## **Deep Revolution**

 2024 Nobel prizes in physics (deep learning) & chemistry (Google DeepMind) shook the scientific world, heralding the new era of AI-enabled science
<u>https://www.nobelprize.org/prizes/physics/2024</u>
https://www.nobelprize.org/prizes/chemistry/2024

In January 2025, DeepSeek sent a shock wave to Wall Street, White House, & Silicon Valley

Al stocks plunge as China's DeepSeek sends shock wave through Wall Street

A Chinese AI company called DeepSeek is sending a shock wave through Wall Street. CBS NEWS 1/28/2025

Trump calls DeepSeek a 'wake-up call' for U.S. tech and welcomes China's AI gains FORTUNE 1/28/2025

Meta is reportedly scrambling 'war rooms' of engineers to figure out how DeepSeek's AI is beating everyone else at a fraction of the price FORTUNE 1/27/2025

# **Key Computational Enablers of DeepSeek?**

- DeepSeek is a large language mode (LLM) that outperforms OpenAI's ChatGPT with less computing
- Multi-head Latent Attention guarantees efficient inference through significantly compressing the Key-Value cache into a latent vector, while DeepSeekMoE (Mixture-of-Experts) enables training strong models at an economical cost through sparse computation [https://arxiv.org/abs/2405.04434]
- DeepSeek-V3 pioneers an auxiliary-loss-free strategy for load balancing and sets a multi-token prediction training objective for stronger performance [https://arxiv.org/html/2412.19437v1]
- Reasoning: DeepSeek-R1 directly applies reinforcement learning to the base model, thereby generating a long chain-of-thoughts [https://arxiv.org/abs/2501.12948]
  My expert friend thinks it's their ingenious engineering, not these known & some new methods
- Will brain-like sparse spiking of neurons solve the AI power catastrophe (*cf.* Google's Pathways)?

#### REECE ROGERS GEAR JUL 11, 2824 6:38 AM

Al's Energy Demands Are Out of Control. Welcome to the Internet's Hyper-Consumption Era

Generative artificial intelligence tools, now part of the everyday user experience online, are causing stress on local power grids and mass water evaporation.

### **Final project?**

https://blog.google/technology/ai/introducing-pathways-next-generation-ai-architecture



## **Scaling Analysis Is Important**

- Understanding scaling laws of LLMs is essential for long-term projection https://arxiv.org/abs/2401.02954
- Use the same scaling exponent analysis (log-log plot & linear fit) as in assignment 2, Part I-2!





(b) Optimal model scaling

(c) Optimal data scaling

Nopt(Mopt) & C Dopt & C	Approach	Coeff. <i>a</i> where $N_{\text{opt}}(M_{\text{opt}}) \propto C^a$	Coeff. <i>b</i> where $D_{\text{opt}} \propto C^b$
OpenAI (OpenWebText2)0.730.27Chinchilla (MassiveText)0.490.51	OpenAI (OpenWebText2) Chinchilla (MassiveText)	0.73	0.27 0.51
Ours (Early Data)     0.450     0.550	Ours (Early Data)	0.450	0.550
Ours (Current Data)     0.524     0.476       Ours (OpenWebText2)     0.578     0.422	Ours (Current Data) Ours (OpenWebText2)	0.524 0.578	0.476 0.422

## **Be Creative & Multi-lingual**

• Necessity is the mother of invention: But DeepSeek found ways to reduce memory usage and speed up calculation without significantly sacrificing accuracy. "The team loves turning a hardware challenge into an opportunity for innovation," says Wang. [*MIT Technology Review*]

https://www.technologyreview.com/2025/01/24/1110526/china-deepseek-top-ai-despite-sanctions/

Dive into new languages if necessary: The breakthrough was achieved by implementing tons of fine-grained optimizations and usage of Nvidia's assembly-like PTX (Parallel Thread Execution) programming instead of Nvidia's CUDA for some functions. [tom's Hardware]
cf. PyTorch → C++ → CUDA → Assembly

https://www.tomshardware.com/tech-industry/artificial-intelligence/deepseeks-ai-breakthrough-bypasses-industry-standard-cuda-uses-assembly-like-ptx-programming-instead

• Our C++/OpenMP-target quantum dynamics code (DCMESH) achieves 22% of theoretical peak performance, while PyTorch-based neural-network molecular dynamics code (Allegro-Legato NNQMD) a fraction of %, on the Intel GPU-based Aurora exaflop/s supercomputer