

AI: The Big Picture

A dive into the past, present and future of frontier AI development

Kai Zuberbühler

March 9, 2025

Table of Contents

1. Scaling
2. Data and Architecture
3. Finance
4. Chips
5. Power
6. Geopolitics

Scaling

Main Factors for Training Better AI Models

- Size (Parameter Count)
 - Increase the size of the model which requires **more compute**
 - Drawback: During inference, the model also requires **more compute**
 - Has not really grown anymore in the past years (GPT-4.5 is a notable exception).
- Quality of Data
 - One approach: Filter and augment data using AI which requires **more compute**
- Quantity of Data
 - Increase the size of the training dataset which requires **more compute**
 - Another approach: Generate new data using AI which requires **more compute**
- Architecture
 - Increase research which requires more researchers and **more compute**
 - Roughly equivalent to tripling compute per year (Ho et al., 2024)

Scaling Law Formulas

Performance = $\log(\text{Model Size} \cdot \text{Data Quantity} \cdot \text{Data Quality} \cdot \text{Architecture Efficiency})$

Figure 1: Oversimplified scaling law for language models

$$\hat{L}(N, D) \triangleq E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

Figure 2: Scaling law for language models where L is the loss of the model, N is the number of model parameters and D is the number of training tokens (Hoffmann et al., 2022)

Predicting Loss with Scaling Laws Is Easy ...

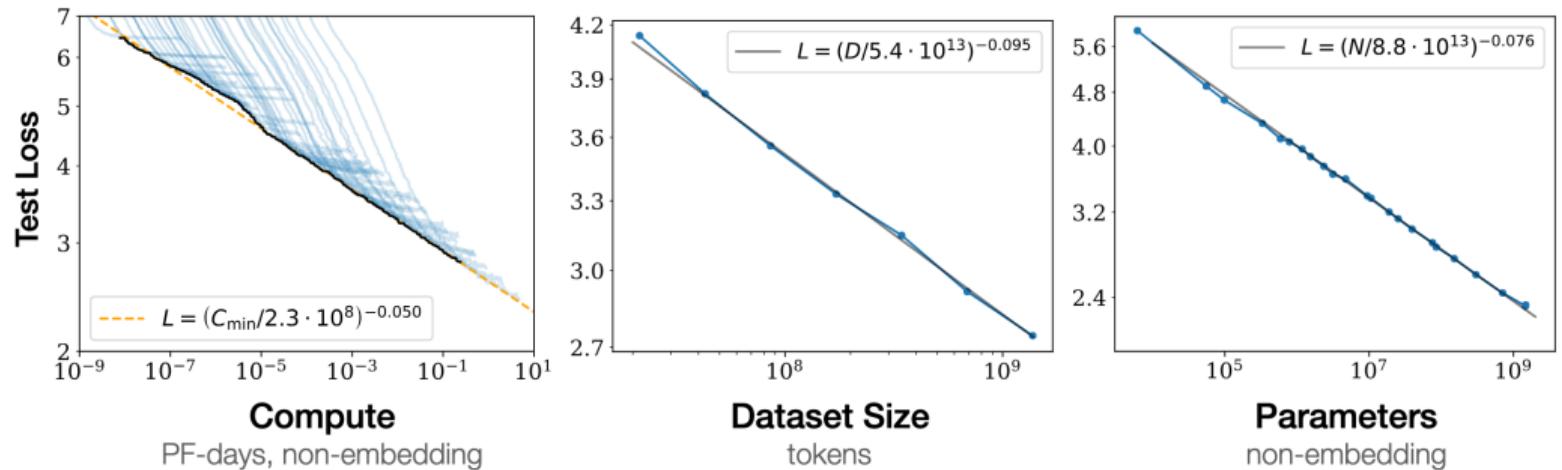


Figure 3: Observed log-linear language model performance improvements with increasing compute, dataset size and parameters (Kaplan et al., 2020)

... But Predicting Benchmark Performance Is Harder

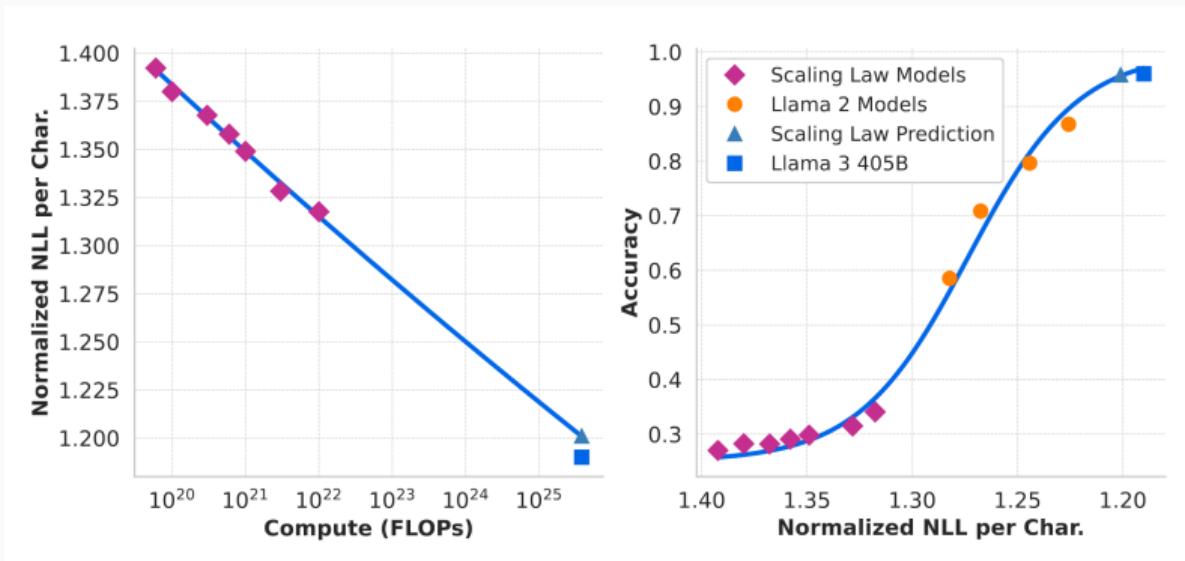


Figure 4: Predicted and observed performance on the ARC Challenge benchmark for the Llama 3 models. At some point, the performance on the benchmark starts to increase faster than an extrapolation would have suggested. (Grattafiori et al., 2024)

Training Compute of Frontier AI Models Grows by About 4-5x Per Year



Figure 5: Summary of compute trends in AI (Sevilla & Roldán, 2024)

Can AI Scaling Continue Through 2030?

- Epoch AI, a research institute, released a lengthy report in August 2024 about whether scaling AI at the current trajectory can continue through 2030 (Sevilla et al., 2024)
- They identify and research **electric power**, **chip manufacturing**, **data** and **latency** as key constraints.
- Much information in this presentation comes from this report or other reports by Epoch AI.
- Due to the emergence of reasoning models and news on the data center build-outs, the report is already pretty outdated.

Scaling: In a Nutshell

- The training compute of frontier AI models continues to grow by 4-5x per year, following a log-linear scaling trend that has persisted for over a decade.
- Improving AI performance requires more compute through scaling model size, data quantity, data quality, or architecture efficiency.
- Scaling laws reliably predict model loss, but benchmark performance is significantly less predictable.
- Epoch AI's analysis identifies electric power, chip manufacturing, data availability, and latency as key constraints that must be overcome to maintain the current scaling trajectory through 2030.

Data and Architecture

Training Phases of Language Models

- Pre-Training (Radford et al., 2019)
 - Model learns to predict the next token based on the given context
 - Dataset of text and other modalities mostly from publicly available data
 - Can be seen as “System 1” (fast thinking)
- Instruction-Tuning / RLHF (Ouyang et al., 2022; Wei et al., 2021)
 - Model learns to converse and follow instructions in a specific way
 - Dataset with user-assistant conversations or human preference data
 - Key to making language models usable for everyday people (see ChatGPT)
- Reinforcement Learning to Scale Chain-of-Thought Reasoning (OpenAI, 2024b)
 - Model learns to use chain-of-thought to “think” longer to solve problems
 - Dataset with problems and corresponding solutions through which the model develops successful chain-of-thought reasoning strategies on its own
 - Can be seen as “System 2” (slow thinking)
 - First introduced in September 2024 by OpenAI and has become the main focus of AI labs in the past few months

Non-Reasoning Models Significantly Under-Perform Humans In Key Areas

Key domains with benchmark examples:

- Reasoning
 - ARC-AGI-Pub (Chollet, 2024)
 - Semi-Private Eval: 14% (Claude 3.5 Sonnet) vs. 77% (MTurkers)
 - Public Eval: 21% (Claude 3.5 Sonnet) vs. 64% (Humans) (LeGris et al., 2024)
- Planning
 - PlanBench (Valmeekam et al., 2024)
 - Mystery Blocksworld (0-shot): 1% (Llama 3.1 405B)
- Autonomy
 - OSWorld (Xie et al., 2024)
 - 22% (Claude 3.5 Sonnet, with Screenshots) vs. 72% (Humans)
 - GAIA (Mialon et al., 2023)
 - 65% (h2oGPTe Agent using Claude 3.5 Sonnet) vs. 92% (Humans)

Will We Run Out Of Human-Made Data for Pre-Training?

- The indexed web is estimated to contain roughly 500 trillion tokens of human-generated text and will grow by approximately 50% by 2030
 - About 30 times larger than the largest datasets currently used (about 15 trillion tokens). With the current trajectory, we would hit the “data wall” in about five years.
- It's possible to train several times over the same data (multiple epochs)
- Inclusion of other modalities (like images, video and audio), the “deep web” (that includes many social media platforms) and private data can further expand the effective data pool
 - Converting multimodal content into text-equivalent tokens could roughly triple the available dataset size
- The remaining data is likely of lower quality than the data already used
- The median projection from Epoch AI is that it's possible to train models with 80,000x the compute of GPT-4 without significant diminishing returns.

Can We Just Generate More Data?

- Using AI to generate training data is increasingly used, particularly for model distillation in post-training, but still unproven at large scale for pre-training data
- This brings the risk of producing low-quality data, e.g. with hallucinations or a lack of diversity, potentially resulting in “model collapse” (Shumailov et al., 2023)
- Scaling inference-time compute is a possible way to increase model performance when generating data (Villalobos & Atkinson, 2023), especially now with reasoning models
- Using AI-driven verifiers (and non-AI verifiers, e.g. code compilers) to filter generated data by assessed quality is a promising approach given that verification is easier than generation (Feng et al., 2024)
- Any large-scale application will likely use a sophisticated synthetic data generation pipeline and focus on generating new data out of existing data (Hao et al., 2025)
- It’s unclear how cost-efficient high-quality synthetic data is

The Costs of LLMs Are Plummeting

- The cost of an AI model with GPT-4-level quality has dropped by around 75x between Q1 and Q4 2024. (Artificial Analysis, 2024)
- This drop is due to a variety of factors including training on more and higher quality data, post-training techniques, distillation, architectural improvements (e.g. increased sparsity through mixture-of-experts), software improvements and potentially price dumping.
- An analysis of open-weight models shows that the “capability density” of LLMs (based on the total parameter count of models and their results on widely used benchmarks) is currently growing exponentially, doubling approximately every three months. (Xiao et al., 2024)

What are “Reasoning Models” Exactly?

