

GSAS@NIDA : Deep Learning Introduction

Assoc.Prof.Thitirat Siriborvornratanakul, Ph.D.

Email: thitirat@as.nida.ac.th
Website: <http://as.nida.ac.th/~thitirat/>

1

OUTLINE

01 Course's Introduction

Schedule, evaluation, books, and FAQs

02 History of Modern DL

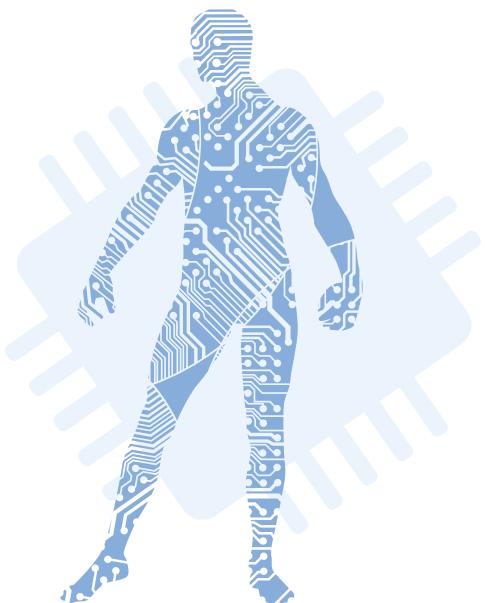
The rapid growth of Deep Learning in the modern era

03 Hardware Acceleration

Prepare necessary hardware to accelerate our deep learning project

04 DL Frameworks

Which frameworks to use in our deep learning project



2

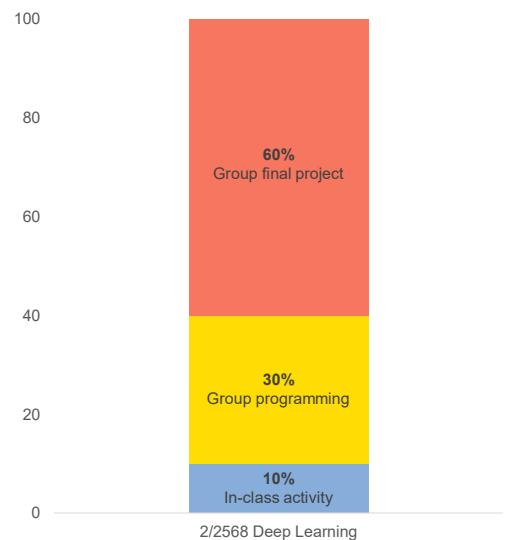
A Course's Introduction

Schedule, evaluation, books, and FAQs

3

Schedule and Evaluation

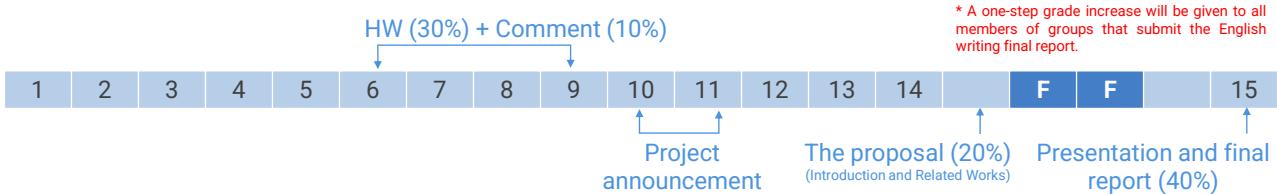
- **Subject:**
 - DADS 7202: Deep Learning (3 credits)
 - CI 7310/7105: Deep Learning (3 credits)
- **Class material and schedule:**
 - MS Teams > General channel > Files (Shared) tab > เอกสารประกอบของคลาส (Class Materials) folder
- **Grading policy:**
 - For students with “Audit” registration, **fulfilling the 80% class attendance does not guarantee the “pass” grade.**
- **Important note!!!**
 - **Participation in the class is only permitted for students who have completed the registration process.**
 - **Video recording is not allowed.**



4

Schedule and Evaluation

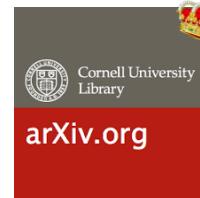
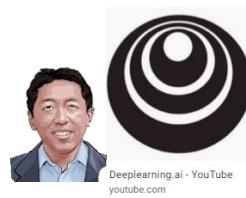
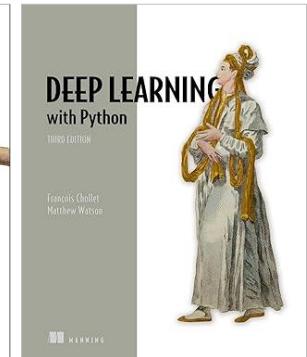
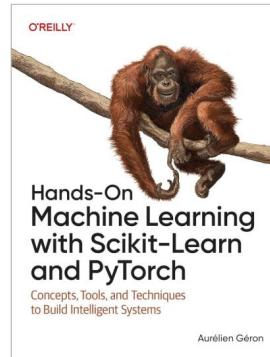
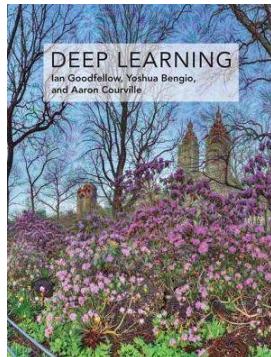
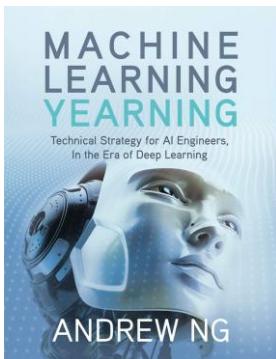
- **10% of in-class individual activities**
- **30% of Group Programming Homework:**
 - 20% for the work
 - 10% for point-by-point responses to comments
- **60% of Group Project:**
 - 20% for the proposal (Introduction and Related Works)
 - 40% for presentation and the final report*



