

Week 5: Recap, outlook, and reminders

Portland State University
 Department of Electrical and Computer Engineering (ECE)
www.teuscher-lab.com
teuscher@pdx.edu



teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY



teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

What did you learn last week?

teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

Recap

- GPU execution model: SIMD (Single Instruction, Multiple Threads)
- GPU overhead:
 - Malloc
 - Copy data to GPU
 - Free
- A CUDA kernel is a function that runs on the GPU.
- Warp (32 threads, SIMD), thread block, kernel grid
- Parallelism is not always easy to find
- $\text{kernel} \ll \langle M, T \rangle \gg$: The kernel launches with a grid of M thread blocks. Each thread block has T parallel threads.



Portland State
 UNIVERSITY

Christof Teuscher teuscher@pdx.edu

[View Online](#)

SETH'S BLOG

The steps vs. the concept

If you memorize the steps, you have a direct, simple and fast path to obtain the result.
 Until the world changes.
 Even the tiniest shift in the system will render your memorization useless.
 On the other hand, if you understand the concept, you'll be able to produce the steps whenever you need them.

teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

Trends in AI Supercomputers

Konstantin F. Pilz¹ James Sanders² Robi Rahman² Lennart Heim^{2,3}

Abstract

Frontier AI development relies on powerful AI supercomputers, yet analysis of these systems is limited. We create a dataset of 500 AI supercomputers from 2019 to 2025 and analyze key trends in performance, power needs, hardware cost, ownership, and global distribution. We find that the computational performance of AI supercomputers has doubled every nine months, while hardware acquisition cost and power needs both doubled every year. The leading system in March 2025, XIAI's Colossal-2025, has a hardware cost of \$78.4M and uses 300,000 cores of power-hungry A100 AI chips, a hardware cost of \$78.4M and uses 300,000 cores of power-hungry A100 AI chips, a hardware cost of \$230 billion, and require 9 GW of power. Our analysis provides visibility into the AI supercomputer landscape, allowing policymakers to assess key AI trends like resource needs, ownership, and national competitiveness.

<https://arxiv.org/pdf/2504.16026.pdf>

teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

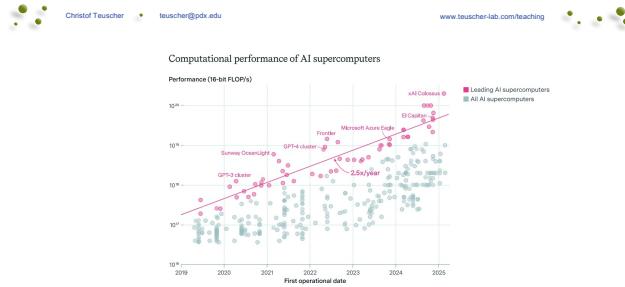


Figure 1: The performance of leading AI supercomputers (in FLOPs, for 16-bit precision) has doubled every 9 months (a rate of 2.5x per year).

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

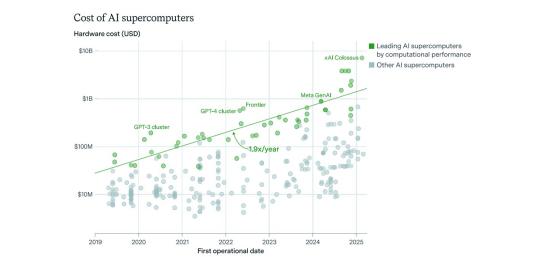


Figure 2: The cost of leading AI supercomputers (in 2025 USD) has doubled roughly every year.

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

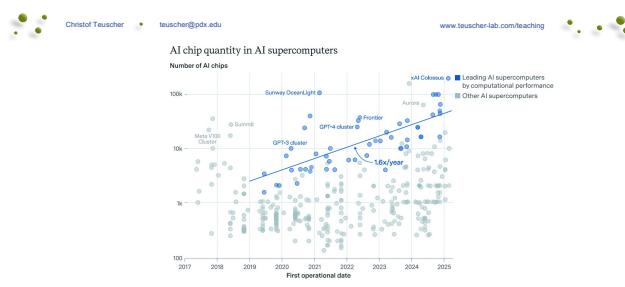


Figure 6: The number of AI chips in the leading AI supercomputers grew by 1.6x per year (90% CI: 1.5–1.8x). We start the regression in 2019, but gathered data further back to determine which 2019 AI supercomputers were in the top-10. Our pre-2019 data is too limited to include in the regression. See Section 2 for full methods.¹¹

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

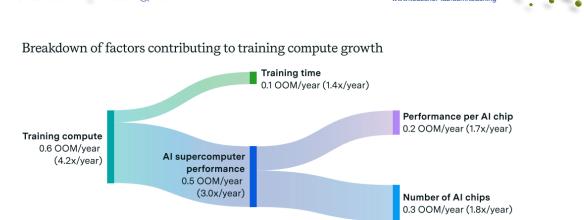


Figure 9: Overview of drivers of increasing training compute. OOM stands for orders of magnitude. AI supercomputer metrics are based on private sector systems and the highest computational performance across precisions.