- The model “thinks” using long chain-of-thought before giving a final answer.
- Reinforcement learning (RL) is used on an existing instruction-tuned language model to scale and improve this chain-of-thought “thinking”.
- Each answer is graded by a reward function to steer the model towards longer and more efficient reasoning. The reward function could e.g. be an AI model comparing the answer to a sample solution or unit tests for code. (Du et al., 2025)
- Such problem-answer pairs might come from a specialized dataset or be extracted from the web using AI. (Yeo et al., 2025)
- Through this training, the model discovers and learns better and longer reasoning chains that result in better answers.
- Small-scale domain-specific experiments show that even models as small as 1.5B parameters can develop long chain-of-thought abilities, even on consumer GPUs. (Pan et al., 2025; Unsloth, 2025)

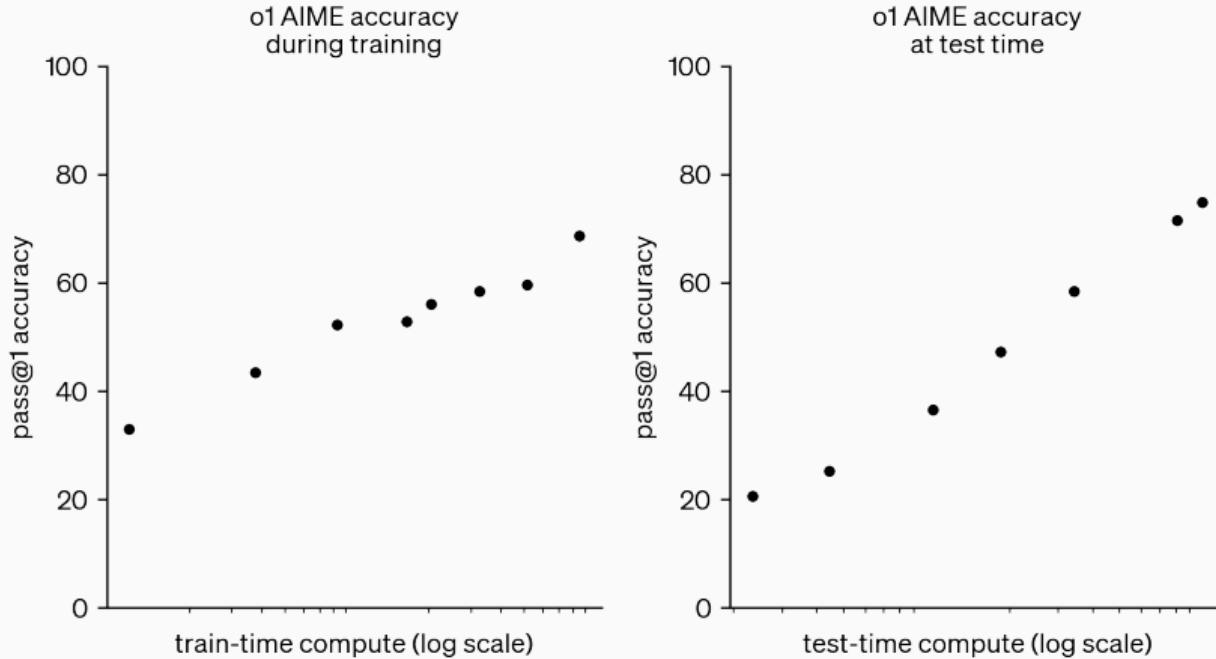


Figure 6: Reasoning models introduce two new compute scaling paradigms: Scaling chain-of-thought reasoning with reinforcement learning during post-training and scaling the length of the chain-of-thought during inference. (OpenAI, 2024a)

Benefits of Reasoning Models

- Reasoning models outperform non-reasoning models in many tasks, most notably and significantly in reasoning, mathematics and coding.
 - Top reasoning models are now reaching Olympiad-level performance in competitive coding and competitive mathematics.
- Reinforcement learning (RL) seems to generalize significantly better than supervised fine-tuning (SFT) for LLM post-training. (Chu et al., 2025)
- Therefore, RL training is likely significantly more data efficient than pre-training, although the details have yet to be properly researched.
- This might significantly push the “data wall” into the future.
- While the model usually “thinks” longer for harder and shorter for easier problems automatically, it’s also possible to manually steer how long the model “thinks” .
- RL can actually be creative and result in novel solutions.

Drawbacks and Limitations of Reasoning Models

- Reasoning models generate many more tokens per query than non-reasoning models, requiring **more compute** and resulting in slower response times (up to several minutes) and higher costs.
- Reinforcement learning (RL) brings the risk of “reward hacking” where the model exploits loopholes in the reward function without achieving the desired behavior.
- RL requires problems where answers can be subjectively graded, ideally in a cheap and fully automated way.
- It’s currently unclear whether significant improvements in agency, creative writing, instruction following, long-context handling, vision, physical world understanding and some other domains can be achieved.

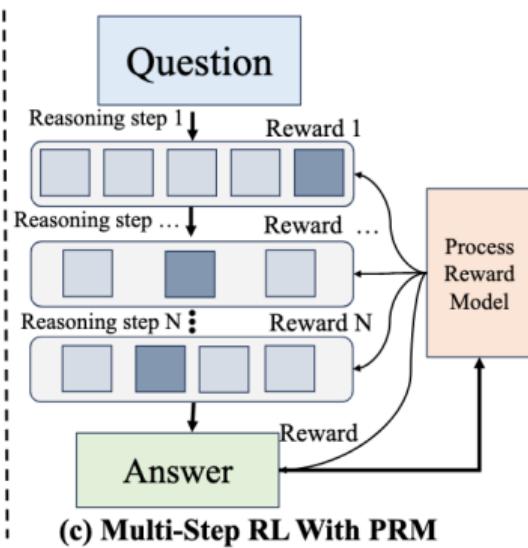
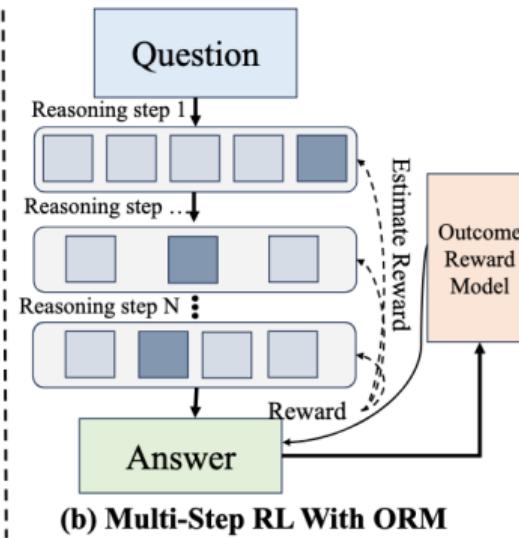
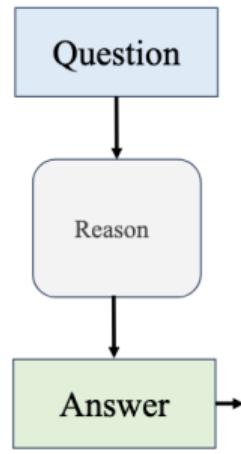


Figure 7: Comparison of different approaches to rewarding models during RL training. The model reasons using chain-of-thought before giving a final answer. Reward models grade the answer and/or the reasoning steps of the model being trained. (F. Xu et al., 2025)

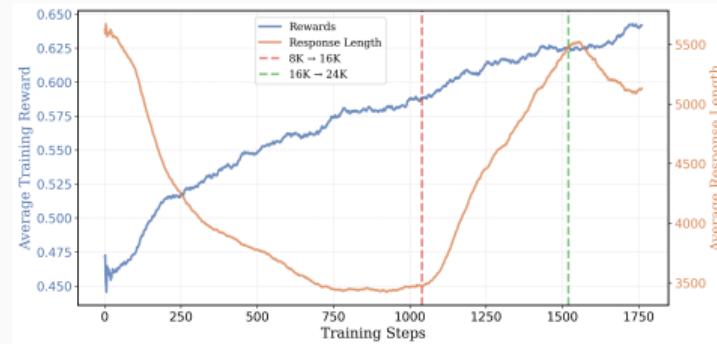
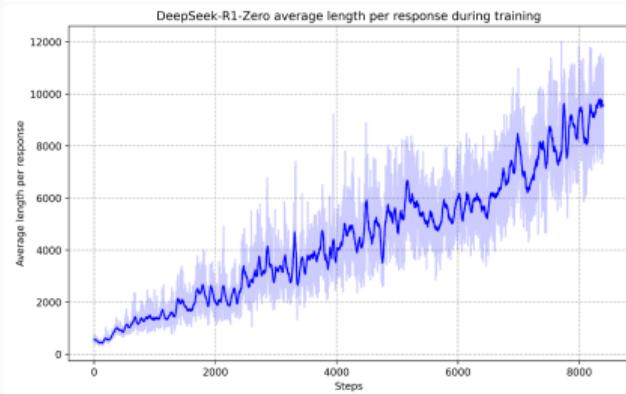
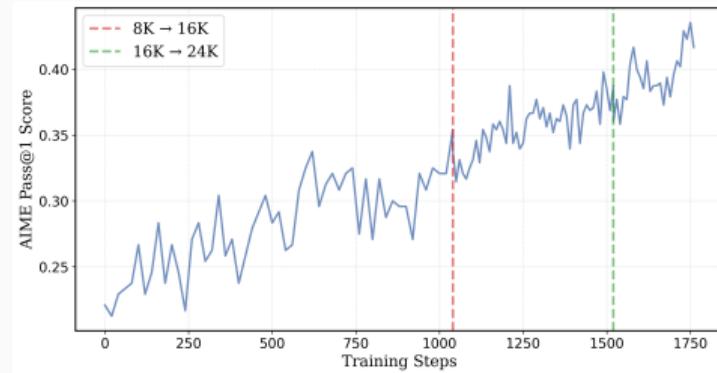
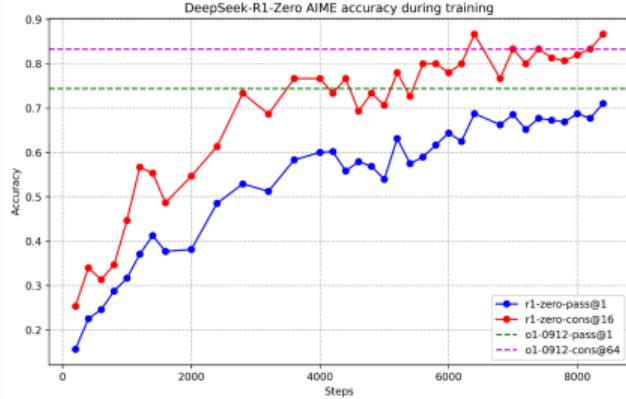


Figure 8: Accuracy on AIME (top) and response length (bottom) of Deepseek-R1-Zero (left) and DeepScaleR-1.5B (right) during RL training. (Guo et al., 2025; Luo et al., 2025)

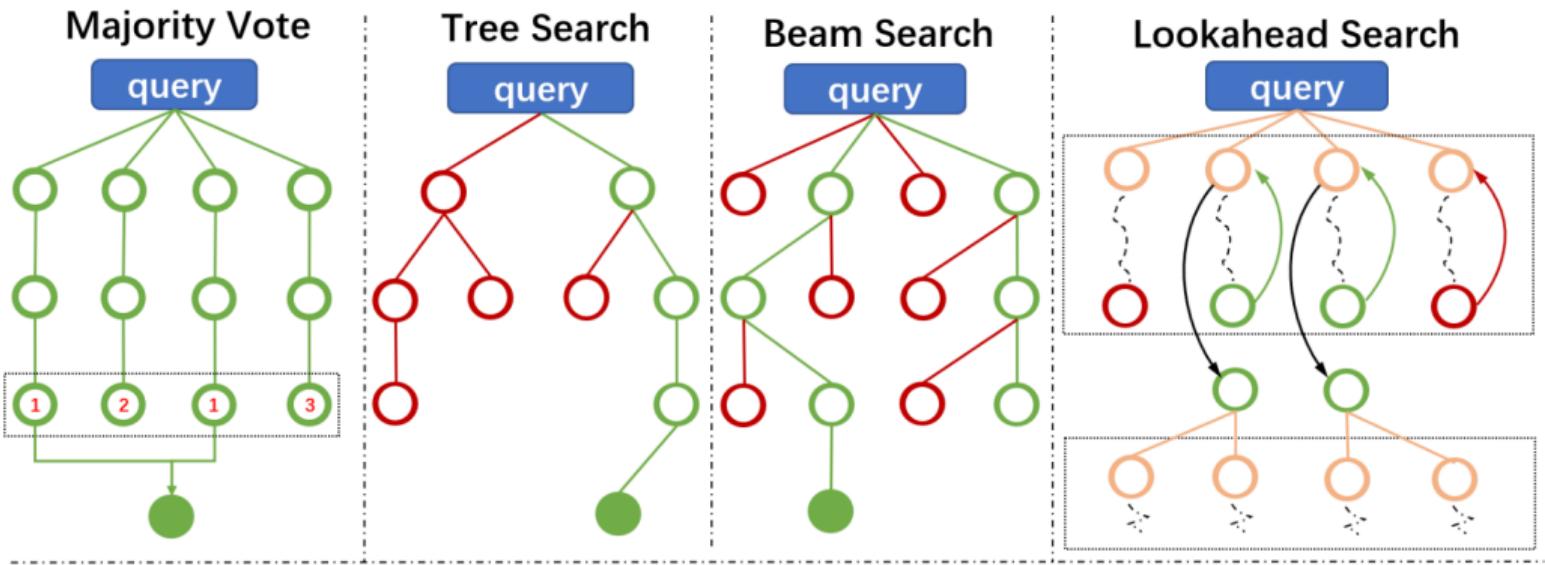


Figure 9: Comparison of different search algorithms to further scale inference-time compute. Such a method is likely used for the "Pro" mode of OpenAI o1. (F. Xu et al., 2025)

Reasoning Models Make Big Jumps in Reasoning and Planning

- Reasoning
 - ARC-AGI-Pub (Chollet, 2024)
 - Semi-Private Eval: 76% (OpenAI o3) vs. 77% (MTurkers)
 - Public Eval: 83% (OpenAI o3) vs. 64% (Humans) (LeGris et al., 2024)
- Planning
 - PlanBench (Valmeekam et al., 2024)
 - Mystery Blocksworld (0-shot): 53% (OpenAI o1-preview)
- Autonomy
 - OSWorld (Xie et al., 2024)
 - 38% (OpenAI CUA) vs. 72% (Humans)
 - GAIA (Mialon et al., 2023)
 - 67% (OpenAI Deep Research) (OpenAI, 2025) vs. 92% (Humans)

How Far Have Reasoning Models Been Scaled so Far?

- DeepSeek R1 is estimated to have only used 6.1×10^{23} FLOP during RL training on 2 trillion tokens (Erdil, 2025)
 - 15,405 H800 days or roughly 1/34th of the compute used to train GPT-4
 - If current compute scaling trends can be sustained, a training run requiring more than 300,000x the compute would be possible in 2030
- It is unknown what type of data and how much data was used for training any state-of-the-art reasoning model
- State-of-the-art reasoning models like OpenAI o3-mini and DeepSeek R1 generate up to several thousands of reasoning tokens for each response
 - Current state-of-the-art language models support context lengths ranging from 128,000 to 2,000,000 tokens

Could Reinforcement Learning Solve Autonomy?

- AI agents aim to make AI able to act in digital environments like a virtual machine or a web browser to solve multi-step tasks
- AI agents currently still have major issues like difficulty interacting with GUIs, difficulty solving long-horizon tasks and a lack of commonsense and social skills (F. F. Xu et al., 2024)
- Applying reinforcement learning to AI agents with verifiers checking the final results of multi-step tasks might be a promising approach (Pan et al., 2024)
- The tasks and the environments might be human-made or even AI-generated (Hu et al., 2024) and need to overcome the simulation-to-reality gap

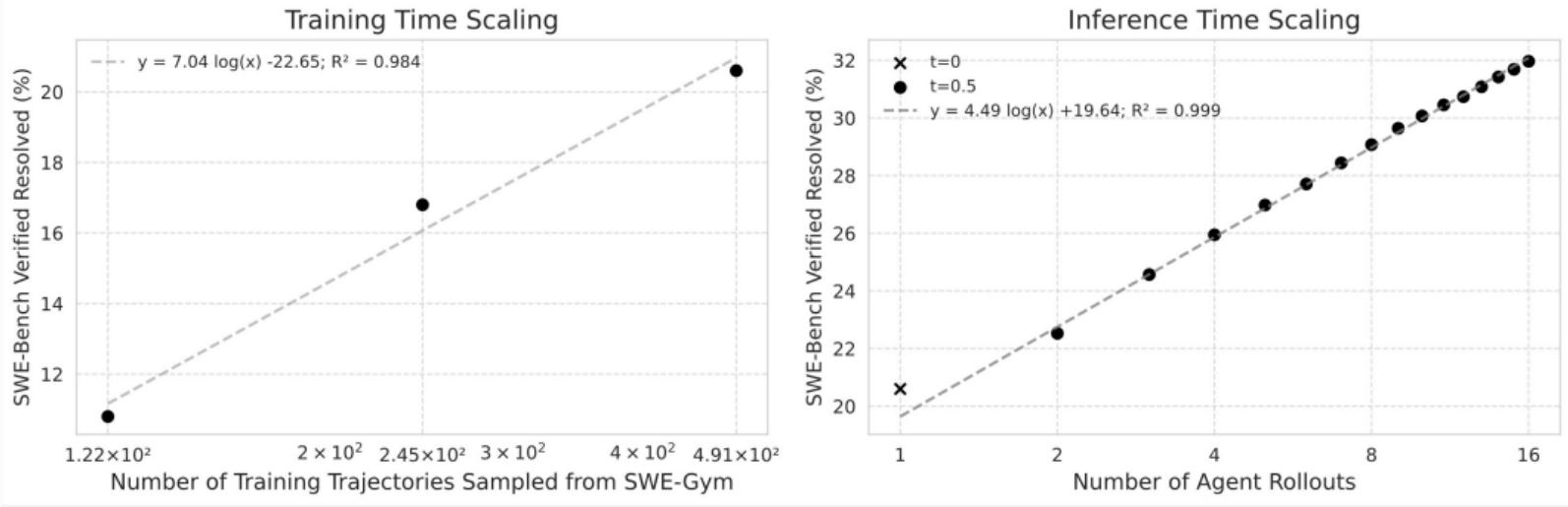


Figure 10: Log-linear increase in performance of software engineering agents through training compute scaling and inference compute scaling. The agents were trained on 491 agent trajectories sampled from agents solving tasks in the SWE-Gym training environment. (Pan et al., 2024) Large-scale agent training, especially the potential of reinforcement learning by scaling chain-of-throught reasoning, still has to be properly researched.

Data and Architecture: In a Nutshell

- Language model training involves three distinct phases: pre-training, instruction-tuning, and now reinforcement learning (RL) to scale reasoning.
- The industry faces a “data wall” within five years at current scaling trends.
- AI-generated training data is increasingly used but brings risks, making verification of generated content crucial for maintaining quality.
- Reasoning models represent a fundamental shift, enabling breakthrough performance in areas where previous models significantly underperformed humans.
- These models require more inference compute but may help overcome the looming “data wall” since RL training appears more data-efficient than pre-training.
- Despite impressive gains in reasoning, planning, math and coding, reasoning models have barely been scaled so far due to them only being discovered recently.
- The potential application of RL to agent training could transform AI autonomy, though significant challenges and uncertainties remain.

Finance

How Much Would Autonomous Generalist AI Be Worth?

- Imagine a system that can autonomously complete basically any job that a human can perform on a computer.
- Let's make a simple calculation to calculate how much such a system might be worth:
 - Let's very conservatively assume that such an AI could perform work equivalent to 1% of world GDP (imagine a 1% increase in global productivity).
 - Such AI would therefore make around \$1.15 trillion annually (1% of world GDP).
 - A tech company with such revenues could be valued at \$6.3 trillion (using current P/S ratio of the NASDAQ 100).
- Even using very conservative estimates, it's clear that such AI would be worth an astronomical amount.
- Investing trillions of dollars into AI research and development seems worth it once it seems possible that autonomous generalist AI can be achieved.

How Much Do Leading Tech Companies Make?

Company	Country	Revenue	Net Income
Amazon	USA	\$514B	-\$3B
Apple	USA	\$394B	\$100B
Google	USA	\$283B	\$60B
Microsoft	USA	\$198B	\$73B
Alibaba	China	\$135B	\$10B
Meta	USA	\$117B	\$23B
ByteDance	China	\$85B	\$25B
Tencent	China	\$82B	\$28B
SAP	Germany	\$31B	\$2B
Baidu	China	\$18B	\$1B

Table 1: Revenue and net income of selected technology companies in 2022

How Much Are Leading Tech Companies Spending Now?

Company	Country	Q1 2023	Q3/Q4 2024		
Amazon	USA	\$57B	\$111B	+\$55B	+96%
Apple	USA	\$12B	\$12B	+\$0B	+1%
Google	USA	\$25B	\$57B	+\$32B	+127%
Microsoft	USA	\$26B	\$63B	+\$37B	+139%
Alibaba	China	\$1B	\$10B*	+\$8B	+560%
Meta	USA	\$27B	\$58B	+\$30B	+111%
ByteDance	China	?	?	?	?
Tencent	China	\$2B	\$14B*	+\$12B	+497%
SAP	Germany	€1B	€1B	+€0B	+5%
Baidu	China	\$1B	\$1B*	+\$0B	+24%

Table 2: Annualized capital expenditures (excluding leases) of selected technology companies

How Much Are Leading Tech Companies Planning to Spend?

- Amazon expects total capital expenditures (capex) in 2025 to be at around \$105 billion with “the vast majority of that capex spend is on AI for AWS” for this “once-in-a-lifetime type of business opportunity”. (CNBC, 2025)
- Google expects total capex to “continue to increase in 2025”, hitting \$75 billion with the majority of it for data centers. (Data Center Dynamics, 2025b)
- Microsoft is “on track to invest approximately \$80 billion” on AI infrastructure in fiscal year 2025 while noting “in many ways, artificial intelligence is the electricity of our age”. (Microsoft, 2025)
 - Equivalent to 29% of total expected revenue in 2025 (Yahoo Finance, 2025)
- Meta expects total capex in 2025 to be “in the range of \$60 billion to \$65 billion” while noting the “hundreds of billions of dollars that we will invest into AI infrastructure over the long term” (Meta, 2025)
- For comparison: Peak annual cost of Apollo and related programs was \$48 billion (1965, adjusted for inflation) (Planetary Society, 2019)

Finance: In a Nutshell

- Amazon, Google, Microsoft and Meta are each making their most expensive bet yet with each of them roughly doubling their total capital expenditures in the past two years in order to spend unprecedented amounts on AI infrastructure.
- The annualized spending on AI infrastructure of each of these companies is now reaching the annual costs of the most expensive megaprojects ever (such as the Apollo program).

Chips

Chips Used For Training AI in the Past, Present and Future

- Graphics Cards (2012-2016)
 - Krizhevsky et al. (2012) championed the use of graphics cards to train their groundbreaking AlexNet model
 - Speeds up calculations with parallel computing and faster memory
- AI Accelerators (2015-present)
 - The first AI accelerators included Google's TPU v1 and NVIDIA's P100
 - Designed for ML and HPC, usually featuring tensor cores, support for lower-precision arithmetic, high-bandwidth memory (HBM) and 2.5D packaging
- ASICs and In-Memory Computing (2019-present)
 - Examples include the Cerebras WSE (Cerebras, 2024), the Groq LPU (Groq, 2024) (only for inference) and the Etched Sohu (Etched, 2024) (only for inference)
- Photonic Computing (Research Stage)
- Analog Neuromorphic Computing (Research Stage)

Computational performance improves 12x when switching from FP32 to INT8

Aggregated trends in AI accelerator performance went through improvements from tensor cores and compact number formats. Arrows show averaged trend improvements since a number format was introduced.

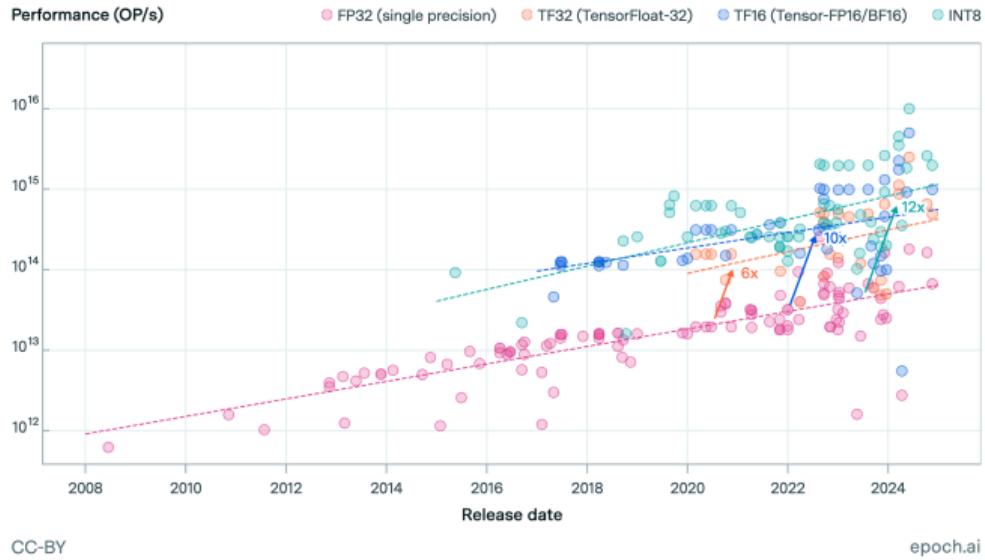


Figure 11: AI accelerator performance is increasing rapidly, especially when including support for lower-precision arithmetic. The NVIDIA B200 (released in 2024) using FP8 offers 2,812x the peak performance of the NVIDIA GeForce GTX 580 (released in 2010) using FP32 (used by Krizhevsky et al. (2012) to train AlexNet). (Epoch AI, 2024b)

NVIDIA: King of the AI Hardware Market

- Dominates the global AI hardware market, selling the overwhelming majority of AI accelerators
- 10x data center revenue growth in just two years
- Big profit margins on AI accelerators (produced for thousands of dollars, sold for tens of thousands of dollars)
- Ramped up production of Hopper GPUs during 2023 and 2024
- Currently ramping up production of the new Blackwell GPUs which are sold out for the next 12 months. (TweakTown, 2024a)
 - The B200 offers roughly a 2.5x performance improvement over the H100.
 - The new NVL72 racks connect 72 GPUs, ideal for large-scale AI training.

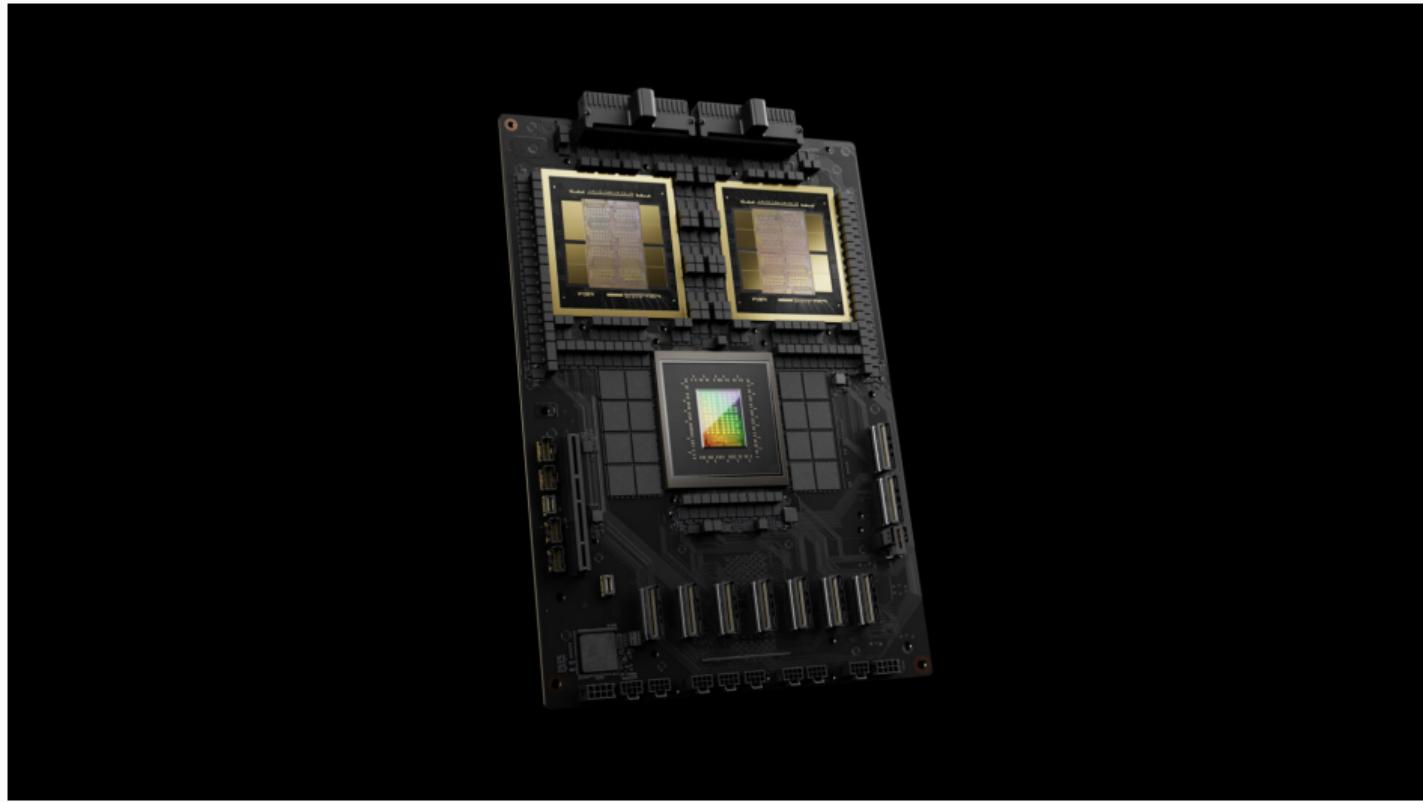


Figure 12: NVIDIA GB200 “superchip” featuring 2 GPUs, 1 CPU, 384 GB HBM3e memory, 480 GB LPDDR5X memory and 2700 W TDP



Figure 13: Two NVIDIA GB200 NVL72 racks, each featuring 72 GPUs, 36 CPUs, 13.5 TB HBM3e memory, 17 TB LPDDR5X memory and 120 kW TDP

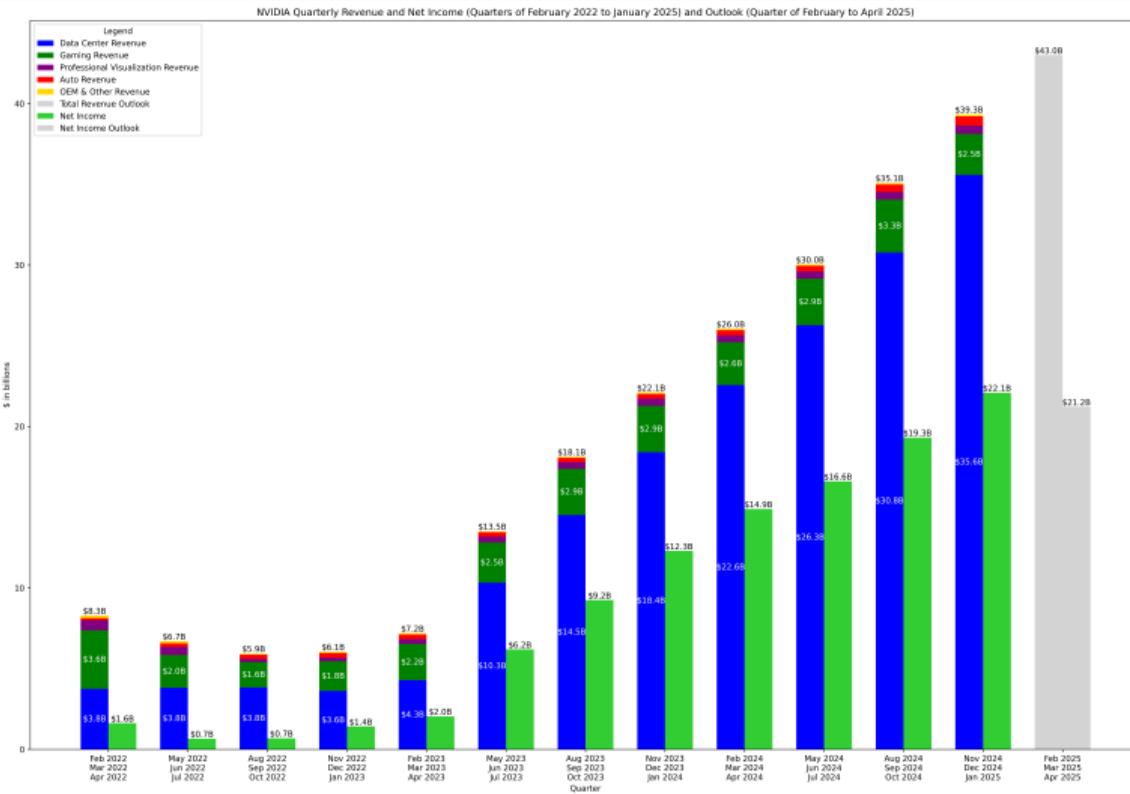
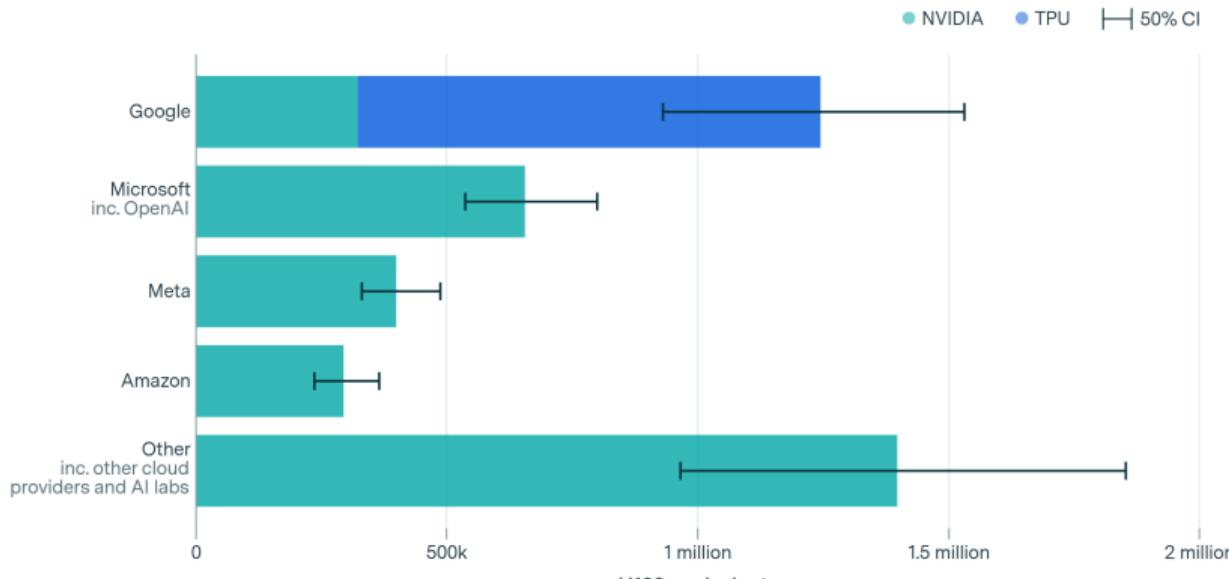


Figure 14: NVIDIA's quarterly data center revenue increased 10-fold from \$3.6B in Q4 2022 to \$35.6B in Q4 2024. In the same time, total profits increased 15-fold from \$1.4B to \$22.1B

AI computing capacity for leading tech companies

Estimates based on 2022 to mid-2024 NVIDIA revenue and TPU shipments.