AI

5

Recommended Materials



6



FAQ 1: Programming skill



7



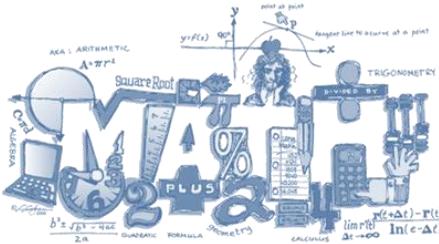
FAQ 2: Math skill

Quora Home Answer Spaces Notifications Search C

Learning Machine Learning +6

I do not have strong mathematics background, what should I learn in mathematics to be able to master Machine Learning and AI?

Answer Follow 253 Request



[14SEP2017] <https://www.quora.com/I-do-not-have-strong-mathematics-background-what-should-I-learn-in-mathematics-to-be-able-to-master-Machine-Learning-and-AI>

Quora Home Answer Spaces Notifications Search C

Andrew Ng, Co-founder of Coursera; Adjunct Professor of Stanford
Answered Sep 14, 2017 · Featured on Quora Sessions's Twitter · Upvoted by Ozan Ozegen, PhD Machine Learning, Ryerson University (2022) and Mir Junaid, Ph.D. Big Data & Machine Learning, University of Technology of Troyes (2020)

Originally Answered: I do not hold strong maths background, what all should I learn in Maths to be able to be master in Machine Learning and AI?

I think the most important areas of math for machine learning are, in decreasing order:

1. Linear algebra
2. Probability and statistics
3. Calculus (including multivariate calculus)
4. Optimization

After that, I think it falls off quickly. I've also found Information Theory helpful. You can find courses on all of these on Coursera or at most universities.

While it's hard to argue against knowing more math, I think the level of math needed to do machine learning effectively, or to get a PhD in machine learning, has decreased over the years. This is because machine learning has become more empirical (based on experiments) and less theoretical, especially with the rise of deep learning.

As a PhD student I had loved real analysis, and also studied differential geometry, measure theory, and algebraic geometry. While you're certainly be better off knowing these areas than not, in a world in which you have limited time, consider just spending more time studying machine learning itself, and even studying some of the other technical foundations for building AI systems, such as the algorithms that underly building big data systems and how to organize giant databases, plus HPC (high performance computing).

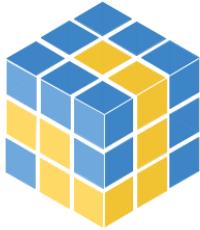
Best of luck!

400.6k views · View Upvoters · View Sharers · Answer requested by Sanjay M.S and Vijendr Gaorh

Upvote 9.6k Share 73

8

AI FAQ 3: Prerequisite

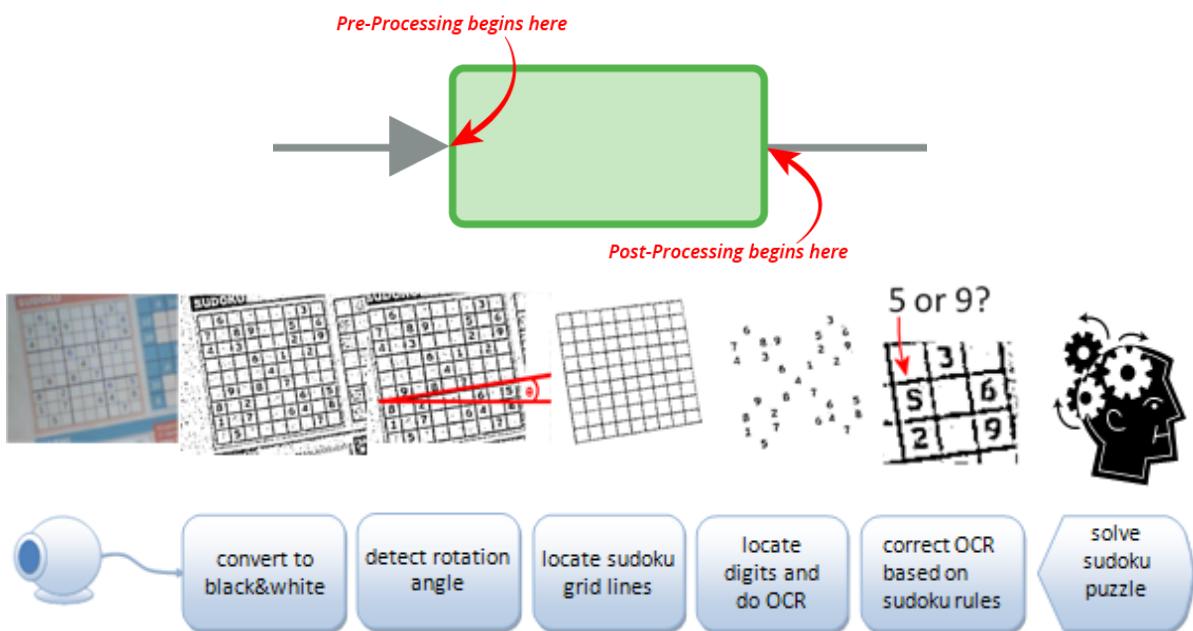


NumPy



Image generated from <https://imgflip.com/meme/generator/179756507/Skipping-steps>

9



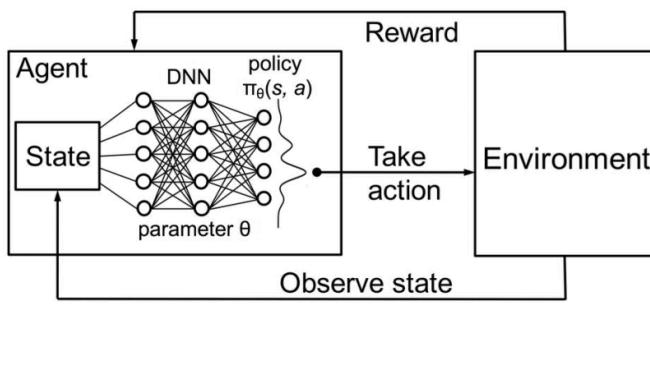
AI

Image source: <https://www.codeproject.com/Articles/238114/Realtime-Webcam-Sudoku-Solver>

10

AI FAQ 4: Not included

Deep Reinforcement Learning



Geometric Deep Learning

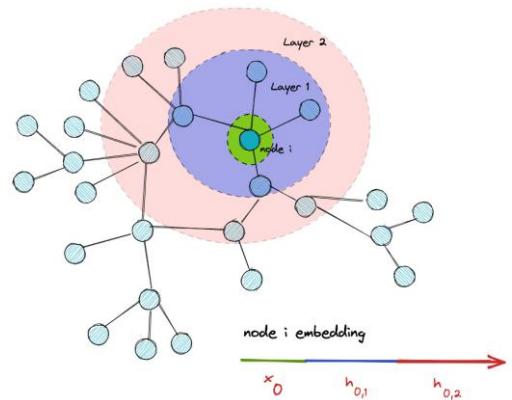


Image credit: <https://medium.com/@vishnuvijayanpv/deep-reinforcement-learning-artificial-intelligence-machine-learning-and-deep-learning-e52cb5974420>

