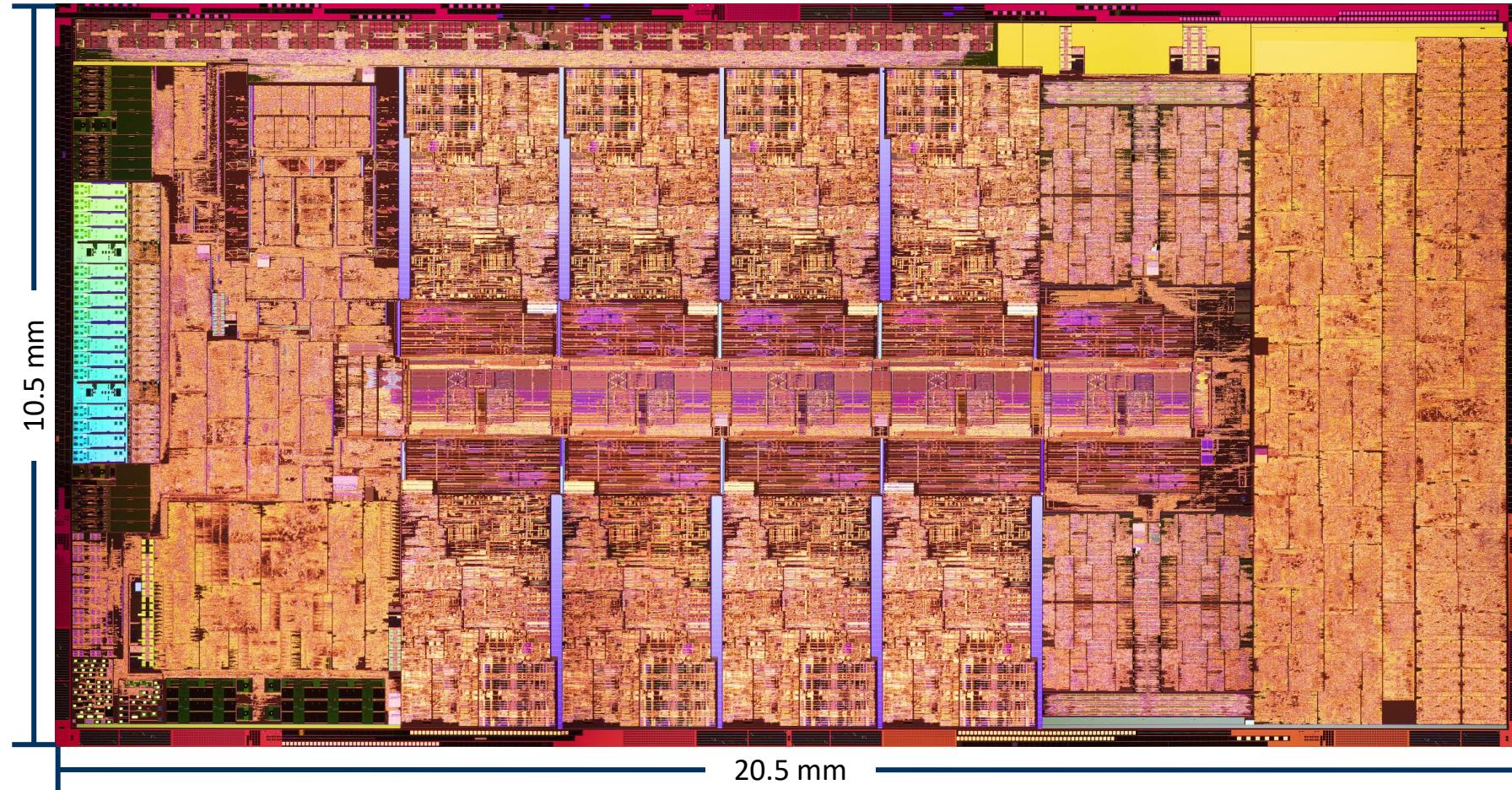
Industry-Leading Memory & I/O

DDR5, PCIe Gen5, Thunderbolt™ 4, Wi-Fi 6E

intel.

79

Alder Lake Die Shot



Alder Lake Die Shot Annotated



Alder Lake Take-Aways

- ▶ Affinity will gain significance
 - ▶ fast P-cores and efficient E-cores on the same chip
(note that ARM has been doing this for years with big.LITTLE)
 - ▶ OS scheduler and programs need to be aware
 - ▶ e.g. previous images show an 8P+8E configuration
- ▶ L3 cache layout leads to on-chip NUMA
 - ▶ all L3 cache is accessible to every core but not with the same performance
 - ▶ has been the case at least since Haswell (~2014)
- ▶ Vectorization units take up a lot of transistor space