CC-BY

epoch.ai

Figure 15: Estimated AI computing capacity of tech companies in mid-2024. Four companies are estimated to have purchased about half of NVIDIA's data center GPUs. (Epoch AI, 2024a)

Major Tech Companies Have Developed In-House AI Accelerators

Latest in-house chips for AI training:

- Amazon: AWS Trainium3 (TweakTown, 2024b)
- Google: Trillium (Google, 2024a)
 - Sixth generation of TPU series that began in 2015
 - Fully covers internal use for AI training and inference
 - Allows them to train AI models for roughly a third of the compute cost compared to labs using NVIDIA hardware. (Epoch AI, 2025)
- Microsoft: Azure Maia 100 (Microsoft, 2023b)
 - First chip designed by Microsoft
- Meta: MTIA v2 (Meta, 2024a)

TPU AI accelerators

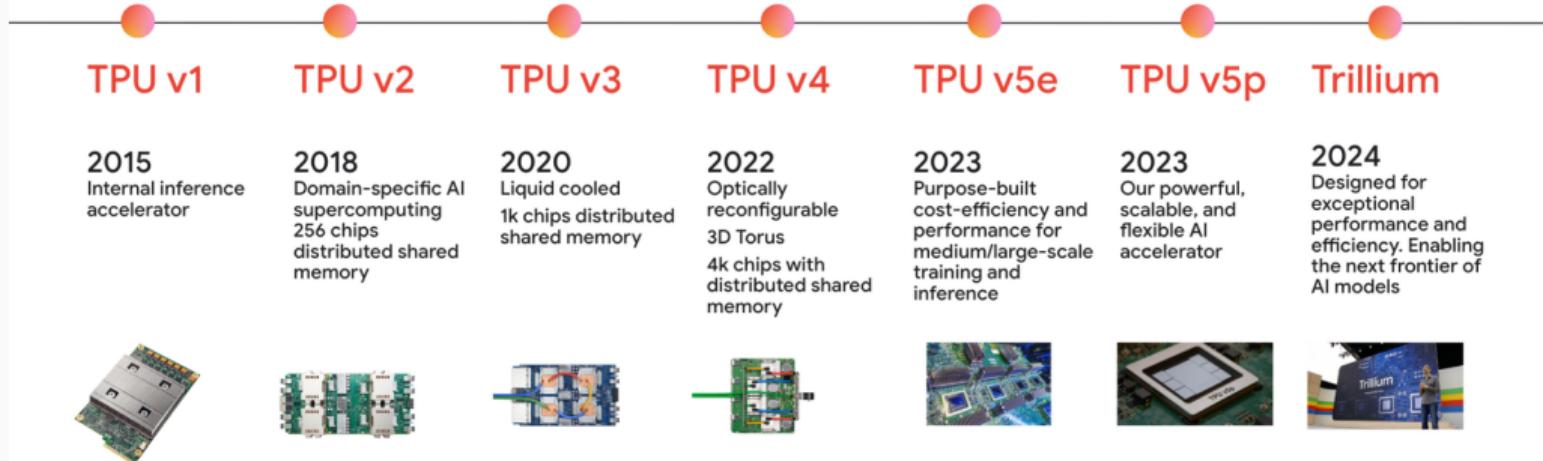


Figure 16: Timeline of Google's TPU AI accelerators from TPU v1 in 2015 to Trillium, the sixth generation, in 2024 (Google, 2024a)

Current Situation of AI Accelerator Supply Chain

- The production of AI accelerators is constrained by advanced chip packaging and high-bandwidth memory (HBM) availability.
 - TSMC's CoWoS process—critical for integrating logic and HBM—is expanding, but new packaging facilities require significant capital investment and highly specialized equipment.
 - HBM supply remains tight, with current volumes nearly fully committed until 2026.
- Projections indicate an annual growth in GPU production between 30% and 100%, yet estimates vary considerably.
- Competition among AI labs means that individual actors may only access a limited fraction of the total chip supply.
 - Projections suggest manufacturing capacity sufficient for roughly 100 million H100-equivalent GPUs, with estimates ranging from 20 to 400 million, to be dedicated to training. Considering all factors, this would be **roughly 5 times the estimated minimum requirements**.

Chips: In a Nutshell

- Hardware used for AI training has evolved from graphics cards to specialized AI accelerators, with future directions including ASICs, in-memory computing, and experimental technologies like photonic and neuromorphic computing.
- NVIDIA dominates the AI hardware market, experiencing a 10x increase in data center revenue over two years while maintaining huge profit margins.
- Major tech companies including Google, Amazon, Meta and Microsoft have developed their own AI accelerators to reduce reliance on NVIDIA and lower costs, with Google's TPU series showing particular maturity.
- The computing capacity for AI is highly concentrated, with these four companies estimated to have purchased about half of all NVIDIA data center GPUs.
- While manufacturing constraints exist, particularly in chip packaging and high-bandwidth memory, current projections (with high uncertainty) suggest the chip industry can produce enough GPUs for AI training through 2030.

Power

What the Hell Is a Gigawatt?

- 2 to 4 MW: Capacity of a wind turbine
- 10 to 30 MW: Mid-sized data center (non-AI)
- 50 MW: Large factory
- 5 to 300 MW: Small modular reactor (SMR)
- 1 GW: Nuclear power plant with single mid-sized reactor
- 2.1 GW: Peak load of Berlin grid (3.8 million residents) (Stromnetz Berlin, 2024)
- 7.6 GW: Peak load of Switzerland (9 million residents) in 2024 (Swissgrid, 2025)

How Much Power Could a Training Run in 2030 Require?

- Epoch estimates a 24x increase in energy efficiency by 2030
 - 4x through hardware efficiency improvements
 - 2x through switching from FP16 to FP8 for training
 - 3x through increasing the lengths of training runs
- To sustain the current trajectory, interconnected data centers with a capacity of approximately **6 gigawatts** working on one training run would be required
- Problem: Gigawatt-scale data centers can require significant upgrades to the local power infrastructure which can take many years
- Solution 1: Bringing data center and power generation as close together as possible, either by using existing or building new power generation
- Solution 2: Training on geographically distributed data centers
 - It's currently unclear how well this works due to bandwidth constraints

How Much Power Might AI Require in Total?

- A report released by the U.S. Department of Energy in December 2024 finds that data centers consumed about 4.4% of total U.S. electricity in 2023 and are expected to consume approximately 6.7 to 12% of total U.S. electricity by 2028. (DOE, 2024)
- The PJM Interconnection, the largest regional transmission organization (RTO) in the U.S., forecasted in January 2025 that its peak load would grow 36% in the next 10 years. As recently as 2022, they forecasted no peak load growth over the same time frame. (PJM, 2025)
- Federal, state and local governments, utility companies, energy companies, RTOs and technology companies are well aware of this challenge and taking action.
- While the U.S. is accustomed to flat electricity demand in the past few decades, a ramp-up of this size appears manageable with these actions.

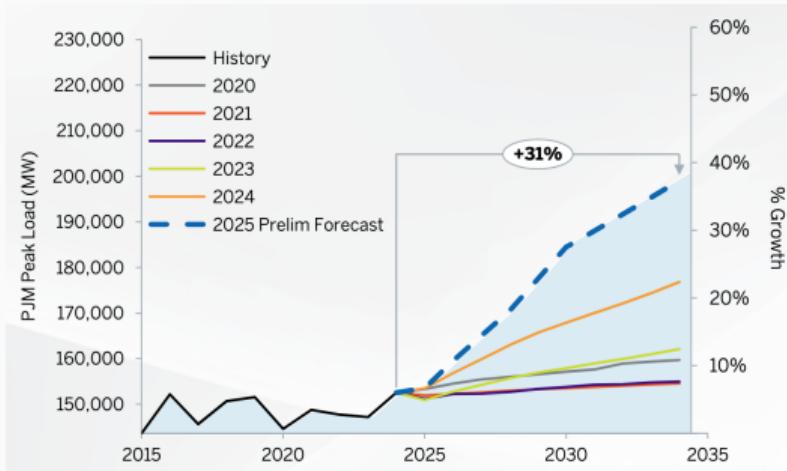
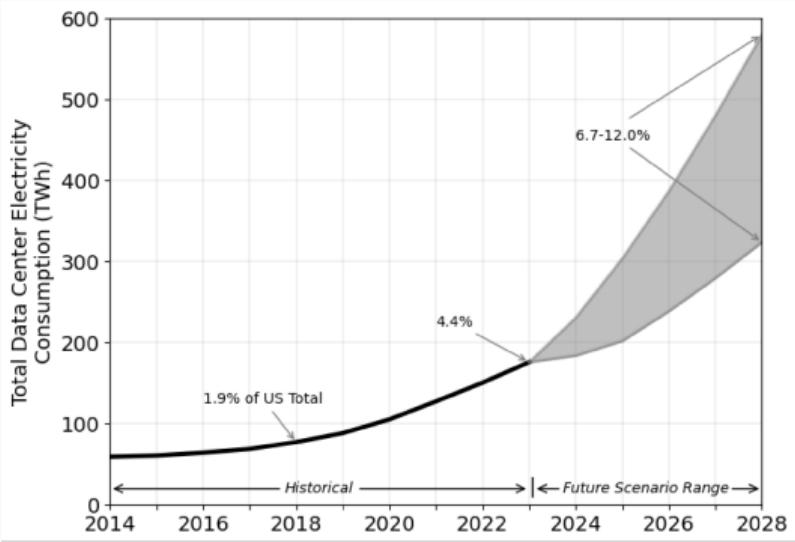


Figure 17: Expected total U.S. data center electricity use from 2023 to 2028 (left) and forecasted peak load growth in the PJM Interconnection from 2024 to 2034 (right).

How Big Are The Biggest Data Centers Currently?

- For context: GPT-4 was trained on around 25,000 NVIDIA A100 (roughly 20 MW) in 2022 (Epoch AI, 2025)
- Scaling beyond this presents novel challenges, including large-scale construction, networking, hardware failures, bandwidth constraints, cooling and power
- Leading tech companies recently finished data centers with around 100,000 NVIDIA H100 (roughly 150 MW), offering roughly 12x compute (Patel, 2025)
- Meta is training their Llama 4 models on a cluster with more than 100,000 NVIDIA H100s (Wired, 2024)
- No models trained on these were released in 2024 due to the time needed for experimentation, training, post-training and safety evaluations, but are being released during 2025 (GPT-5, Claude 4, Llama 4, Grok 3).

How Big Are the Data Centers Planned for 2025?

- Microsoft is said to be building 5 interconnected data centers for AI training with 100,000 NVIDIA Blackwell GPUs each for a total of 500,000 GPUs using 1 GW while other tech giants have similar plans. (Patel, 2024a)

Building New Nuclear Reactors for Data Centers?

- Amazon
 - Looked for a “Principal Nuclear Engineer” to “build internal and external nuclear product and fuel strategy roadmaps [...] for AWS data centers” (Amazon, 2024b)
 - Signed three agreements to support “nuclear energy projects—including enabling the construction of several new Small Modular Reactors (SMRs)” (Amazon, 2024a)
- Google
 - Signed an agreement to “purchase nuclear energy from multiple small modular reactors (SMR) to be developed by Kairos Power” (Google, 2024b)
- Microsoft
 - Looked for a “Principal Program Manager Nuclear Technology” to explore SMRs to power data centers (Microsoft, 2023a) and hired directors for “Nuclear Technologies” and “Nuclear Development Acceleration” (Data Center Dynamics, 2024)
- Meta
 - Request for proposals to “identify nuclear energy developers [...] targeting 1-4 gigawatts (GW) of new nuclear generation capacity in the U.S.” (Meta, 2024b)

Gigawatt-Scale Data Centers Next to Nuclear Power Plants?

- Constellation Energy and Vistra are in discussions with large tech companies to power data centers with up to several gigawatts of capacity co-located with existing nuclear power plants through behind-the-meter agreements. (Bloomberg, 2024; Data Center Frontier, 2024; Wall Street Journal, 2024)
- Under such an agreement, a data center would be build next to an existing nuclear power plant and directly connected with it, requiring a significant portion or even all of the output of the plant.

What Are the Biggest Data Center Currently Planned for the Future?

- The biggest known plans for data centers are from Amazon, Microsoft, OpenAI (through Stargate) and Meta, each having roughly 1 to 2 gigawatts of capacity.
- These projects will either use existing large substations and power plants or build new dedicated power plants or use a combination of these two approaches.
- The hardware costs of a 2 GW project alone (when using NVIDIA hardware) will be on the order of \$50 billion.
- Each of these projects will rank among the most expensive factories ever built and among the most expensive megaprojects in U.S. history.

X.AI Builds Gigawatt-Scale Data Center in Memphis, Tennessee

- X.AI is building a gigawatt-scale data center inside a closed Electrolux factory in Memphis, Tennessee
- Located next to a 1.1 GW natural gas plant and a tapped-off natural gas line.
- Improvised using mobile gas generators, battery packs and chillers parked outside the building.
- Currently being built out to 200,000 NVIDIA Hopper GPUs (roughly 300 MW).
- Build-out to 1 million NVIDIA Blackwell GPUs (roughly 2 GW) and construction of a natural gas plant planned.

Amazon's Nuclear-Powered Data Center Plan

- In March 2024, AWS acquired the Cumulus campus intended to have 960 MW capacity once finished, right besides the 2.5 GW Susquehanna nuclear power plant in Pennsylvania.
- The site was initially intended for crypto mining companies and will be expanded in 120 MW increments into a gigawatt-scale campus.
- The project faces pushback from utility companies and consumers and regulatory challenges, including the FERC blocking the proposed expansion in a first decision.
 - Talen Energy is challenging the decision. (Data Center Dynamics, 2025a)



Figure 18: Rendering from Talen Energy of how the Cumulus data center campus next to the Susquehanna nuclear power plant could look like in the future. Amazon's plans might differ.

Microsoft's Gigawatt-Scale Data Center in Wisconsin

- Microsoft is constructing data centers in Mount Pleasant, Wisconsin that are expected to have a total capacity of 1.5 gigawatts by mid-2027. (Patel, 2024b)
- Announced in May 2024 by Microsoft President Brad Smith and U.S. President Joe Biden. (Microsoft, 2024)

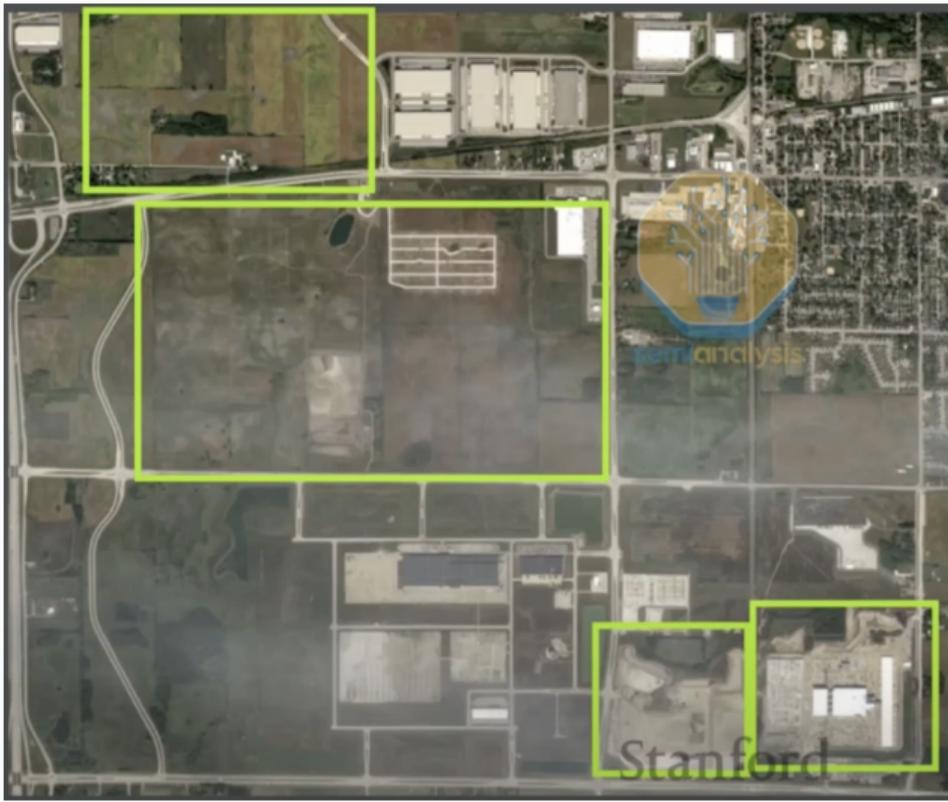


Figure 19: Satellite footage of Microsoft's data center sites in Mount Pleasant, Wisconsin with construction sites and expansion plans highlighted. (Patel, 2024b)

OpenAI's Project Stargate

- Initially leaked in March 2024 as a \$100 billion supercomputer project planned by OpenAI and Microsoft featuring millions of chips and requiring up to 5 gigawatts of power (The Information, 2024)
- Officially announced at the beginning of U.S. President Trump's first press conference back in office (White House, 2025e)
 - Joint venture created by OpenAI, SoftBank, Oracle and MGX
 - No funding from the U.S. government or Microsoft
 - Plans to invest up to \$500 billion in AI infrastructure for OpenAI, but so far only \$100 billion have been committed and the rest is questionable
- First data center is under construction in Abilene, Texas and is reportedly 2.2 GW with estimated hardware costs of roughly \$50 billion (Patel, 2025)
 - Located at a site for a large data center project planned since 2021 and close to two natural gas pipelines.
 - Two substations built and eight data hall buildings under construction currently.



Figure 20: Satellite footage from February 27, 2025 of the construction of Stargate Site 1 in Abilene, Texas. Clearing of additional 2 km^2 of land began in January 2025. (Copernicus, 2025)

Meta's Plans for Gigawatt-Scale Data Centers

- Scrapped Plan for Data Center at Nuclear Power Plant (Financial Times, 2024)
 - Multiple complications including the discovery of a rare bee species at the location
- 2+ GW Data Center in Richland Parish, Louisiana (Zuckerberg, 2024)
 - Three natural gas turbines with a combined capacity of 2.26 GW, six substations, nearly 100 miles of 500kV transmission lines, eight new 230kV transmission lines and 1.5 GW of solar and storage resources will be built. (Entergy, 2024)
 - Announced by Louisiana Governor Landry in December 2024, who said it "is expected to be the largest private capital investment announcement in the state's history" (Governor of Louisiana, 2024)
- Possibly a 1+ GW Data Center in West Feliciana, Louisiana (Advocate, 2025)
 - Hut 8 is building an AI data center for an undisclosed tenant
 - 300 MW capacity in first phases and possibly up to and beyond 1 GW in future phases with possible construction of a power plant



Figure 21: Rendering of Meta's planned data center in Richland Parish, Louisiana, featuring a 1 km^2 main complex, additional buildings and cooling ponds on a 6 km^2 site (Meta, 2024c)



Figure 22: Satellite footage from March 5, 2025 of the construction of Meta's data center in Richland Parish, Louisiana. Clearing of 2.5 km² of land for the construction of two natural gas plants began in November 2024, visible in the top center. Clearing for the main complex and cooling ponds began in February 2025, visible in the center. (Copernicus, 2025)

Power: In a Nutshell

- Frontier AI models in 2030 could require interconnected data centers with a capacity of approximately 6 gigawatts for a single training run, posing unprecedented power infrastructure challenges.
- The U.S. Department of Energy and the PJM Interconnection project a significant, yet likely manageable, increase in total electricity demand, a major shift from previously flat electricity demand forecasts.
- Major tech companies are pursuing gigawatt-scale data centers connected to existing power plants, building new power plants, and exploring nuclear energy options through agreements with energy companies.
- The largest planned data centers will each rank among the most expensive factories ever built and among the most expensive megaprojects in U.S. history.

Geopolitics

U.S. Leads AI Investment; U.S. and China Lead AI Research

- 61% of AI patents granted in 2022 were from China, while 21% were from the U.S. and 2% were from the EU/UK.
- 73% of foundation models released in 2023 were from the U.S., while 13% were from China and 10% were from the EU/UK.
- 39% of GitHub stars on AI projects until 2023 were on projects from the U.S., while 17% were from the EU/UK, 8% from China and 7% from India.
- 70% of global private investment in AI in 2023 occurred in the U.S., while 11% occurred in the EU/UK and 8% in China.
 - In generative AI, 89% occurred in the U.S., 3% in the EU/UK and 3% in China.
- 897 AI companies were funded in the U.S. in 2023, while 368 were funded in the EU/UK and 122 were funded in China.
- All data from the 2024 AI Index Report (Stanford University, 2024)

Regulatory Restrictions		Unreleased	No Licence Required	NAC License Required	Presumption of Denial
NVIDIA GPU Architecture	Model	Pre-Controls	October 2022 Controls ²	October 2023 Controls ^{3,4}	AI Diffusion Rules ⁵
Blackwell	Announced		7-Oct-22	17-Oct-23	13-Jan-25
	Effective ¹		21-Oct-22	17-Nov-23	15-May-25
Hopper	B200				
	B100				
Lovelace	H100				
	H200				
Ampere	H800				
	H20				
Consumer GPUs	L40S				
	L4				
Ampere	L40				
	L20				
Ampere	L2				
	A100				
Consumer GPUs	A800				
	A40				
Consumer GPUs	A30				
	RTX 6000 Ada				
Consumer GPUs	RTX 4090				
	RTX 4090D				
Consumer GPUs	RTX 3090				

Figure 23: U.S. chip export restrictions to China and affected NVIDIA GPUs. Export of the H100 was banned in October 2022, leading NVIDIA to develop the H800. Export of the H800 was banned in October 2023, leading NVIDIA to develop the H20. (Artificial Analysis, 2025)

Key Actions and Statements of the Biden Administration on AI

- AI National Security Memorandum in October 2024 (White House, 2024)
 - Declares the U.S. “must lead the world’s development of safe, secure, and trustworthy AI”, instructing federal agencies to “streamline permitting, approvals, and incentives for the construction of AI-enabling infrastructure” including “clean energy generation, power transmission lines, and high-capacity fiber data links”
- Order on AI Export Restrictions on January 13, 2025 (White House, 2025a)
 - Each country except “18 key allies and partners” is subjected to a new import compute cap equivalent to around 20,000 NVIDIA Blackwell GPUs until 2027.
 - Small purchases don’t count against that total and companies can, under strict requirements, apply for increased caps.
- Order on Advancing AI Infrastructure on January 14, 2025 (White House, 2025b)
 - “Lease federal sites [...] to host gigawatt-scale AI data centers.”
- Farewell Address on January 15, 2025
 - “Meanwhile, artificial intelligence is the most consequential technology of our time — perhaps of all time.” (White House, 2025c)

Key Actions and Statements of the Trump Administration on AI

- Senate Hearing of Interior Secretary Burgum (U.S. Senate, 2025)
 - “So we’re competing against someone who’s going to create more electricity, produce more AI, and this could be how we lose the Cold War with [China]”
- National Energy Emergency Declaration on Day 1 (White House, 2025d)
 - Declared the first ever national energy emergency for several reasons including “high demand for energy and natural resources to power the next generation of technology”
- Announcement of “Project Stargate” on Day 2 (White House, 2025e)
 - “I’m going to help a lot through emergency declarations because we have an emergency. We have to get this stuff built. So they have to produce a lot of electricity, and we’ll make it possible for them to get that production done very easily at their own plants”
- Speech of Vice President Vance at the AI Action Summit in Paris (UCSB, 2025)
 - “Now, at this moment, we face the extraordinary prospect of a new industrial revolution, one on par with the invention of the steam engine or Bessemer steel.”

Europe, and Particularly France, Step Up On AI

- In February 2025, the EU launched InvestAI, a public-private partnership “to mobilise €200 billion for investment in AI, including a new European fund of €20 billion for AI gigafactories” including “four future AI gigafactories across the EU” with “around 100 000 last-generation AI chips”. (European Commission, 2025)
- French President Macron announced in February 2025 that France has secured €109 billion in “private French and foreign investments” in AI “for the coming years”. (Le Monde, 2025)
- A few days earlier after a meeting with UAE President bin Zayed, he announced €30-50 billion investments from French and Emirati investors into a 1 GW AI data center in France. (Data Center Dynamics, 2025c)
- Mistral AI, the leading AI startup in Europe, has raised almost €1 billion so far and is headquartered in Paris.

Geopolitics: In a Nutshell

- The U.S. leads the world in AI investment, while both the U.S. and China lead in AI research.
- The U.S. has implemented increasingly stringent chip export restrictions to China (and now almost all countries of the world), limiting access to AI accelerators.
- Both the Biden and Trump administrations have elevated AI to a matter of national security and economic competitiveness, with both taking measures to streamline construction of gigawatt-scale data centers and Trump declaring the first ever national energy emergency partly to power AI.
- Europe, led by France, is mounting a significant response with multi-billion euro investments to build data centers, though still lagging behind the U.S.

References

Advocate. (2025, January 6). **\$2.5 Billion West Feliciana Data Center Approved With 'Aggressive Timeline To Build'**. Retrieved January 22, 2025, from

https://www.theadvocate.com/baton_rouge/news/business/25-billion-west-feliciana-data-center-approved/article_d6dce702-cc45-11ef-8ff0-1fd2550a45e5.html

Amazon. (2024a, October 16). **Amazon Signs Agreements For Innovative Nuclear Energy Projects To Address Growing Energy Demands**. Retrieved March 2, 2025, from <https://www.aboutamazon.com/news/sustainability/amazon-nuclear-small-modular-reactor-net-carbon-zero>

References ii

- Amazon. (2024b, November 6). **Principal Nuclear Engineer, Datacenter Engineering, Power Generation Solutions - Job ID: 2741394 — Amazon.jobs.** Retrieved March 2, 2025, from <https://web.archive.org/web/20241106092116/https://www.amazon.jobs/en/jobs/2741394/principal-nuclear-engineer-datacenter-engineering-power-generation-solutions>
- Artificial Analysis. (2024, December 17). **Artificial analysis AI review: 2024 highlights.** <https://artificialanalysis.ai/downloads/ai-review/2024/Artificial-Analysis-AI-Review-2024-Highlights.pdf>
- Artificial Analysis. (2025). **State of AI: China.** Retrieved March 8, 2025, from <https://artificialanalysis.ai/downloads/china-report/2025/Artificial-Analysis-State-of-AI-China-Q1-2025.pdf>

References iii

- Bloomberg. (2024, September 23). **Constellation CEO Says US Should Copy China To Meet AI Power Use.** <https://www.bloomberg.com/news/articles/2024-09-23/constellation-ceo-sees-ai-as-critical-to-us-national-security>
- Cerebras. (2024, March 11). **Cerebras Systems Unveils World's Fastest AI Chip With Whopping 4 Trillion Transistors [Cerebras].** Retrieved February 8, 2025, from <https://cerebras.ai/press-release/cerebras-announces-third-generation-wafer-scale-engine/>
- Chollet, F. (2024, December 20). **OpenAI O3 Breakthrough High Score On ARC-AGI-Pub [ARC prize].** Retrieved February 8, 2025, from <https://arcprize.org/blog/oai-o3-pub-breakthrough>
- Chu, T., Zhai, Y., Yang, J., Tong, S., Xie, S., Schuurmans, D., Le, Q. V., Levine, S., & Ma, Y. (2025). **Sft memorizes, rl generalizes: A comparative study of foundation model post-training.** <https://arxiv.org/abs/2501.17161>

- CNBC. (2025, February 6). **Amazon Plans To Spend \$100 Billion This Year To Capture ‘Once In A Lifetime Opportunity’ In AI [CNBC]**. Retrieved February 23, 2025, from <https://www.cnbc.com/2025/02/06/amazon-expects-to-spend-100-billion-on-capital-expenditures-in-2025.html>
- Copernicus. (2025, January 21). **Copernicus Browser [Copernicus browser]**. Retrieved January 23, 2025, from <https://browser.dataspace.copernicus.eu/>
- Data Center Dynamics. (2024, January 22). **Microsoft Hires Archie Manoharan As Director Of Nuclear Technologies, Joins From Micro Modular Reactor Firm.** Retrieved March 2, 2025, from <https://www.datacenterdynamics.com/en/news/microsoft-hires-archie-manoharan-as-director-of-nuclear-technologies-joins-from-micro-modular-reactor-firm/>

References v

Data Center Dynamics. (2025a, January 25). **Susquehanna Nuclear Challenges FERC Rejection Of Connection Agreement To Power AWS Data Center.**

Retrieved February 9, 2025, from

<https://www.datacenterdynamics.com/en/news/susquehanna-nuclear-challenges-ferc-rejection-of-connection-agreement-to-power-aws-data-center/>

Data Center Dynamics. (2025b, February 5). **Google Expects 2025 Capex To Surge To \$75Bn On AI Data Center Buildout.** Retrieved February 23, 2025, from

<https://www.datacenterdynamics.com/en/news/google-expects-2025-capex-to-surge-to-75bn-on-ai-data-center-buildout/>

Data Center Dynamics. (2025c, February 7). **France And UAE To Invest Billions Into 1GW European AI Data Center.** Retrieved February 9, 2025, from

<https://www.datacenterdynamics.com/en/news/france-and-uae-to-invest-billions-into-1gw-european-ai-data-center/>

Data Center Frontier. (2024, April 29). **The Gigawatt Data Center Campus Is Coming [Data center frontier].** Retrieved March 2, 2025, from

<https://www.datacenterfrontier.com/hyperscale/article/55021675/the-gigawatt-data-center-campus-is-coming>

DOE. (2024, December 20). **DOE Releases New Report Evaluating Increase In Electricity Demand From Data Centers [Energy.gov].** Retrieved March 2, 2025, from
<https://www.energy.gov/articles/doe-releases-new-report-evaluating-increase-electricity-demand-data-centers>

- Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., Li, C., Xiao, C., Du, C., Liao, C., Tang, C., Wang, C., Zhang, D., Yuan, E., Lu, E., Tang, F., Sung, F., Wei, G., Lai, G., Guo, H., ... Yang, Z. (2025). **Kimi k1.5: Scaling reinforcement learning with llms.** <https://arxiv.org/abs/2501.12599>
- Entergy. (2024, December 5). **Entergy Louisiana To Power Meta's Data Center In Richland Parish.** Retrieved January 22, 2025, from <https://www.entropynewsroom.com/news/entropy-louisiana-power-meta-s-data-center-in-richland-parish/>
- Epoch AI. (2024a, October 23). **Data On Machine Learning Hardware.** Retrieved January 22, 2025, from <https://epoch.ai/data/machine-learning-hardware>

- Epoch AI. (2024b, October 23). **Performance Improves 12X When Switching From FP32 To Tensor-INT8 [Epoch AI]**. Retrieved March 9, 2025, from <https://epoch.ai/data-insights/hardware-performance-trend>
- Epoch AI. (2025, January 31). **Data On Notable AI Models [Epoch AI]**. Retrieved February 1, 2025, from <https://epoch.ai/data/notable-ai-models>
- Erdil, E. (2025, January 31). **What Went Into Training DeepSeek-R1? [Epoch AI]**. Retrieved February 7, 2025, from <https://epoch.ai/gradient-updates/what-went-into-training-deepseek-r1>
- Etched. (2024, June 25). **Etched Is Making The Biggest Bet In AI**. Retrieved February 8, 2025, from <https://www.etched.com/announcing-etched>

- European Commission. (2025, February 11). **EU Launches InvestAI Initiative To Mobilise €200 Billion Of Investment In Artificial Intelligence [European commission - european commission]**. Retrieved March 1, 2025, from
https://ec.europa.eu/commission/presscorner/detail/en/https://%5C/%5C/ec.europa.eu%5C/commission%5C/presscorner%5C/detail%5C/en%5C/ip_25_467
- Feng, Y., Dohmatob, E., Yang, P., Charton, F., & Kempe, J. (2024). **Beyond model collapse: Scaling up with synthesized data requires verification.**
<https://arxiv.org/abs/2406.07515>
- Financial Times. (2024, November 4). **Meta's Plan For Nuclear-powered AI Data Centre Thwarted By Rare Bees.** Retrieved January 22, 2025, from
<https://www.ft.com/content/ed602e09-6c40-4979-aff9-7453ee28406a>

References x

- Google. (2024a, July 31). **TPU Transformation: A Look Back At 10 Years Of Our AI-specialized Chips [Google cloud blog]**. Retrieved February 9, 2025, from <https://cloud.google.com/transform/ai-specialized-chips-tpu-history-gen-ai>
- Google. (2024b, October 14). **New Nuclear Clean Energy Agreement With Kairos Power [Google]**. Retrieved March 2, 2025, from <https://blog.google/outreach-initiatives/sustainability/google-kairos-power-nuclear-energy-agreement/>
- Governor of Louisiana. (2024, December 4). **Landry Announces Meta Selects North Louisiana As Site Of \$10 Billion Artificial Intelligence Optimized Data Center — Office Of Governor Jeff Landry**. Retrieved January 22, 2025, from <https://gov.louisiana.gov/news/4697>

- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., ... Ma, Z. (2024, November 23). **The llama 3 herd of models.**
<https://doi.org/10.48550/arXiv.2407.21783>
- Groq. (2024, June 28). **GroqRack - Groq Is Fast AI Inference.** Retrieved February 8, 2025, from <https://groq.com/groqrack/>
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., ... Zhang, Z. (2025). **Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.** <https://arxiv.org/abs/2501.12948>

- Hao, X., Shen, K., & Li, C. (2025). **Maga: Massive genre-audience reformulation to pretraining corpus expansion.** <https://arxiv.org/abs/2502.04235>
- Ho, A., Besiroglu, T., Erdil, E., Owen, D., Rahman, R., Guo, Z. C., Atkinson, D., Thompson, N., & Sevilla, J. (2024, March 9). **Algorithmic progress in language models.** <https://doi.org/10.48550/arXiv.2403.05812>
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A., Hennigan, T., Noland, E., Millican, K., Driessche, G. v. d., Damoc, B., Guy, A., Osindero, S., Simonyan, K., Elsen, E., ... Sifre, L. (2022, March 29). **Training compute-optimal large language models.** <https://doi.org/10.48550/arXiv.2203.15556>

References xiii

- Hu, M., Zhao, P., Xu, C., Sun, Q., Lou, J., Lin, Q., Luo, P., & Rajmohan, S. (2024). **Agentgen: Enhancing planning abilities for large language model based agent via environment and task generation.** <https://arxiv.org/abs/2408.00764>
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., & Amodei, D. (2020, January 23). **Scaling laws for neural language models.** <https://doi.org/10.48550/arXiv.2001.08361>
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). **Imagenet classification with deep convolutional neural networks.** In F. Pereira, C. Burges, L. Bottou, & K. Weinberger (Eds.), *Advances in neural information processing systems* (Vol. 25). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf

- Le Monde. (2025). **Intelligence artificielle : Emmanuel Macron annonce des investissements en France de '109 milliards d'euros dans les prochaines années**. Retrieved February 9, 2025, from
https://www.lemonde.fr/pixels/article/2025/02/09/intelligence-artificielle-emmanuel-macron-annonce-des-investissements-en-france-de-109-milliards-d-euros-dans-les-prochaines-annees_6539115_4408996.html
- LeGris, S., Vong, W. K., Lake, B. M., & Gureckis, T. M. (2024). **H-arc: A robust estimate of human performance on the abstraction and reasoning corpus benchmark.** <https://arxiv.org/abs/2409.01374>
- Luo, M., Tan, S., Wong, J., Shi, X., Tang, W. Y., Roongta, M., Cai, C., Luo, J., Zhang, T., Li, L. E., Popa, R. A., & Stoica, I. (2025). **Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl [Notion Blog].**

- Meta. (2024a, April 10). **Our Next-generation Meta Training And Inference Accelerator.** Retrieved February 9, 2025, from
<https://ai.meta.com/blog/next-generation-meta-training-inference-accelerator-AI-MTIA/>
- Meta. (2024b, December 3). **Accelerating The Next Wave Of Nuclear To Power AI Innovation [Meta sustainability].** Retrieved March 2, 2025, from
<https://sustainability.atmeta.com/blog/2024/12/03/accelerating-the-next-wave-of-nuclear-to-power-ai-innovation/>
- Meta. (2024c, December 4). **Richland Parish Data Center.** Retrieved January 22, 2025, from <https://www.facebook.com/RichlandParishDataCenter>

References xvi

- Meta. (2025, January 29). **Meta Platforms (META) Q4 2024 Earnings Call Transcript [The motley fool]**. Retrieved February 1, 2025, from <https://www.fool.com/earnings/call-transcripts/2025/01/29/meta-platforms-meta-q4-2024-earnings-call-transcri/>
- Mialon, G., Fourrier, C., Swift, C., Wolf, T., LeCun, Y., & Scialom, T. (2023). **Gaia: A benchmark for general ai assistants.**
- Microsoft. (2023a, September 22). **Principal Program Manager Nuclear Technology — Microsoft Careers.** Retrieved March 2, 2025, from <https://web.archive.org/web/20230922194154/https://jobs.careers.microsoft.com/global/en/job/1627555/Principal-Program-Manager-Nuclear-Technology>

Microsoft. (2023b, November 15). **With A Systems Approach To Chips, Microsoft Aims To Tailor Everything ‘From Silicon To Service’ To Meet AI Demand [Microsoft]**. Retrieved February 9, 2025, from

<https://news.microsoft.com/source/features/ai/in-house-chips-silicon-to-service-to-meet-ai-demand/>

Microsoft. (2024, May 8). **Microsoft Announces \$3.3 Billion Investment In Wisconsin To Spur Artificial Intelligence Innovation And Economic Growth.**