11

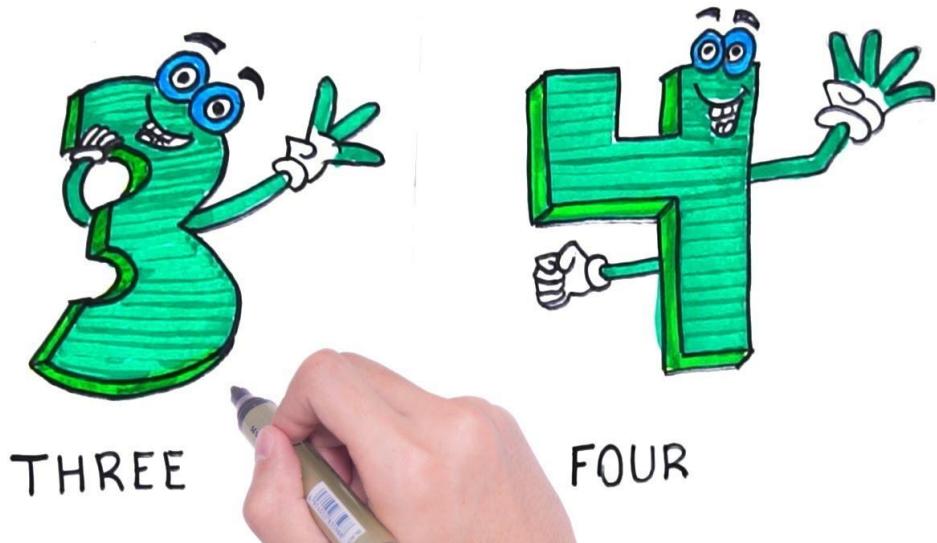
AI FAQ 5: Hardware and tools



12



FAQ 6: Group members



13

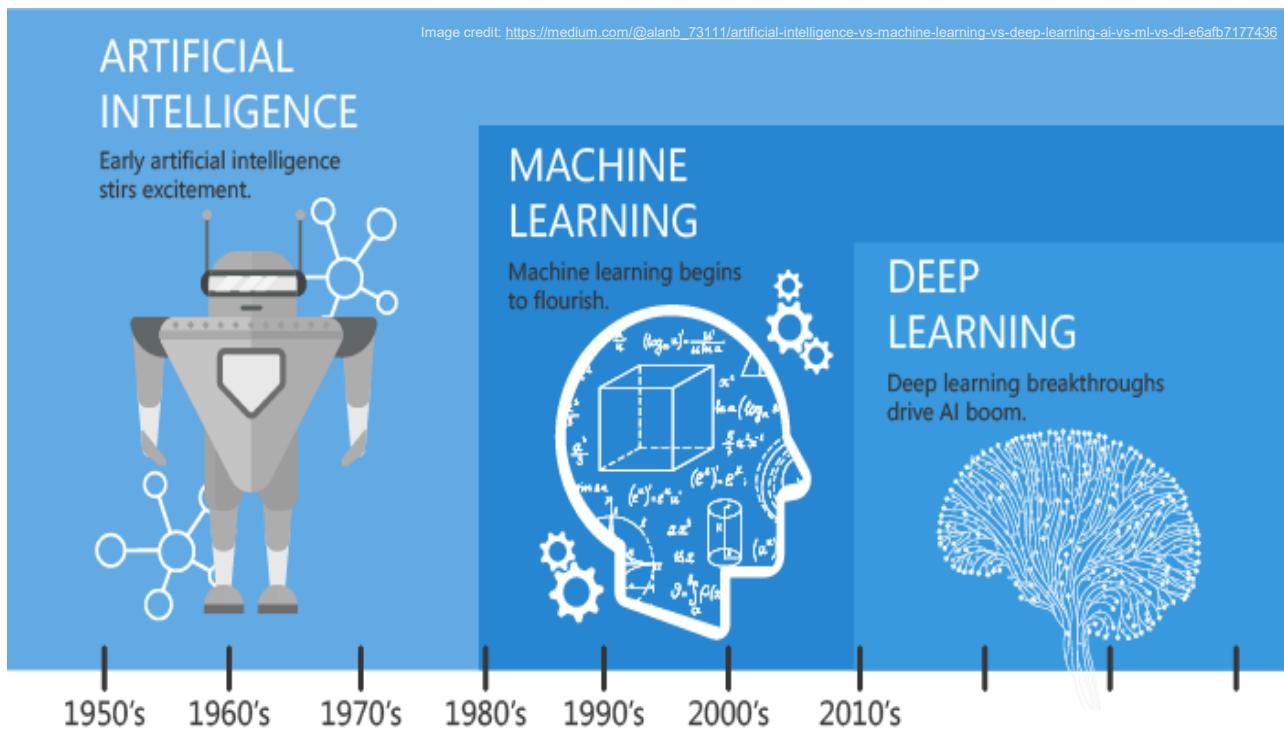


14

AI History of Modern DL

The rapid growth of Deep Learning in the modern era

15



16

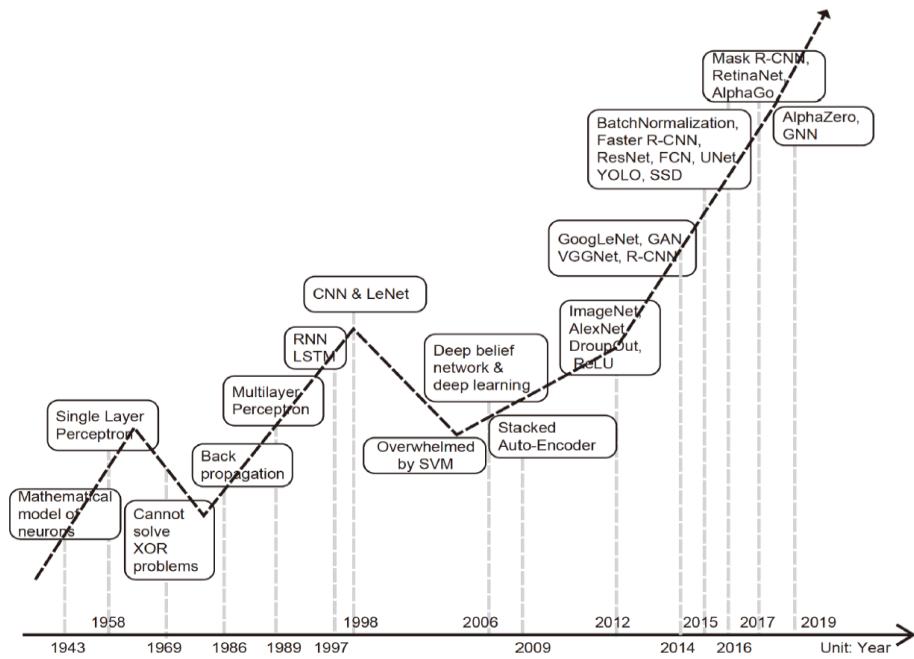


The universal approximation theorem

Feed-forward neural networks with at least one hidden layer can approximate any continuous function.

Papers that prove the theorem <https://ai.stackexchange.com/questions/13317/where-can-i-find-the-proof-of-the-universal-approximation-theorem>
A visual proof that neural nets can compute any function <http://neuralnetworksanddeeplearning.com/chap4.html>

17



Guo et al., "Application of deep learning in ecological resource research: Theories, methods, and challenges," Science China Earth Sciences, 2020 <https://doi.org/10.1007/s11430-019-9584-9>

18

AI THE OLD HISTORY

- **1943:** The first [neural network](#) was proposed by McCulloch and Pitts.
- **1955:** John McCarthy first coined the term [Artificial Intelligence](#).
- **1959:** Samuel coined and popularized the term [Machine Learning](#).