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

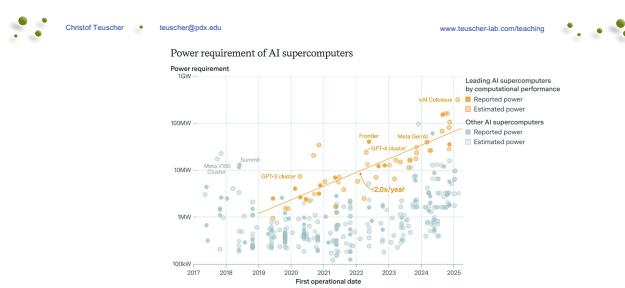


Figure 10: Peak data center power requirements of leading 10 AI supercomputers doubled every year (90% CI: 1.6–2.2x per year). We display reported power requirements whenever available. If not, we estimate capacity based on the number and type of chips used.

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

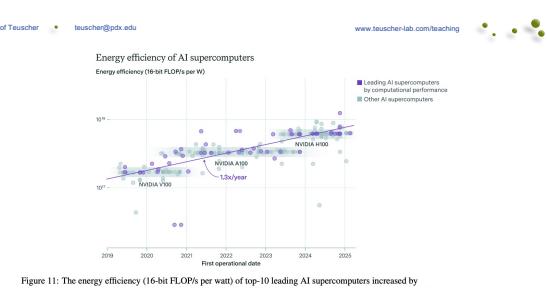


Figure 11: The energy efficiency (16-bit FLOPs per W) of leading AI supercomputers increased by 1.3x per year between 2019 and 2025 (90% CI: 1.25–1.43x). Adoption of new chips was the primary driver of energy efficiency improvements, with data center infrastructure efficiency only playing a minor role. We use reported power requirement whenever available and estimated power otherwise. For detailed methods and limitations, see Appendix B.4.

<https://arxiv.org/pdf/2504.16026>

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

Christof Teuscher teuscher@pdx.edu

<https://arxiv.org/pdf/2504.16026.pdf>

Aggregate AI supercomputer performance in our dataset by country (as of March 2025)

Country	Performance (H100-equivalents)
USA	850,000
China	700,000
European Union	60,000
Rest of the world	0

Our dataset covers approximately 10-20% of global aggregate AI supercomputer performance as of March 2025, so the true supercomputer capacity per country is likely 10-100x higher than what we report here. Note that some countries due to their reporting practices may have more or less coverage than others. The share of aggregate performance shown may change dramatically as new countries are added to our dataset.

Figure 15: Total AI supercomputer performance by country in H100-equivalents. To convert a system's performance to H100-equivalents, we first take the performance in the lowest precision its AI chips support, considering 32-bit, 16-bit, and 8-bit. We then divide by the 8-bit performance of the H100.

www.teuscher-lab.com/teaching Christof Teuscher teuscher@pdx.edu

<https://doi.org/10.1007/s43503-025-00053-x>

Share of aggregate AI supercomputer performance by country over time

Year	USA (%)	China (%)	Other (%)
2019	40	30	30
2020	42	32	26
2021	45	35	20
2022	48	38	14
2023	50	40	12
2024	52	42	10
2025	55	45	10

Our dataset covers approximately 10-20% of global aggregate AI supercomputer performance as of March 2025. While coverage varies across countries, sectors, and hardware types due to uneven public reporting, we believe the overall trend remains accurate. The chart shows a clear trend of increasing dominance by the USA and China, while other countries remain relatively stable or show minor growth over time.

Figure 16: Share of AI supercomputer computational performance by country over time. We are visualizing all countries that hold more than a 3% share at some point in time.

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

<https://arxiv.org/pdf/2504.16026.pdf>

Week	Monday	Wednesday (Codefest)
2	HW/AI/ML overview + codesign overview	Start main project: pick workload, start analysis, benchmark, ...
3	GPU architecture and programming for AI	Drafting a HW architecture, creating a model
4	Deep neural networks on GPUs	Coding HW description
5	Transformers on GPUs	First simulation + refinement
6	In-memory computation	Improving initial design
7	Neuromorphic chips: TrueNorth, Loihi, Akida	Simulation + refinement
8	Neuromorphic computing with mem-devices	Synthesizing design + benchmarking
9	Hardware accelerators for embedded systems <small>Non-codefest day</small>	Final improvements
10	Emerging technologies and future directions	Final tests, validation, verification, benchmarking

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

<https://doi.org/10.1007/s43503-025-00053-x>

This week

Monday

- CUDA
- Mapping deep neural nets (DNN/CNN) on GPUs
- TPUs
- Transformers on GPUs

Wednesday

- Codefest #5
 - Going to the gate level
 - Getting first performance numbers

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY