# 703365 Sustainability in Computer Science:
# Green HPC: Paving the Way for Sustainable Supercomputing

Philipp Gschwandtner, Research Center HPC, 13.01.2025

# What is Green HPC?
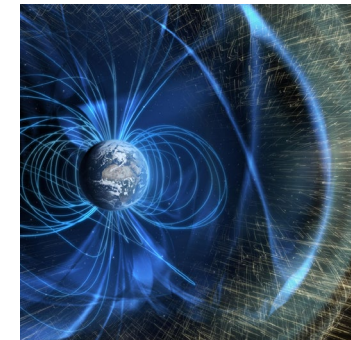
"environmentally friendly"

## High Performance Computing

"[…] uses supercomputers and computer clusters to solve advanced computation problems"
(Wikipedia)

# Green HPC

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner
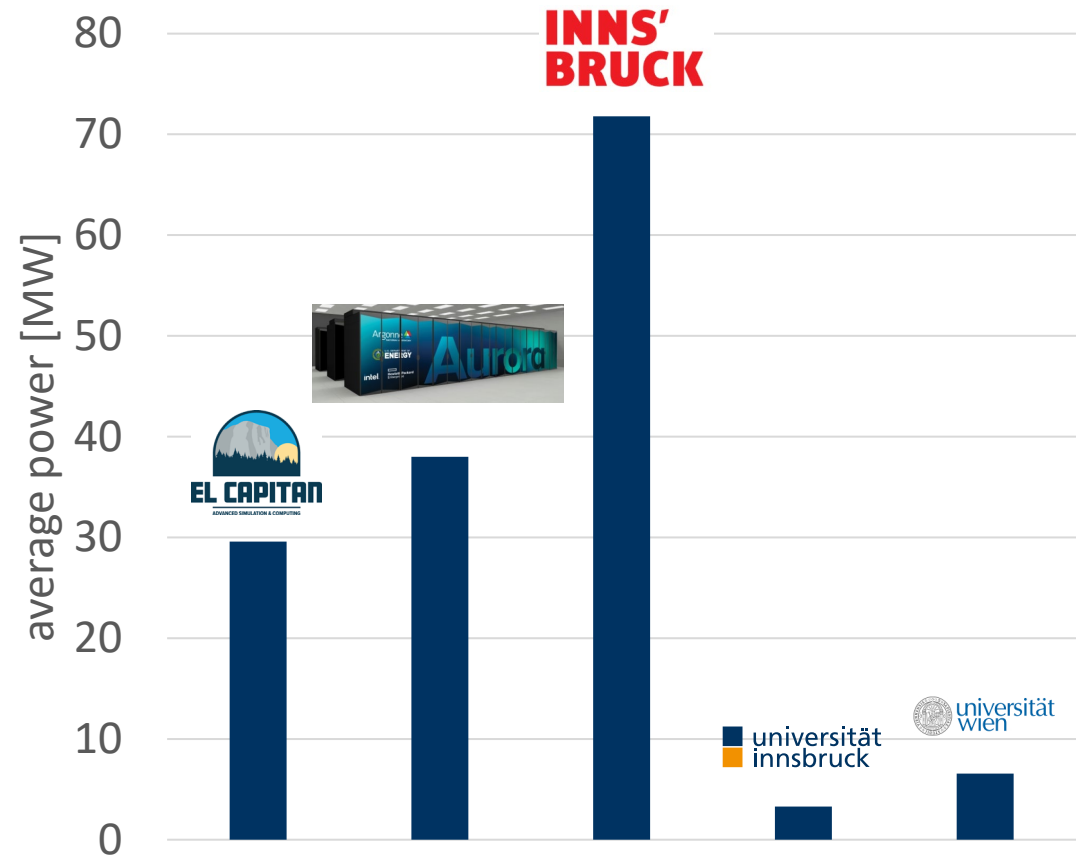
# Main sustainability concern in HPC is energy – but why?

- Example: El Capitan
  - currently the fastest supercomputer world-wide, located in Livermore, California, USA



- Extreme computing and storage capacities
  - 11 million cores (~44k AMD MI300A CPUs + GPUs)
  - 5,4 petabytes of main memory (RAM)
  - 1,3 exabytes of storage (planned)
  - measured peak performance of ~1,7 exaflops (>$10^{18}$ arithmetic operations per second)

- High electrical and thermal operating requirements
  - 30 megawatts of power
  - the **main limiting design factor** when building supercomputers

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# How much is 30 megawatts?

▶ 100% El Capitan supercomputer

▶ 78% Aurora supercomputer (No. 3 world-wide)

▶ 41% of Innsbruck
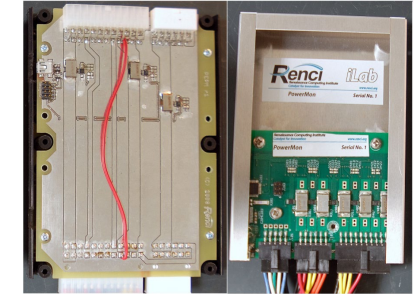
▶ 9x University of Innsbruck

▶ 4,5x University of Vienna



Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# To be able to optimize, we need to measure first!

▸ Need power/energy instrumentation
  ▸ "homemade" examples on the right, good for experimental research but does not scale to large systems (also: fire hazard!)
  ▸ modern supercomputers have this built-in
  ▸ note that there is also a plethora of power/energy models – some better, some worse
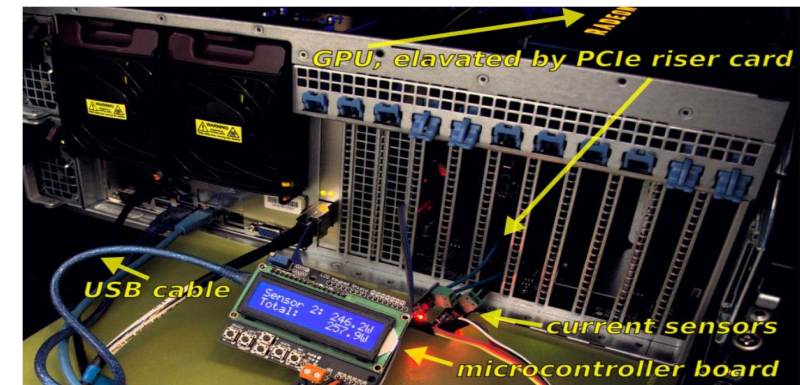


Voltech PM1000+



PowerMon2

▸ Need metrics to be able to evaluate, compare and optimize
  ▸ e.g. Power Usage Effectiveness

$$PUE = \frac{\text{total facility energy}}{\text{IT equipment energy}}$$

  ▸ e.g. Energy-Delay Product (EDP)

$$EDP = \text{energy} \times \text{walltime}$$



PowerSensor 2

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# TOP500 List

▸ **List of the fastest supercomputers world-wide**
  ▸ released twice every year
  ▸ high performance linpack (HPL) benchmarking (linear algebra stress-testing)
  ▸ very interesting analyses and statistics around supercomputing and HPC
  ▸ https://www.top500.org

▸ **Currently November 2024 edition**
  ▸ Power consumption reported for many (but not all!) systems
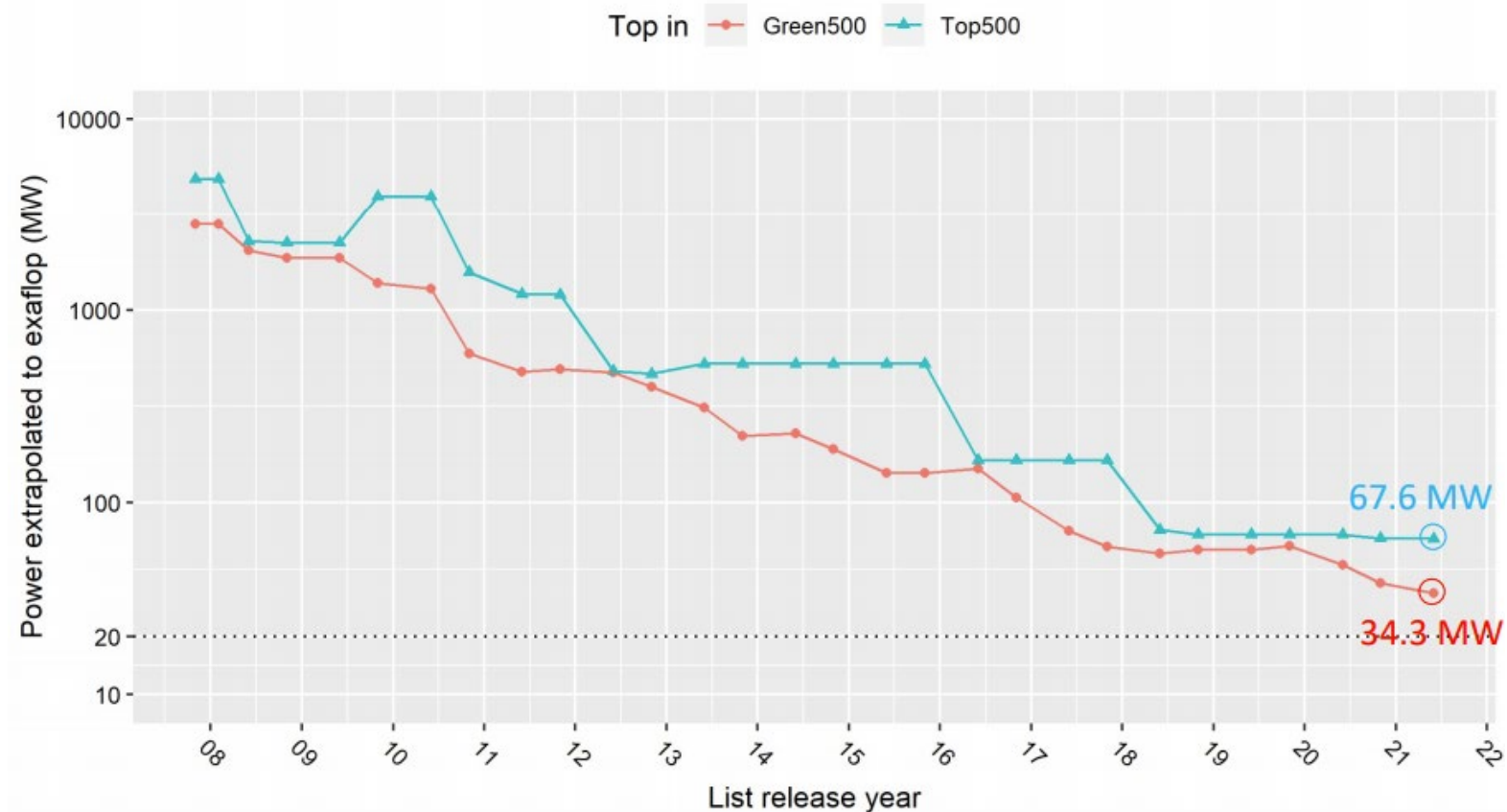
▸ **Also: Green500**
  ▸ Performance-per-energy ranking

| Rank | System | Cores | Rmax (PFlop/s) | Rpeak (PFlop/s) | Power (kW) |
|---|---|---|---|---|---|
| 1 | **El Capitan** - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS, HPE DOE/NNSA/LLNL United States | 11,039,616 | 1,742.00 | 2,746.38 | 29,581 |
| 2 | **Frontier** - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE Cray OS, HPE DOE/SC/Oak Ridge National Laboratory United States | 9,066,176 | 1,353.00 | 2,055.72 | 24,607 |
| 3 | **Aurora** - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel DOE/SC/Argonne National Laboratory United States | 9,264,128 | 1,012.00 | 1,980.01 | 38,698 |
| 4 | **Eagle** - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft Azure Microsoft Azure United States | 2,073,600 | 561.20 | 846.84 | |
| 5 | **HPC6** - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, RHEL 8.9, HPE Eni S.p.A. Italy | 3,143,520 | 477.90 | 606.97 | 8,461 |

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# 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 `std::max({2 kW, 10% of the system, 15 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 `std::max({10 kW, 12% of the system, 15 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!)

| Rank | TOP500 Rank | System | Cores | Rmax (PFlop/s) | Power (kW) | Energy Efficiency (GFlops/watts) |
|---|---|---|---|---|---|---|
| 1 | 222 | **JEDI** - BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200, ParTec/EVIDEN EuroHPC/FZJ Germany | 19,584 | 4.50 | 67 | 72.733 |
| 2 | 122 | **ROMEO-2025** - BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200, Red Hat Enterprise Linux, EVIDEN ROMEO HPC Center - Champagne-Ardenne France | 47,328 | 9.86 | 160 | 70.912 |
| 3 | 440 | **Adastra 2** - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, RHEL, HPE Grand Equipement National de Calcul Intensif - Centre Informatique National de l'Enseignement Suprieur (GENCI-CINES) France | 16,128 | 2.53 | 37 | 69.098 |
| 4 | 155 | **Isambard-AI phase 1** - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE University of Bristol United Kingdom | 34,272 | 7.42 | 117 | 68.835 |

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner
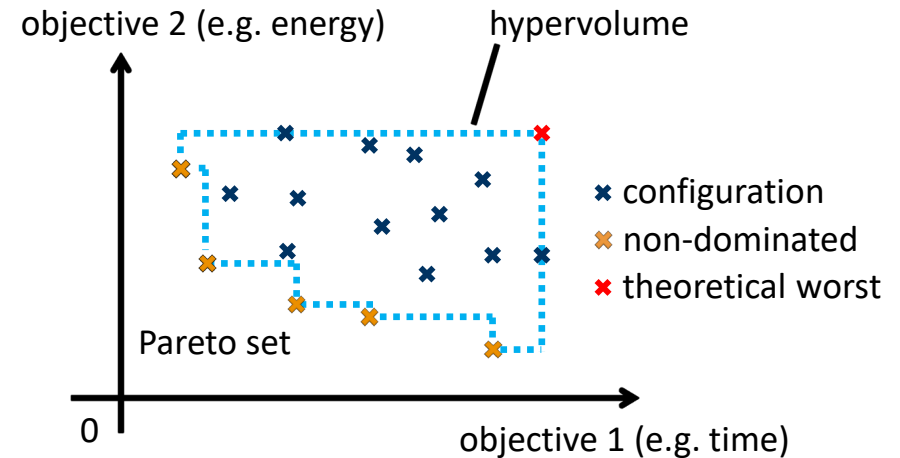
# Power consumption projected to 1 exaflop



https://www.hpcwire.com/2021/07/15/15-years-later-the-green500-continues-its-push-for-energy-efficiency-as-a-first-order-concern-in-hpc/

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner
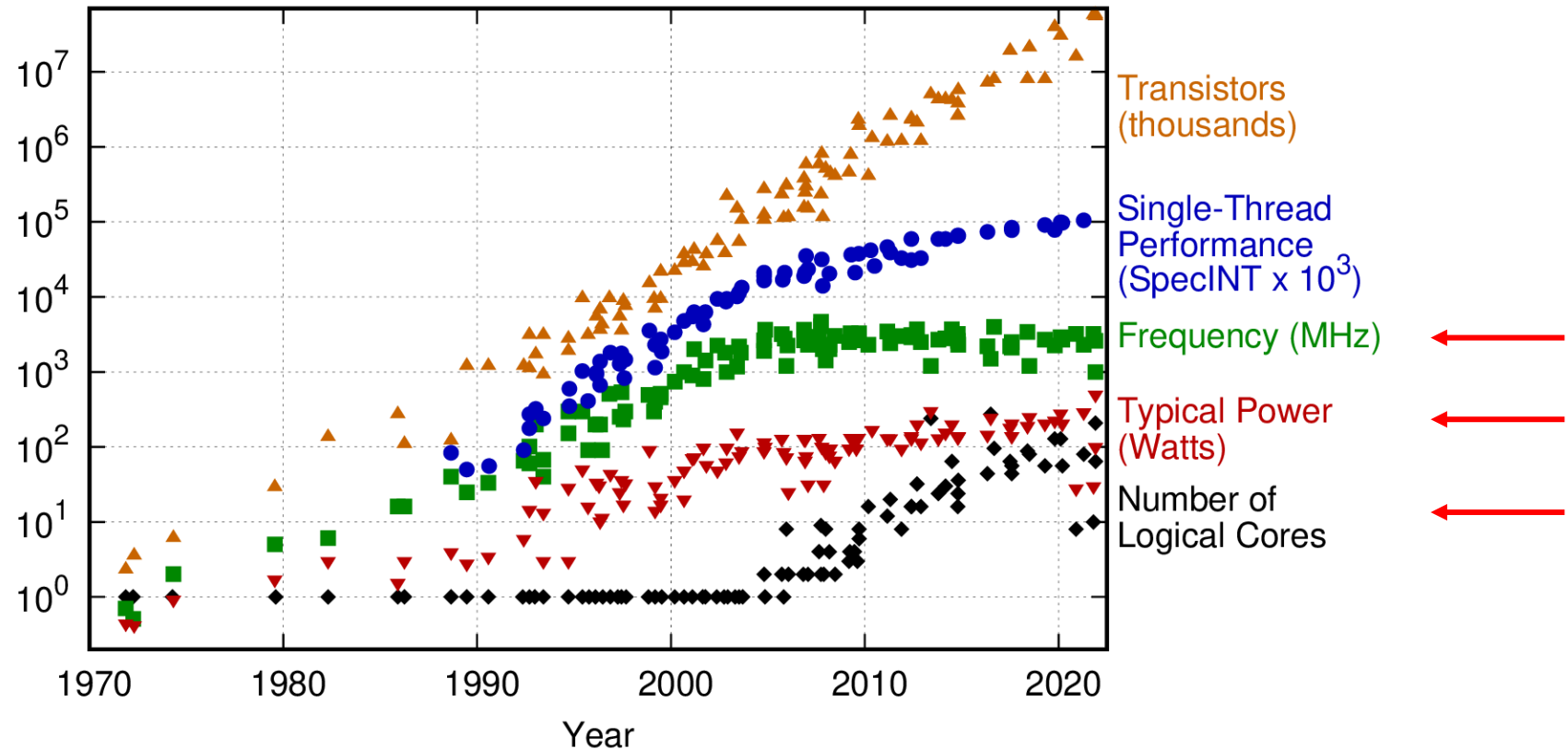
# Avenues of optimization

▶ **Multiple attack points for making HPC more energy-efficient**

  ▶ increased parallelism

  ▶ cooling

  ▶ what and how the hardware is used

▶ **When working with energy-efficient HPC, it's always a multi-objective problem**

  ▶ optimizing for power and/or energy often means sacrificing (a little bit of) performance

  ▶ e.g. Pareto-optimality



objective 2 (e.g. energy)    hypervolume

✖ configuration
✖ non-dominated
✖ theoretical worst

Pareto set

0    objective 1 (e.g. time)

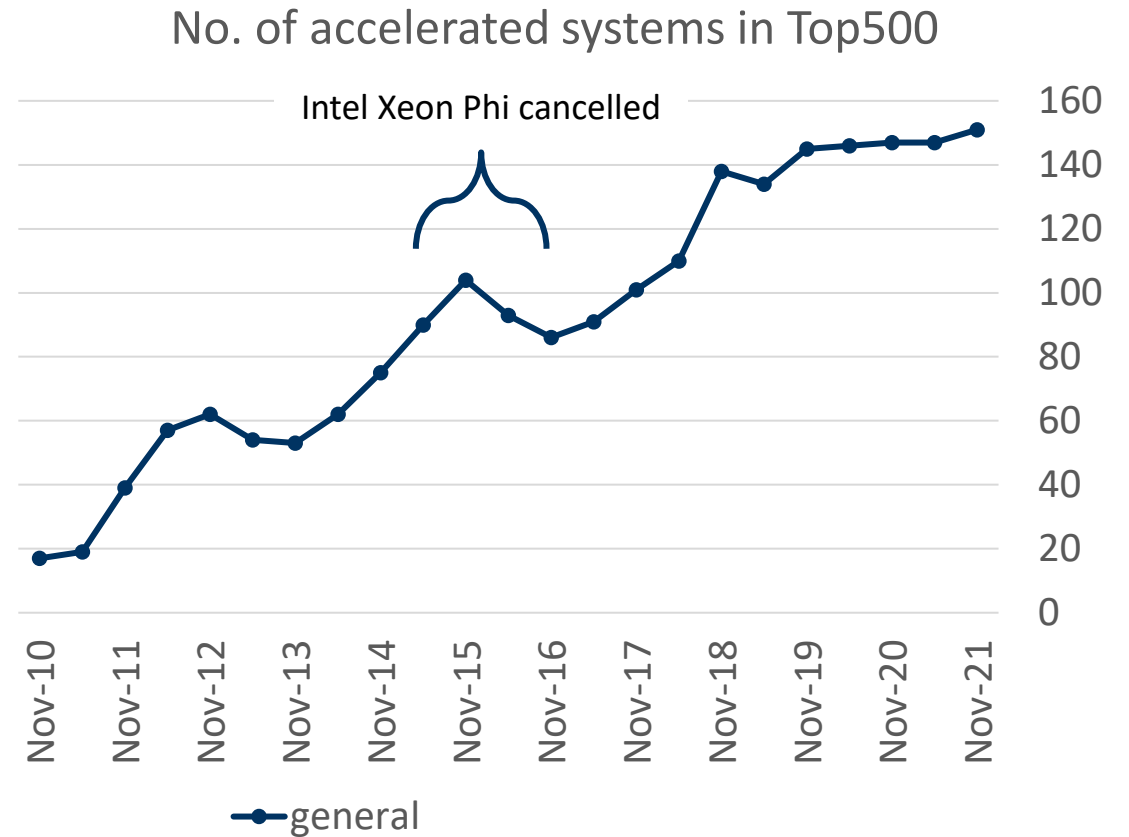# Reducing energy in computing: Parallelism!



50 Years of Microprocessor Trend Data

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
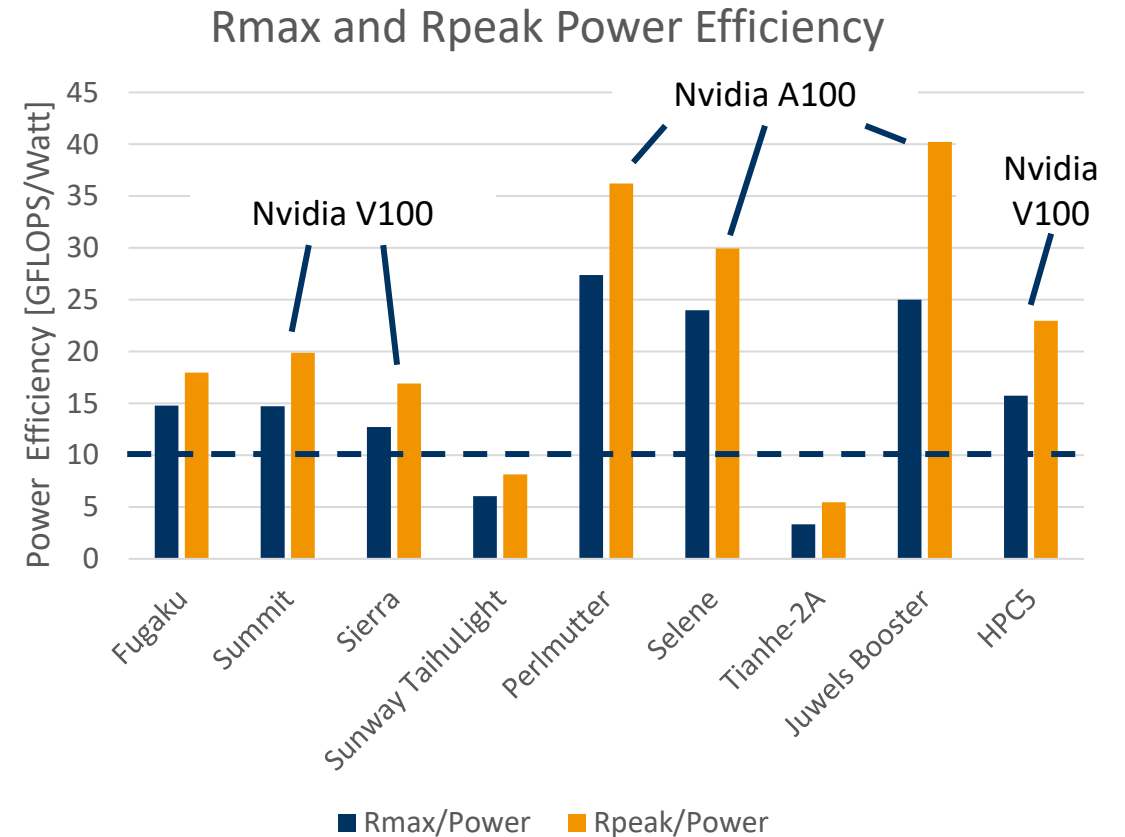New plot and data collected for 2010-2021 by K. Rupp

# Reducing energy in computing: Accelerators!

▶ **Accelerator market share in HPC has been steadily increasing and will likely continue to do so**

  ▶ Why? Distributed memory clusters with accelerators provide some of the best cost- and energy-efficiency in HPC

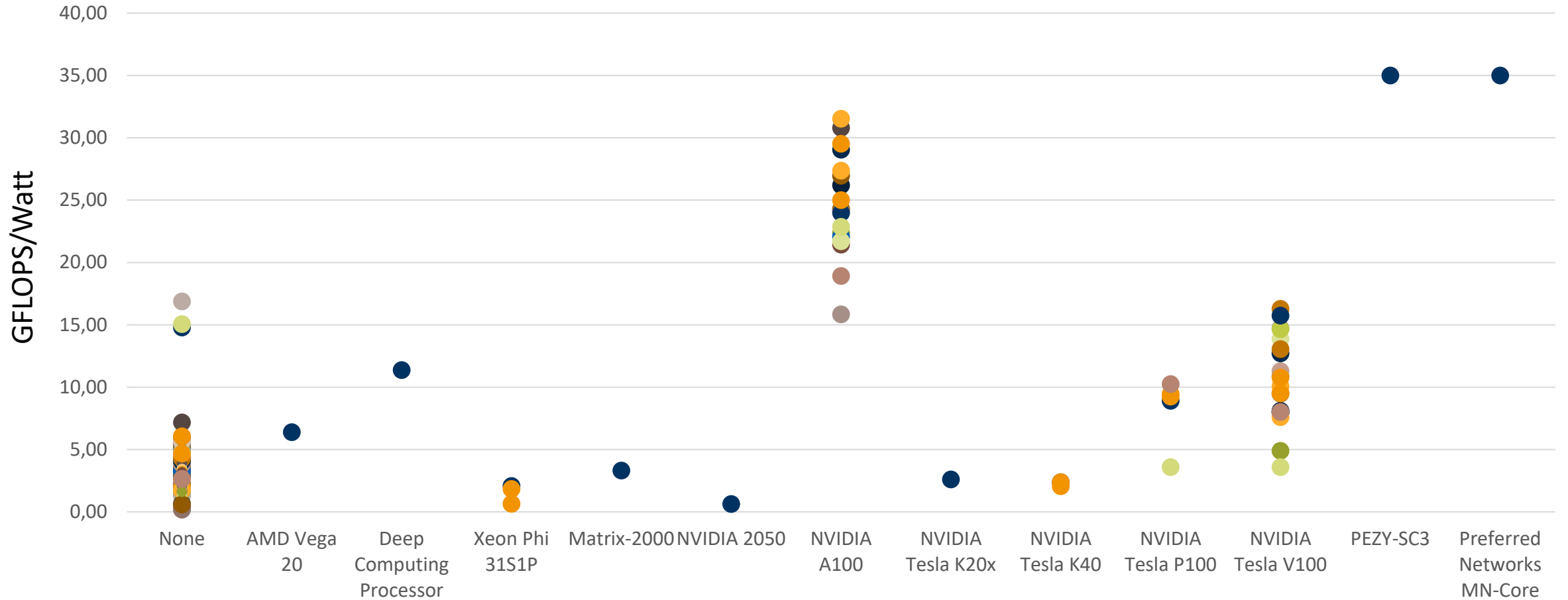  ▶ All 10 out of the top 10 entries in the November 2022 "Green 500" list are accelerator clusters (9/10 GPUs)

No. of accelerated systems in Top500



Intel Xeon Phi cancelled

general

# A closer look at power efficiency

- **8 of the top 10 systems above 10 GFLOPS/Watt use accelerators (2021 data)**

- **Exceptions:**
  - Fugaku: ARM-based, no accelerators
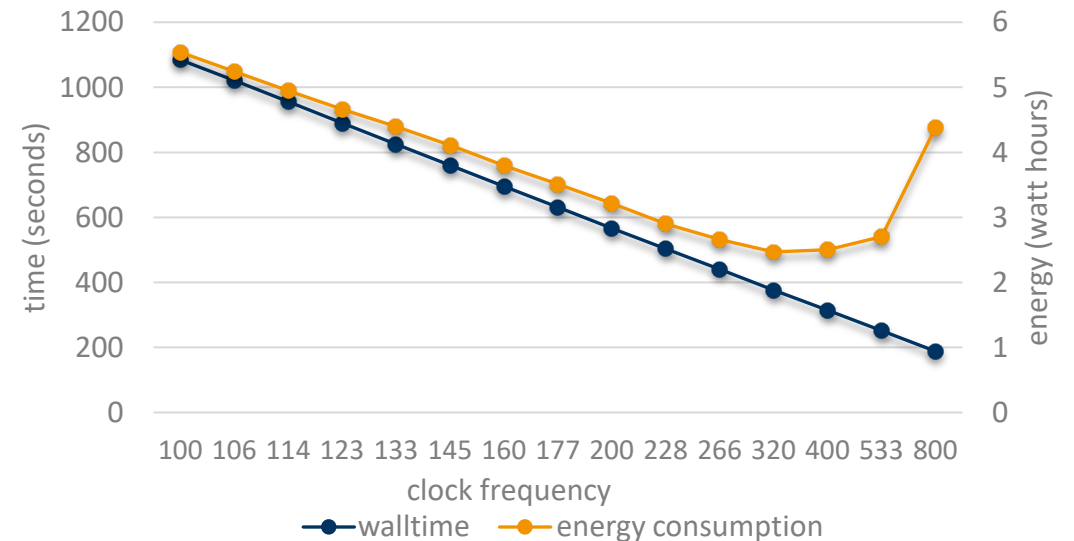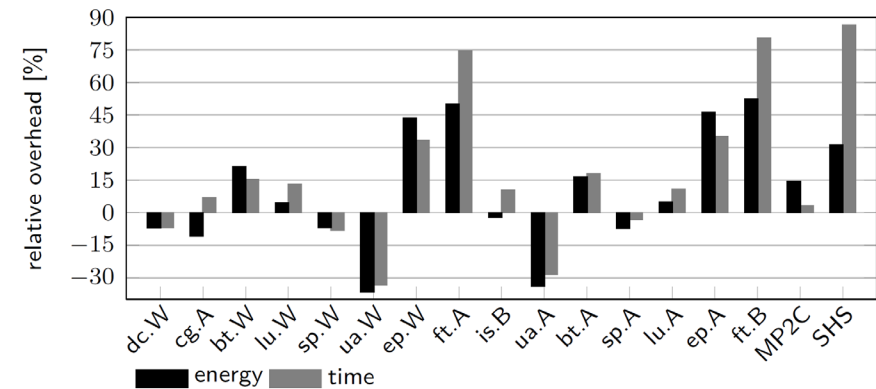  - Tianhe-2A: Matrix 2000 accelerators (128 core RISC CPUs)

### Rmax and Rpeak Power Efficiency



Legend: ■ Rmax/Power ■ Rpeak/Power

Labels: Nvidia V100, Nvidia A100, Nvidia V100

Y-axis: Power Efficiency [GFLOPS/Watt]

X-axis: Fugaku, Summit, Sierra, Sunway TaihuLight, Perlmutter, Selene, Tianhe-2A, Juwels Booster, HPC5

# Power efficiency of all TOP500 systems (2021)

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Reducing energy in computing: Tuning!

- Lots of research in software means for reducing energy consumption

- Top figure: effects of instruction mix on energy consumption of an IBM POWER7 CPU (using GCC vs. IBM XL compilers)
  - Result: In general, IBM XL produces more efficient binaries, but not always!

- Bottom figure: Dynamic Frequency and Voltage Scaling (DVFS) of the Intel SCC (experimental many-core CPU)
  - reduce clock frequency to save power and often also energy, effect heavily depends on workload
  - used in most CPUs these days (laptops, desktops, server, smartphones, etc.)
  - used on supercomputers (e.g. Energy Aware Runtime)
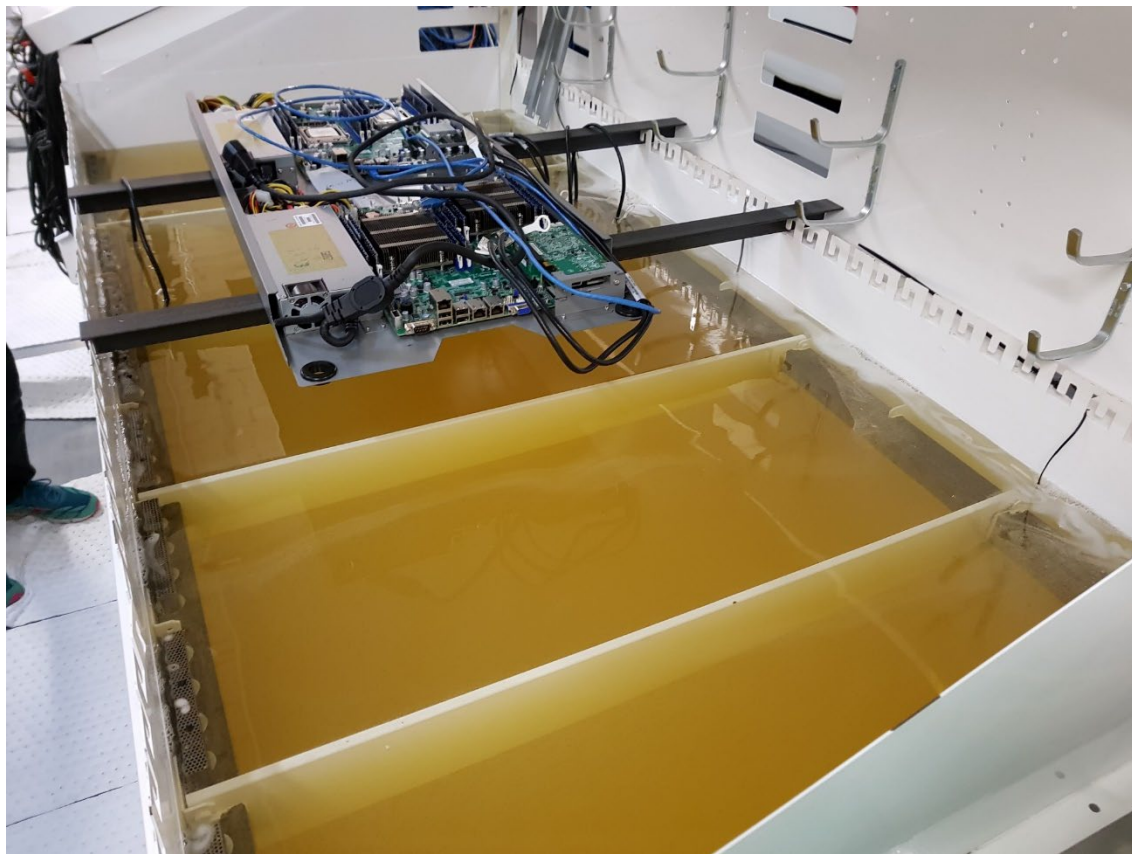
# Reducing energy in cooling: use oil!

▸ **VSC-3, fastest supercomputer in Austria in 2014**

  ▸ ranked 85$^{th}$ world-wide

  ▸ 32.768 cores

  ▸ 450 kW

  ▸ mechanical PUE of **1.02**!

    ▸ compare to VSC-2 (water-cooled): mPUE of 1.18

    ▸ VSC-4 (water-cooled): 1.05
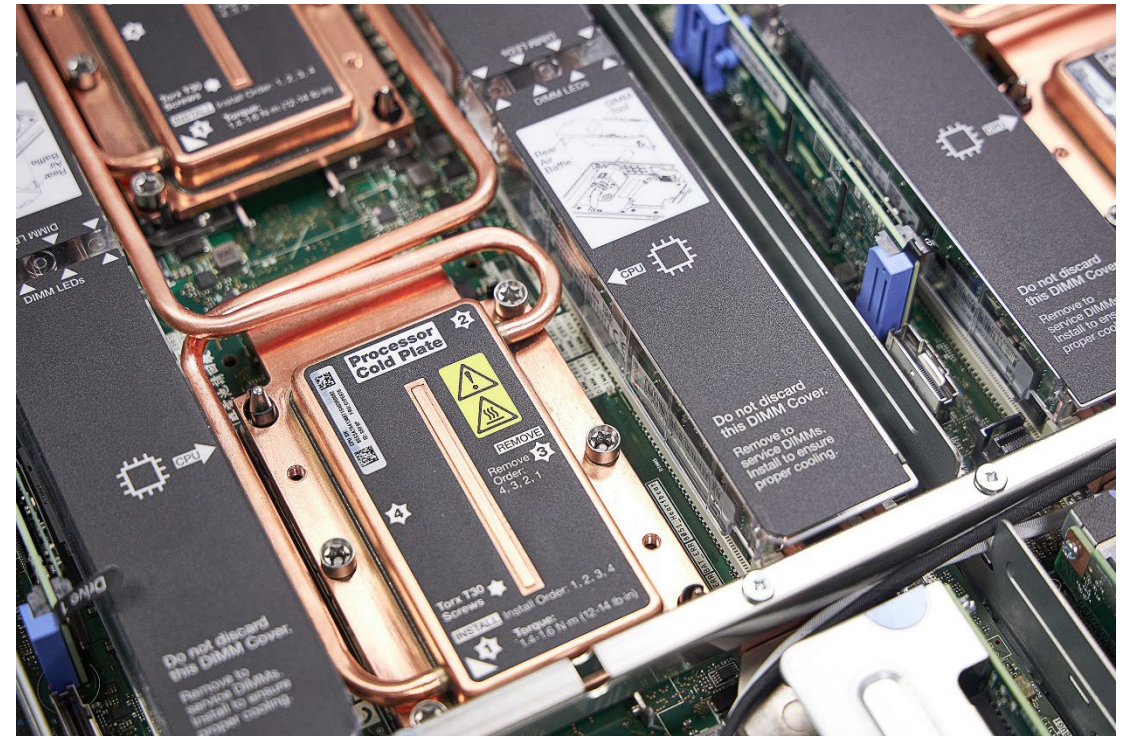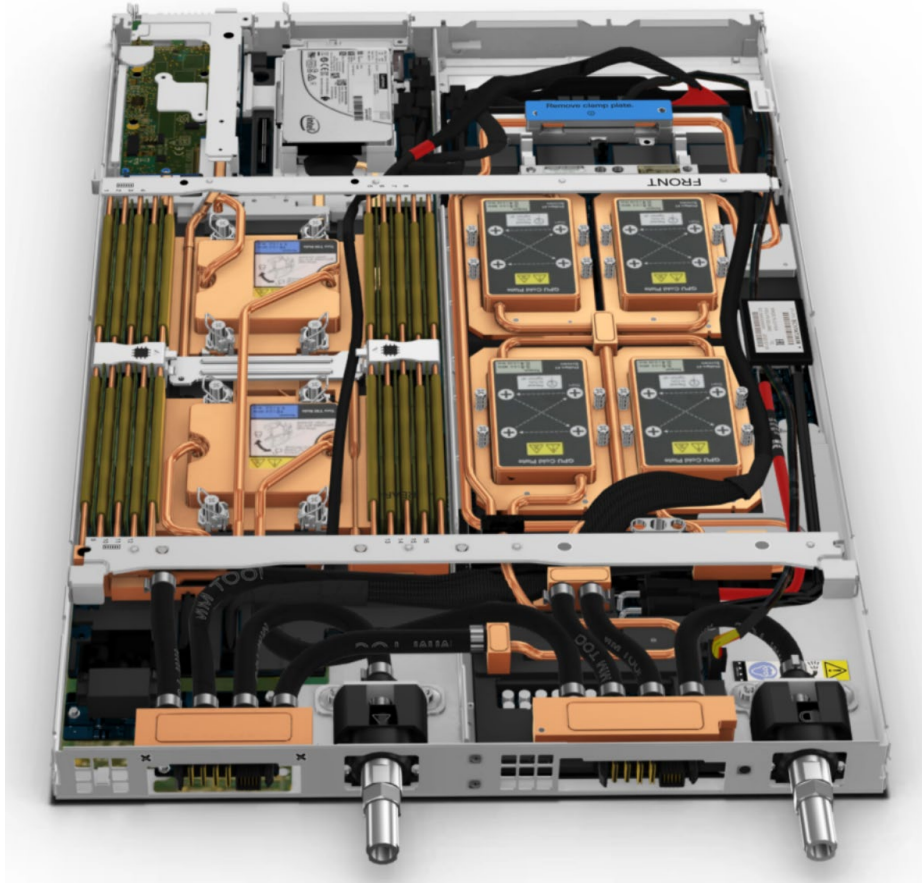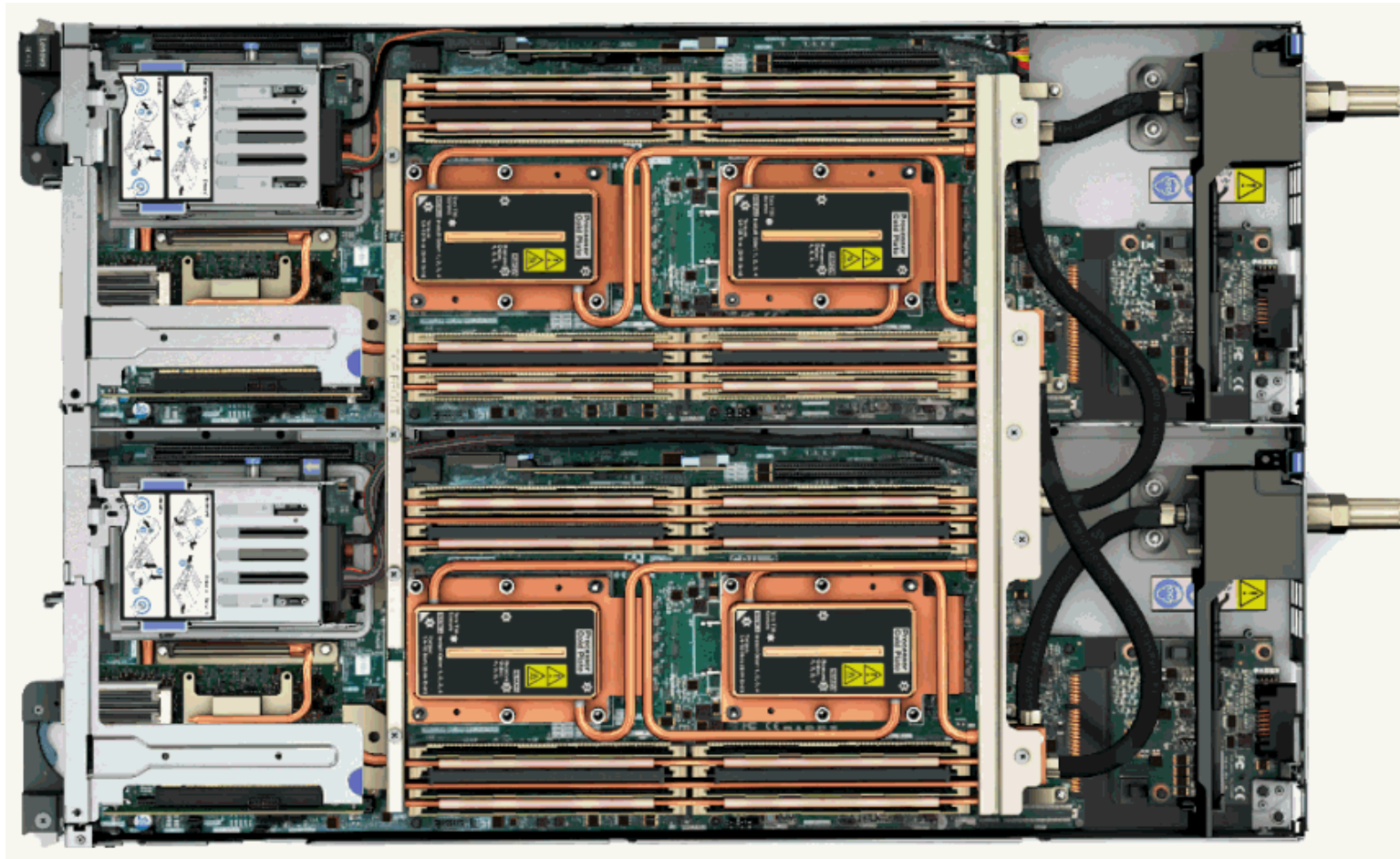
# "Immersion cooling"

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

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

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Lenovo SD650 direct water cooling