- **1982:** Birth of [Hopfield Network](#)
- **1986:** Birth of [backpropagation](#)
- **1986:** Birth of [RNN](#) and [AE](#)

- **1997:** Birth of [LSTM](#)
- **1998:** Birth of [CNN](#)

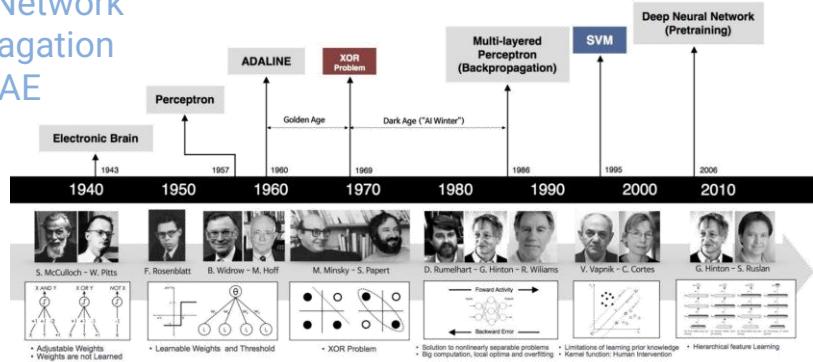


Image credit:

https://beamandrew.github.io/deeplearning/2017_02/23/deep_learning_101_part1.html

19

AI THE MODERN HISTORY

- **2006:** The [pretraining](#) technique by Hinton et al.



- **2012:** ImageNet evolution—[AlexNet](#)
- **2014:** Birth of [GRU](#), [VAE](#), and [GAN](#)
- **2015:** Birth of [diffusion model](#)
- **2017:** The [DeepFake](#) viral
- **2017:** Birth of the groundbreaking [Transformer](#)
- **2018:** NLP's ImageNet moment—[BERT](#)
- **2018:** ACM Turing Award to Deep NN



20

AI THE MODERN HISTORY

- **2020:** Birth of ViT
- **2020:** AlphaFold2 "This will change medicine. It will change research. It will change bioengineering. It will change everything."
- **2020:** Self-supervised learning trend in vision

- **2021:** Bridging the gap between vision and NLP—DALL-E and CLIP
- **2021:** The comeback of diffusion models in image generation

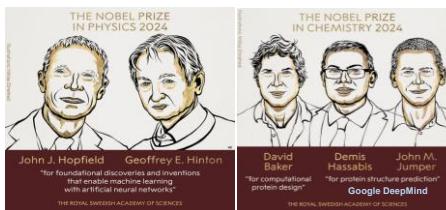
- **2022:** The year of Generative AI (massive adoption): Midjourney (July), Stable Diffusion (August), ChatGPT (November)

- **2023:** Generative AI's breakout year
- **2023:** The homerun year for LLMs
- **2023:** Birth of Flow Matching (Rectified Flow) and Mamba

21

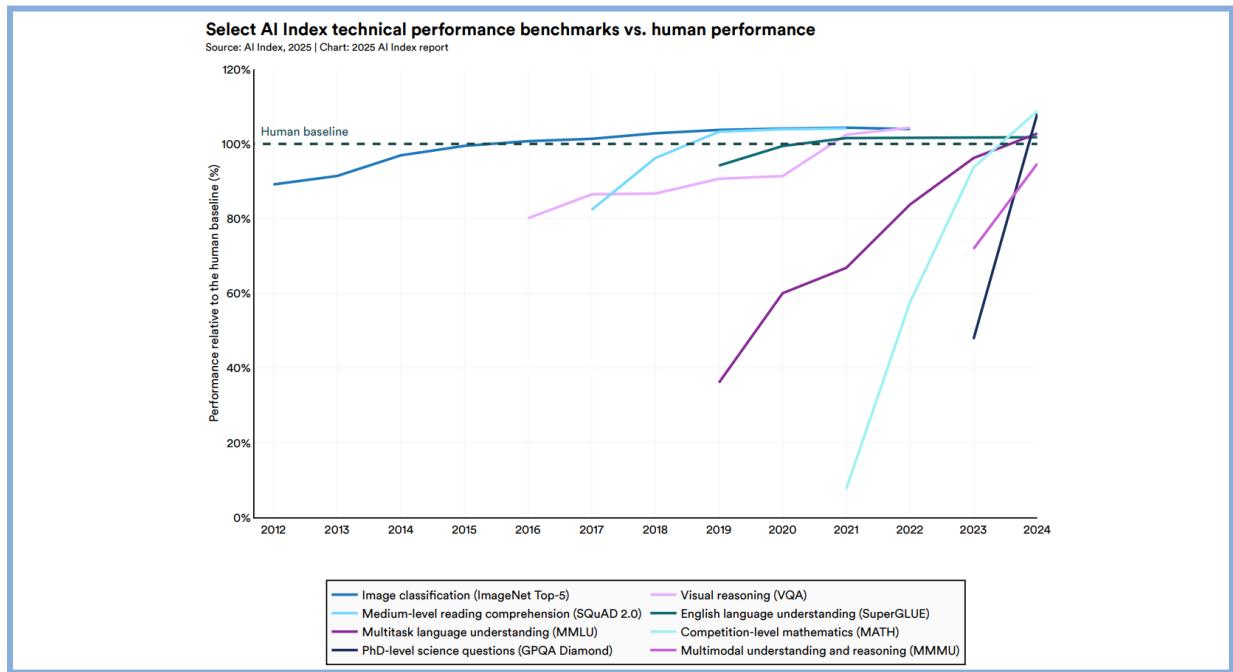
AI THE MODERN HISTORY

- **2024:** Vim, Multimodal LLMs, Agentic AI 🔥, Small Language Models, AI pricing war, Large Reasoning Models (LRMs), Generative video
- **2024:** Nobel Prize in Physics and Chemistry 2024 🏅
- **2024:** ACM Turing Award to RL 🎖️



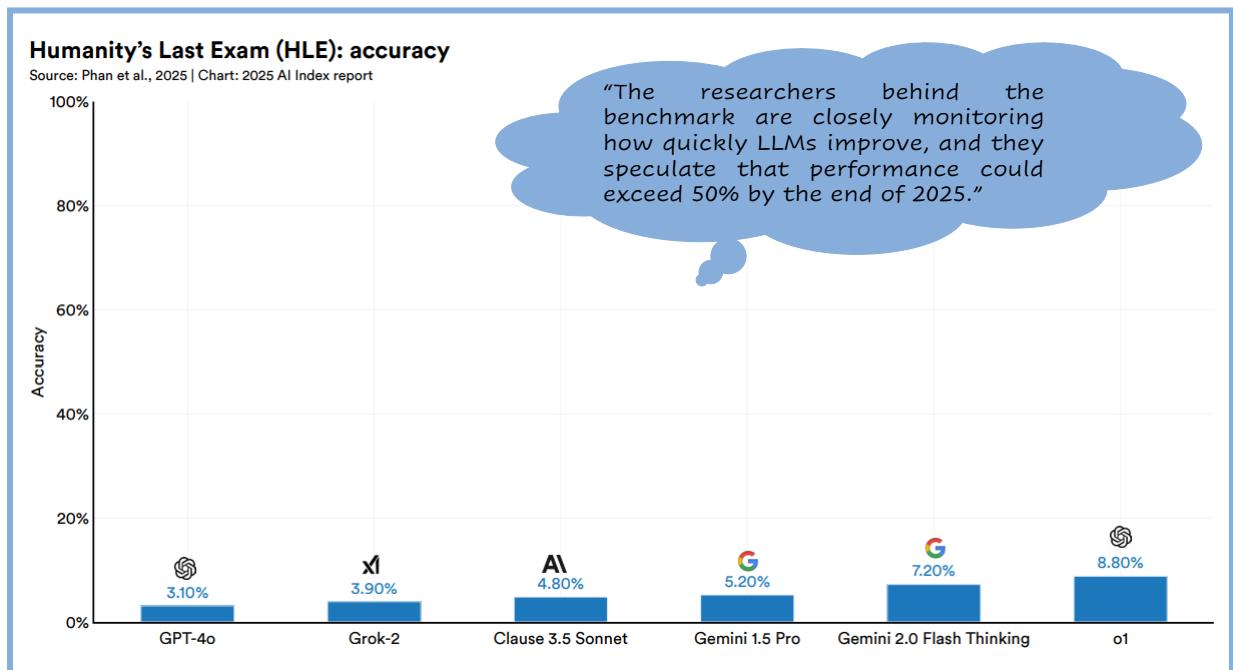
- **2025:** China dominated open-source AI models (DeepSeek, Qwen, Kimi).
 - (January) DeepSeek: the Sputnik moment for RL-based LLM training and Mixture of Experts (MoE) models
- **2025:** The year of agentic AI products (AI coding and more)
- **2025:** The dawn of AI's industrial age