Retrieved February 17, 2025, from

<https://news.microsoft.com/2024/05/08/microsoft-announces-3-3-billion-investment-in-wisconsin-to-spur-artificial-intelligence-innovation-and-economic-growth/>

- Microsoft. (2025, January 3). **The Golden Opportunity For American AI [Microsoft on the issues]**. Retrieved January 28, 2025, from <https://blogs.microsoft.com/on-the-issues/2025/01/03/the-golden-opportunity-for-american-ai/>
- NVIDIA. (2025, February 9). **NVIDIA GB200 NVL72 GPU – Optimized For AI And Data Centers [NVIDIA]**. Retrieved February 8, 2025, from <https://www.nvidia.com/en-us/data-center/gb200-nvl72/>
- OpenAI. (2024a, September 12). **Learning To Reason With LLMs**. Retrieved February 23, 2025, from <https://openai.com/index/learning-to-reason-with-langs/>
- OpenAI. (2024b, December 21). **OpenAI o1 system card**.
<https://doi.org/10.48550/arXiv.2412.16720>
- OpenAI. (2025, February 2). **Introducing Deep Research**. Retrieved February 8, 2025, from <https://openai.com/index/introducing-deep-research/>

References xix

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022, March 4). **Training language models to follow instructions with human feedback.** <https://doi.org/10.48550/arXiv.2203.02155>

Pan, J., Wang, X., Neubig, G., Jaitly, N., Ji, H., Suhr, A., & Zhang, Y. (2024). **Training software engineering agents and verifiers with swe-gym.**

<https://arxiv.org/abs/2412.21139>

Pan, J., Zhang, J., Wang, X., Yuan, L., Peng, H., & Suhr, A. (2025). **Tinyzero.**

Patel, D. (2024a, October 9). **\$10B OpenAI & Microsoft Cluster In 2025 – Dylan Patel & @Asianometry.** Retrieved February 23, 2025, from

<https://www.youtube.com/watch?v=J7yvSkIOZLw>

References xx

- Patel, D. (2024b, November 12). **Dylan Patel - Inference Math, Simulation, And AI Megaclusters - Stanford CS 229S - Autumn 2024.** Retrieved February 17, 2025, from <https://www.youtube.com/watch?v=hobvps-H38o>
- Patel, D. (2025, February 3). **Transcript For DeepSeek, China, OpenAI, NVIDIA, XAI, TSMC, Stargate, And AI Megaclusters — Lex Fridman Podcast #459 [Lex fridman].** Retrieved February 6, 2025, from [https://lexfridman.com深深寻求中国,OpenAI,NVIDIA,XAI,TSMC,Stargate,And AI Megaclusters — Lex Fridman Podcast #459 \[Lex fridman\].](https://lexfridman.com深深寻求中国,OpenAI,NVIDIA,XAI,TSMC,Stargate,And AI Megaclusters — Lex Fridman Podcast #459 [Lex fridman].)
- PJM. (2025, January 24). **PJM long-term load forecast report.** <https://www.pjm.com/-/media/DotCom/library/reports-notices/load-forecast/2025-load-report.pdf>
- Planetary Society. (2019, June 14). **How Much Did The Apollo Program Cost? [The planetary society].** Retrieved March 8, 2025, from <https://www.planetary.org/space-policy/cost-of-apollo>

References xxi

- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). **Language models are unsupervised multitask learners.** Retrieved February 6, 2025, from <https://www.semanticscholar.org/paper/Language-Models-are-Unsupervised-Multitask-Learners-Radford-Wu/9405cc0d6169988371b2755e573cc28650d14dfe>
- Sevilla, J., Besiroglu, T., Cottier, B., You, J., Roldán, E., Villalobos, P., & Erdil, E. (2024). **Can ai scaling continue through 2030?**
<https://epoch.ai/blog/can-ai-scaling-continue-through-2030>
- Sevilla, J., & Roldán, E. (2024). **Training compute of frontier ai models grows by 4-5x per year [Accessed: 2025-02-08].**
<https://epoch.ai/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year>

References xxii

Shumailov, I., Shumaylov, Z., Zhao, Y., Gal, Y., Papernot, N., & Anderson, R. (2023).

The curse of recursion: Training on generated data makes models forget.

<https://arxiv.org/abs/2305.17493>

Stanford University. (2024). **The 2024 AI Index Report — Stanford HAI.** Retrieved March 8, 2025, from <https://hai.stanford.edu/ai-index/2024-ai-index-report>

Stromnetz Berlin. (2024, July 17). **Faktenblatt - stromnetz berlin GmbH.** <https://www.stromnetz.berlin/files/globalassets/dokumente/presse/faktenblatt-stromnetz-berlin.pdf>

Swissgrid. (2025, January 31). **Netzlast.** Retrieved February 6, 2025, from <https://www.swissgrid.ch/de/home/operation/grid-data/load.html>

References xxiii

The Information. (2024, March 29). **Microsoft And OpenAI Plot \$100 Billion Stargate AI Supercomputer.** Retrieved February 6, 2025, from

<https://www.theinformation.com/articles/microsoft-and-openai-plot-100-billion-stargate-ai-supercomputer>

TweakTown. (2024a, October 13). **NVIDIA Blackwell GPUs For AI Are Effectively 'Sold Out' For The Next 12 Months [TweakTown].** Retrieved March 7, 2025, from

<https://www.tweaktown.com/news/101054/nvidia-blackwell-gpus-for-ai-are-effectively-sold-out-the-next-12-months/index.html>

TweakTown. (2024b, December 4). **Amazon Teases Its Next-gen Trainium3 AI Accelerator Is 4X Faster Than Trainium 3, Drops In 2025 [TweakTown].**

Retrieved February 8, 2025, from <https://www.tweaktown.com/news/102010/amazon-teases-its-next-gen-trainium3-ai-accelerator-is-4x-faster-than-trainium-3-drops-in-2025/index.html>

References xxiv

- UCSB. (2025, February 11). **Remarks By The Vice President At The Artificial Intelligence Action Summit In Paris, France — The American Presidency Project.** Retrieved March 1, 2025, from
<https://www.presidency.ucsb.edu/documents/remarks-the-vice-president-the-artificial-intelligence-action-summit-paris-france>
- Unsloth. (2025, February 6). **Train Your Own R1 Reasoning Model With Unsloth (GRPO).** <https://unsloth.ai/blog/r1-reasoning>
- U.S. Senate. (2025, January 16). **Hearing To Consider The Nomination Of The Honorable Doug Burgum To Be Secretary Of The Interior [U.s. senate committee on energy and natural resources].** Retrieved January 22, 2025, from
<https://www.energy.senate.gov/hearings/2025/1/hearing-to-consider-the-nomination-of-the-honorable-doug-burgum-to-be-secretary-of-the-interior>

References xxv

Valmeekam, K., Stechly, K., & Kambhampati, S. (2024). **Llms still can't plan; can Irlms? a preliminary evaluation of openai's o1 on planbench.**

<https://arxiv.org/abs/2409.13373>

Villalobos, P., & Atkinson, D. (2023). **Trading off compute in training and inference.** Retrieved February 6, 2025, from

<https://epoch.ai/blog/trading-off-compute-in-training-and-inference>

Wall Street Journal. (2024, July 1). **Tech Industry Wants To Lock Up Nuclear Power For AI.** Retrieved March 2, 2025, from <https://www.msn.com/en-us/money/other/tech-industry-wants-to-lock-up-nuclear-power-for-ai/ar-BB1pb4e>

Wei, J., Bosma, M., Zhao, V. Y., Guu, K., Yu, A. W., Lester, B., Du, N., Dai, A. M., & Le, Q. V. (2021, September 3). **Finetuned language models are zero-shot learners [version: 1].** <https://doi.org/10.48550/arXiv.2109.01652>

References xxvi

- White House. (2024, October 24). **Memorandum On Advancing The United States' Leadership In Artificial Intelligence; Harnessing Artificial Intelligence To Fulfill National Security Objectives; And Fostering The Safety, Security, And Trustworthiness Of Artificial Intelligence.** Retrieved January 22, 2025, from <https://bidenwhitehouse.archives.gov/briefing-room/presidential-actions/2024/10/24/memorandum-on-advancing-the-united-states-leadership-in-artificial-intelligence-harnessing-artificial-intelligence-to-fulfill-national-security-objectives-and-fostering-the-safety-security/>
- White House. (2025a, January 13). **FACT SHEET: Ensuring U.S. Security And Economic Strength In The Age Of Artificial Intelligence.** Retrieved January 22, 2025, from <https://bidenwhitehouse.archives.gov/briefing-room/statements-releases/2025/01/13/fact-sheet-ensuring-u-s-security-and-economic-strength-in-the-age-of-artificial-intelligence/>

- White House. (2025b, January 14). **Statement By President Biden On The Executive Order On Advancing U.S. Leadership In Artificial Intelligence Infrastructure.** Retrieved January 22, 2025, from
<https://bidenwhitehouse.archives.gov/briefing-room/statements-releases/2025/01/14/statement-by-president-biden-on-the-executive-order-on-advancing-u-s-leadership-in-artificial-intelligence-infrastructure/>
- White House. (2025c, January 16). **Remarks By President Biden In A Farewell Address To The Nation [The white house].** Retrieved February 8, 2025, from
<https://bidenwhitehouse.archives.gov/briefing-room/speeches-remarks/2025/01/15/remarks-by-president-biden-in-a-farewell-address-to-the-nation/>

References xxviii

- White House. (2025d, January 21). **Declaring A National Energy Emergency.** Retrieved January 22, 2025, from <https://www.whitehouse.gov/presidential-actions/2025/01/declaring-a-national-energy-emergency/>
- White House. (2025e, January 21). **President Trump Gives Remarks Regarding U.S. Infrastructure Investment.** Retrieved January 22, 2025, from <https://www.youtube.com/watch?v=X5gMiDnYEds>
- *Wired. (2024). **Meta's next Llama AI models are training on a GPU cluster 'bigger than anything' else.** *Wired*. Retrieved March 2, 2025, from <https://www.wired.com/story/meta-llama-ai-gpu-training/>
- Xiao, C., Cai, J., Zhao, W., Zeng, G., Lin, B., Zhou, J., Zheng, Z., Han, X., Liu, Z., & Sun, M. (2024). **Densing law of llms.** <https://arxiv.org/abs/2412.04315>

- Xie, T., Zhang, D., Chen, J., Li, X., Zhao, S., Cao, R., Hua, T. J., Cheng, Z., Shin, D., Lei, F., Liu, Y., Xu, Y., Zhou, S., Savarese, S., Xiong, C., Zhong, V., & Yu, T. (2024). **Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments.**
- Xu, F., Hao, Q., Zong, Z., Wang, J., Zhang, Y., Wang, J., Lan, X., Gong, J., Ouyang, T., Meng, F., Shao, C., Yan, Y., Yang, Q., Song, Y., Ren, S., Hu, X., Li, Y., Feng, J., Gao, C., & Li, Y. (2025). **Towards large reasoning models: A survey of reinforced reasoning with large language models.** <https://arxiv.org/abs/2501.09686>

References xxx

Xu, F. F., Song, Y., Li, B., Tang, Y., Jain, K., Bao, M., Wang, Z. Z., Zhou, X., Guo, Z., Cao, M., Yang, M., Lu, H. Y., Martin, A., Su, Z., Maben, L., Mehta, R., Chi, W., Jang, L., Xie, Y., ... Neubig, G. (2024). **Theagentcompany: Benchmarking lilm agents on consequential real world tasks.**

<https://arxiv.org/abs/2412.14161>

Yahoo Finance. (2025, January 28). **Microsoft Corporation (MSFT) Analyst Ratings, Estimates & Forecasts [Yahoo finance].** Retrieved January 28, 2025, from <https://finance.yahoo.com/quote/MSFT/analysis/>

Yeo, E., Tong, Y., Niu, M., Neubig, G., & Yue, X. (2025). **Demystifying long chain-of-thought reasoning in llms.** <https://arxiv.org/abs/2502.03373>

Zuckerberg, M. (2024, December 6). **Last Big AI Update Of The Year.** Retrieved January 22, 2025, from <https://www.instagram.com/zuck/reel/DDPm9gqv2cW/?hl=en>