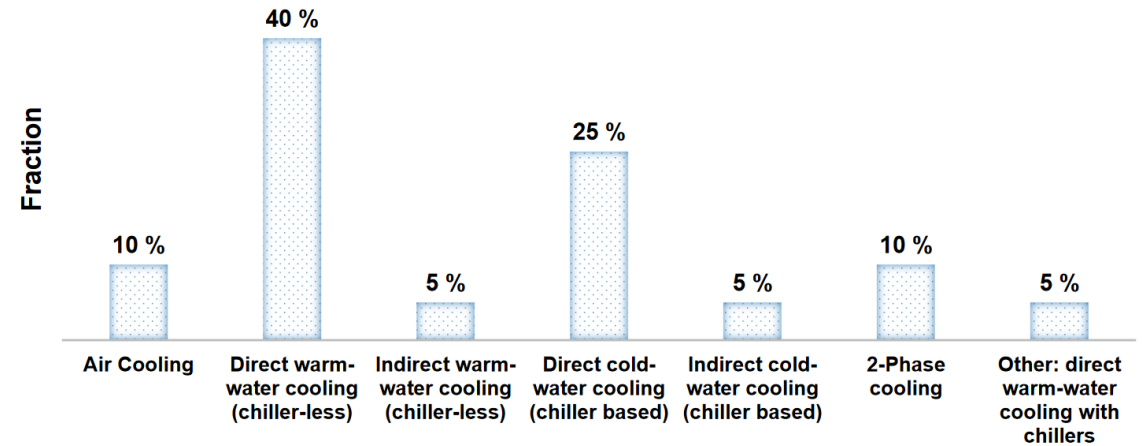Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner
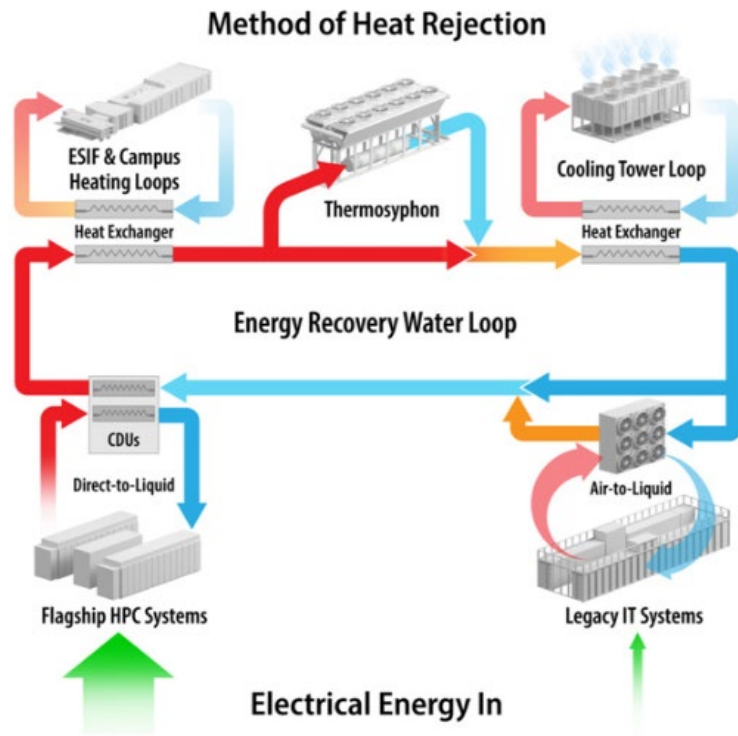
# Cooling technologies

▶ **Air cooling**

   ▶ easy to build and maintain, inefficient

▶ **Direct water cooling**

   ▶ warm: difficult to build and maintain, very efficient, only for cooler climates ("free air cooling")

   ▶ cold: difficult to build and maintain, semi-efficient, for warmer climates



2020 survey among tier-0 and tier-1 HPC sites in Europe

▶ **Indirect cooling**

   ▶ cool hardware with air, cool air with water

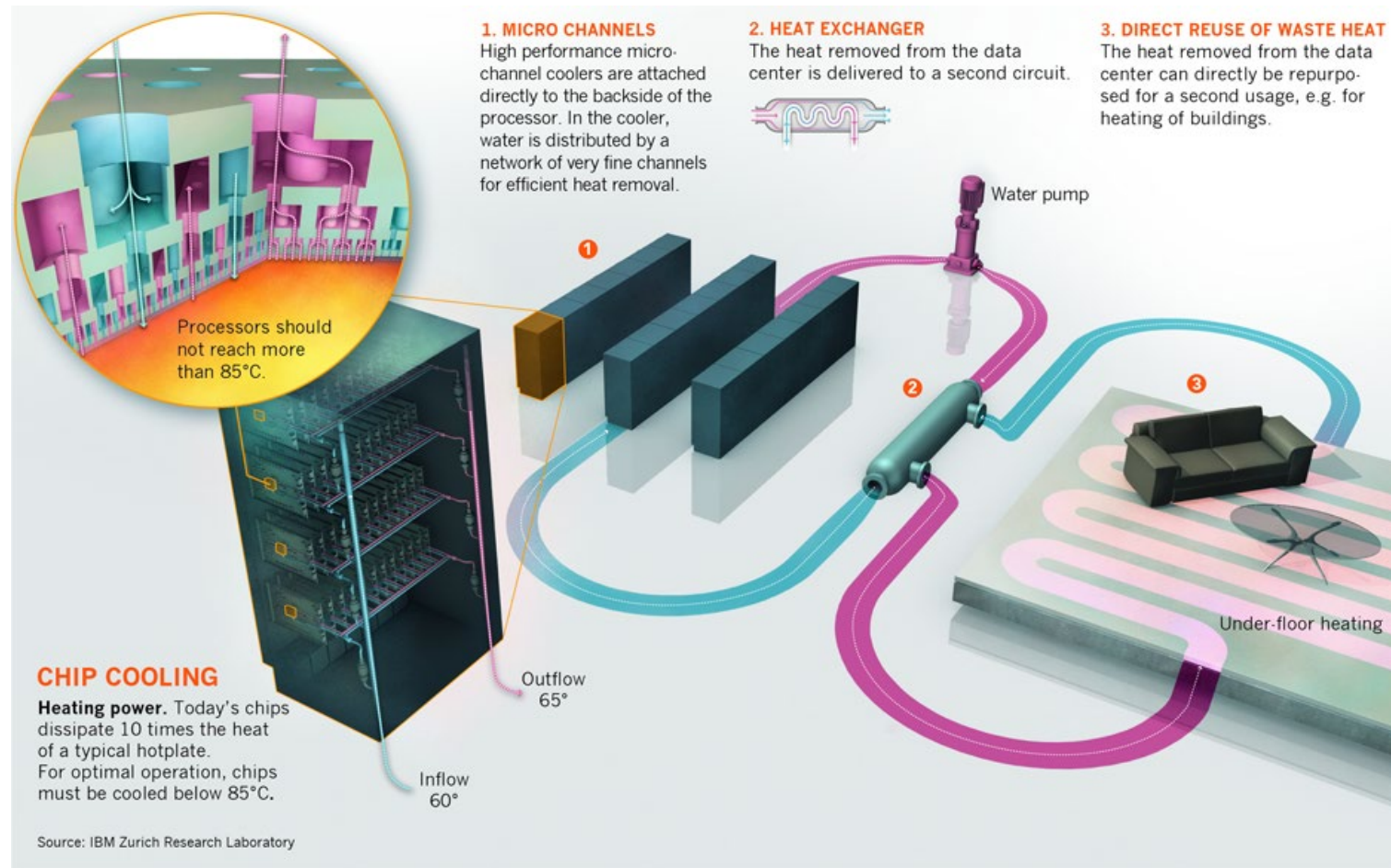Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# ESIF data center, NREL (PUE of 1.06)

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# How can we recycle any remaining energy consumption?



1. MICRO CHANNELS
High performance micro-channel coolers are attached directly to the backside of the processor. In the cooler, water is distributed by a network of very fine channels for efficient heat removal.

2. HEAT EXCHANGER
The heat removed from the data center is delivered to a second circuit.

3. DIRECT REUSE OF WASTE HEAT
The heat removed from the data center can directly be repurposed for a second usage, e.g. for heating of buildings.

Water pump

Processors should not reach more than 85°C.

Outflow 65°

Under-floor heating

CHIP COOLING
Heating power. Today's chips dissipate 10 times the heat of a typical hotplate. For optimal operation, chips must be cooled below 85°C.

Inflow 60°

Source: IBM Zurich Research Laboratory

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Sustainability in HPC vs. HPC for sustainability

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Open issues

- **There is more than just energy and power**
  - Carbon Usage Effectiveness (CUE)
  - Water Usage Effectiveness (WUE)
  - Space Usage Effectiveness (SpUE)

- **There are too many metrics and many are inaccurate**
  - Power Usage Effectiveness (PUE)
    - Partial PUE (pPUE)
  - Energy Reuse Effectiveness (ERE)
  - Energy Reuse Factor (ERF)

- **The metrics are often flawed**
  - e.g. PUE cannot be used to compare HPC sites in different climate zones

- **There are diverging interests**
  - Operator: minimize power/energy, maximize workload throughput
  - User: minimize wall time
  - Taxpayer/politicians: minimize costs

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner
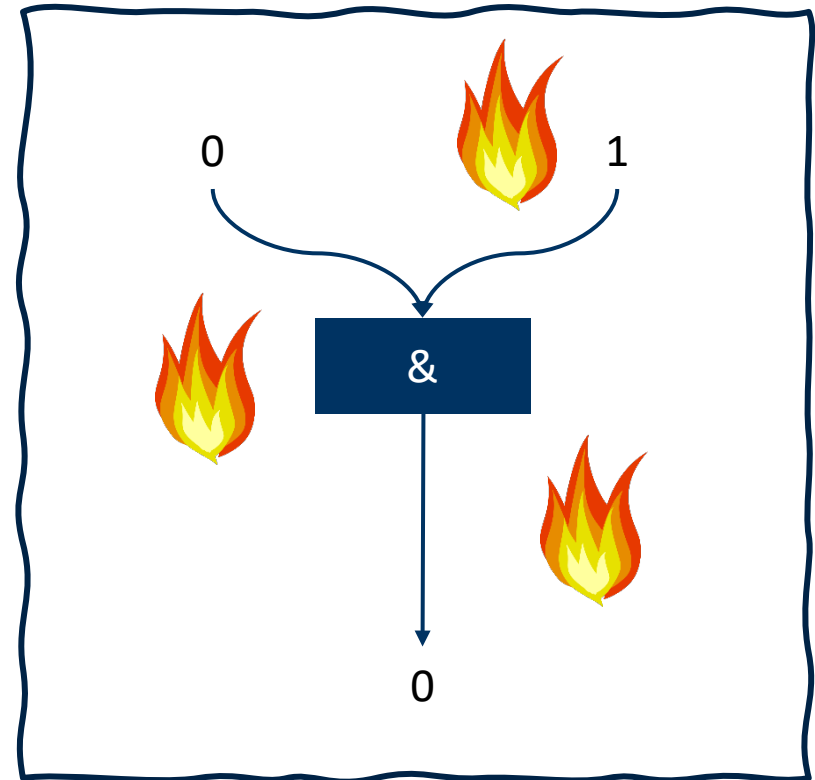
# Future developments and ideas

▸ **High-bandwidth memory (HBM)**
  ▸ Memory and computational units physically as close together as possible, minimize data transport distance

▸ **Fabrication size reduction**
  ▸ Research in new designs and materials (away from silicon) to decrease below ~2 nm threshold

▸ **Near-threshold voltage computing**
  ▸ operate CPUs below power safety limits, accept computational errors and mitigate in software (e.g. iterative solvers)

▸ **Special purpose hardware**
  ▸ Accelerators (scientific computing, AI, etc.)
  ▸ FPGAs
  ▸ Custom hardware designs for domain-specific problems

▸ **Optical computing**
  ▸ Use photons instead of electrons
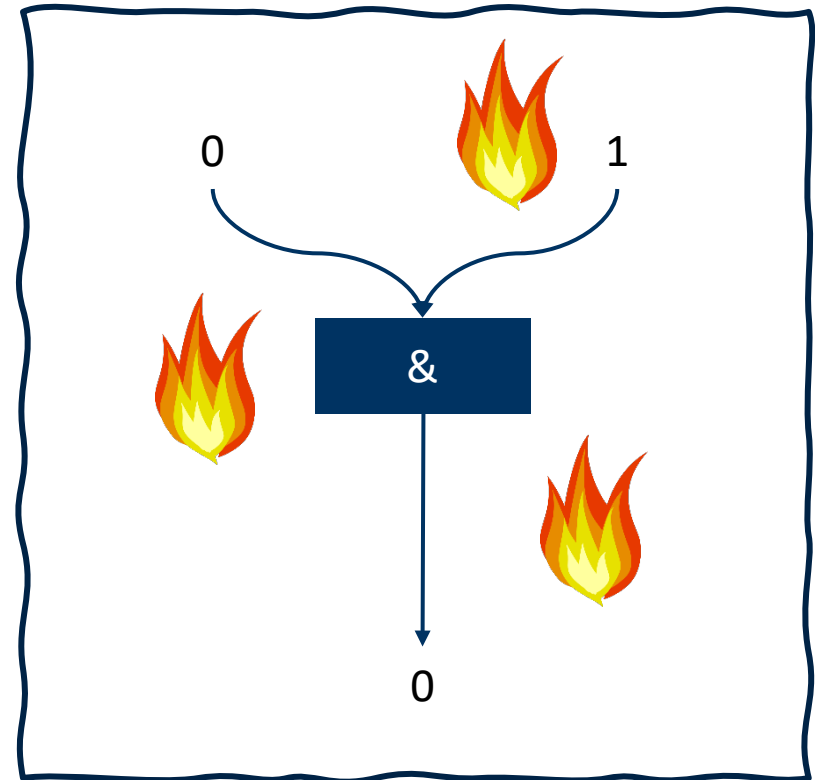  ▸ Various approaches in research, not clear yet if viable alternative

# 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 ¯\_(ツ)_/¯

  ▶ 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…

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Reversible computing and Landauer principle: the future?

▶ **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

# Today's takeaway

▶ **There's a lot of research and engineering going on**
  ▶ in sustainability for HPC
  ▶ in sustainability with HPC

▶ **Power/heat are the main limiting factors in HPC**
  ▶ almost everything uses water cooling these days
  ▶ waste heat is recycled as much as possible and "freely cooled" afterwards (no active chillers)
  ▶ short-term developments quite clear, long-term future very unclear

▶ **How to reach me/us**
  ▶ philipp.gschwandtner@uibk.ac.at
  ▶ https://dps.uibk.ac.at/~philipp
  ▶ https://uibk.ac.at/fz-hpc

Green HPC: Paving the Way for Sustainable Supercomputing | P. Gschwandtner

# Image sources

- Green HPC: https://www.hpcwire.com/2021/07/15/15-years-later-the-green500-continues-its-push-for-energy-efficiency-as-a-first-order-concern-in-hpc/, https://www.chemistryworld.com/features/oil-spill-cleanup/3008990.article, Marcel Ritter (UIBK), https://twitter.com/maven2mars/status/984440044659159040, https://www.nasa.gov/ames/image-feature/nasa-highlights-simulations-at-supercomputing-conference-like-aircraft-landing-gear

- TOP500 Trend: https://www.top500.org/statistics/perfdevel/

- Lenovo SD650 Water Cooling Images and Animation: https://lenovopress.lenovo.com/lp0636-thinksystem-sd650-direct-water-cooled-server-xeon-sp-gen-1

- Cooling Technology Survey: https://events.prace-ri.eu/event/1186/attachments/1587/2924/Shoukourian.pdf

- ESIF Data Center: https://www.nrel.gov/docs/fy21osti/79712.pdf

- IBM Research Energy Reuse: https://www.zurich.ibm.com/st/energy_efficiency/zeroemission.html

- Wind turbine: https://www.nrel.gov/docs/fy22osti/81212.pdf