22



HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

23

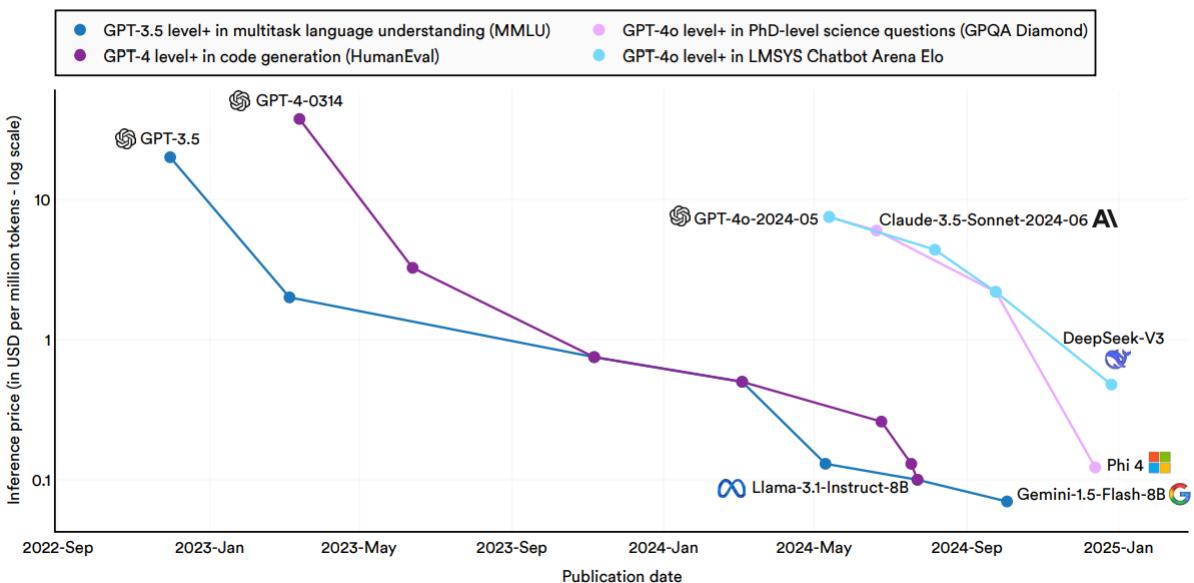


HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

24

Inference price across select benchmarks, 2022–24

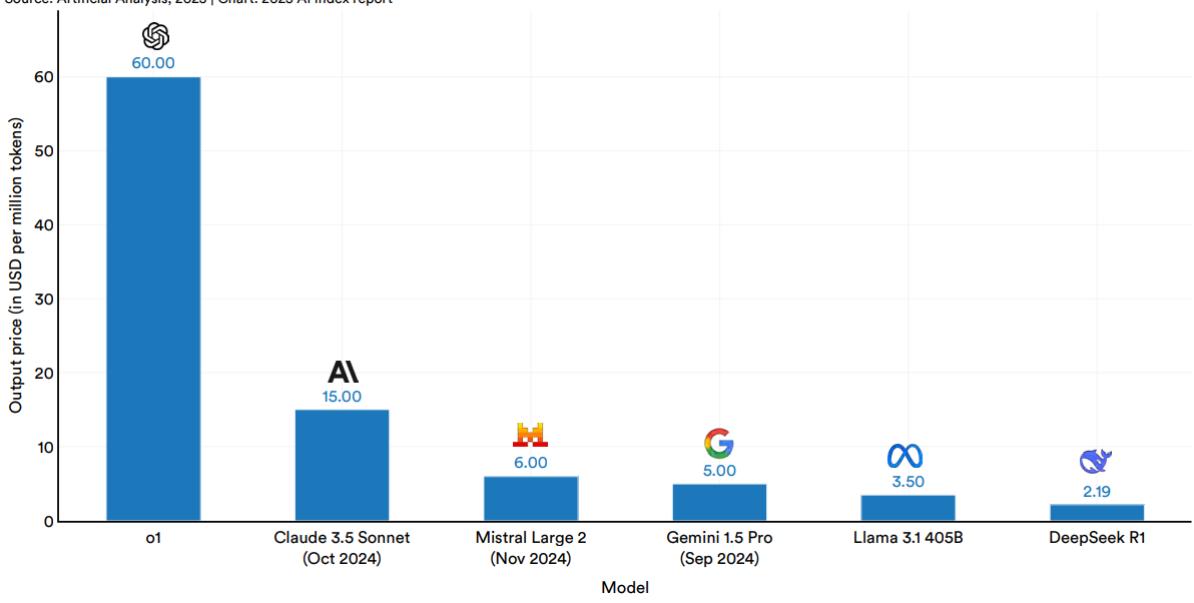
Source: Epoch AI, 2025; Artificial Analysis, 2025 | Chart: 2025 AI Index report

HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

25

Output price per million tokens for select models

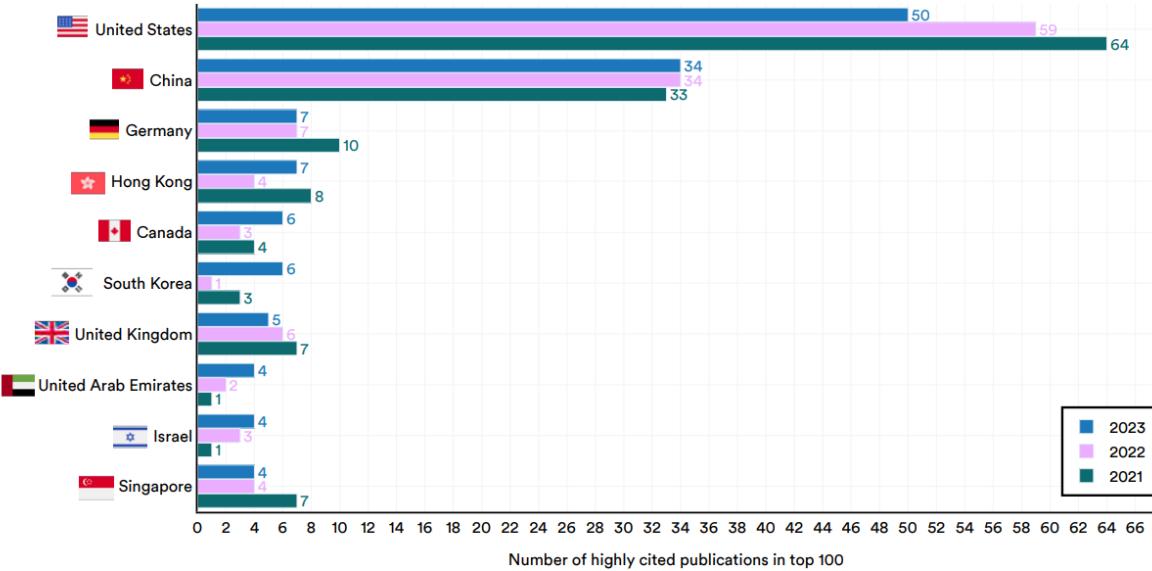
Source: Artificial Analysis, 2025 | Chart: 2025 AI Index report

HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

26

Number of highly cited publications in top 100 by select geographic areas, 2021–23

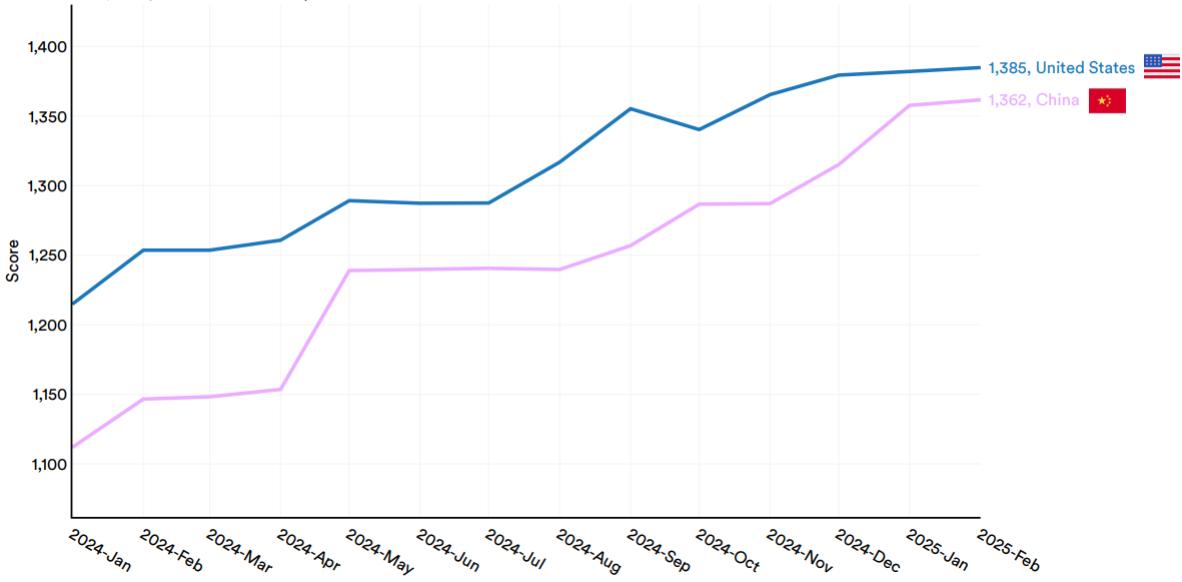
Source: AI Index, 2025 | Chart: 2025 AI Index report

HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

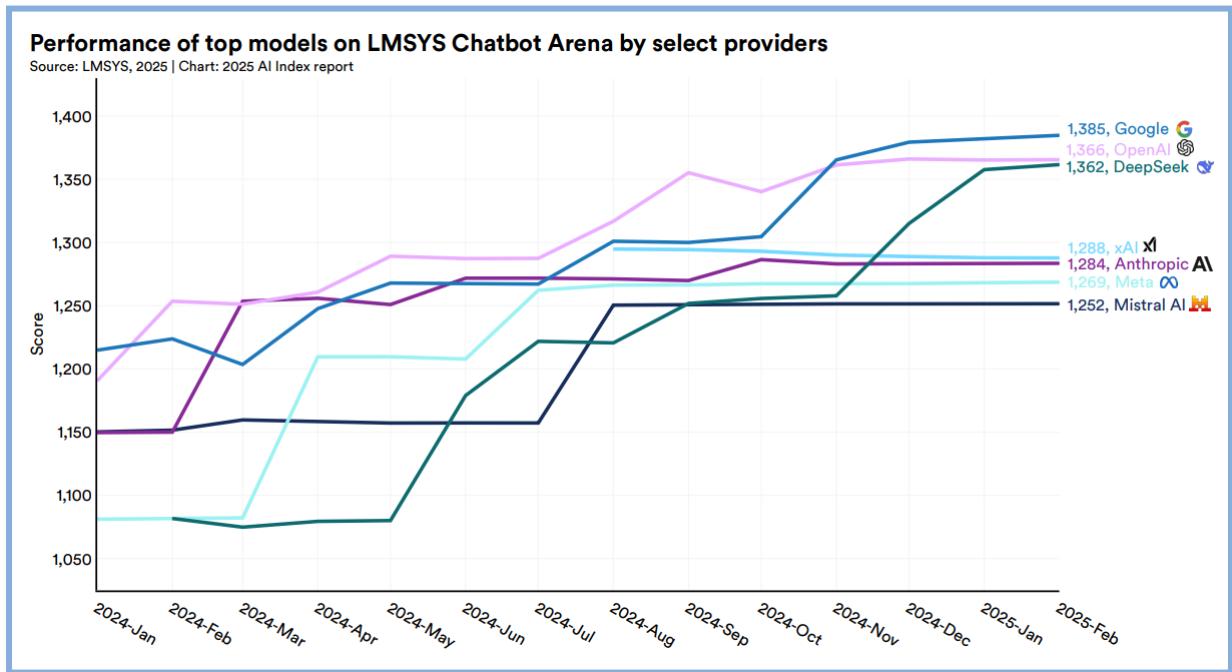
27

Performance of top United States vs. Chinese models on LMSYS Chatbot Arena

Source: LMSYS, 2025 | Chart: 2025 AI Index report

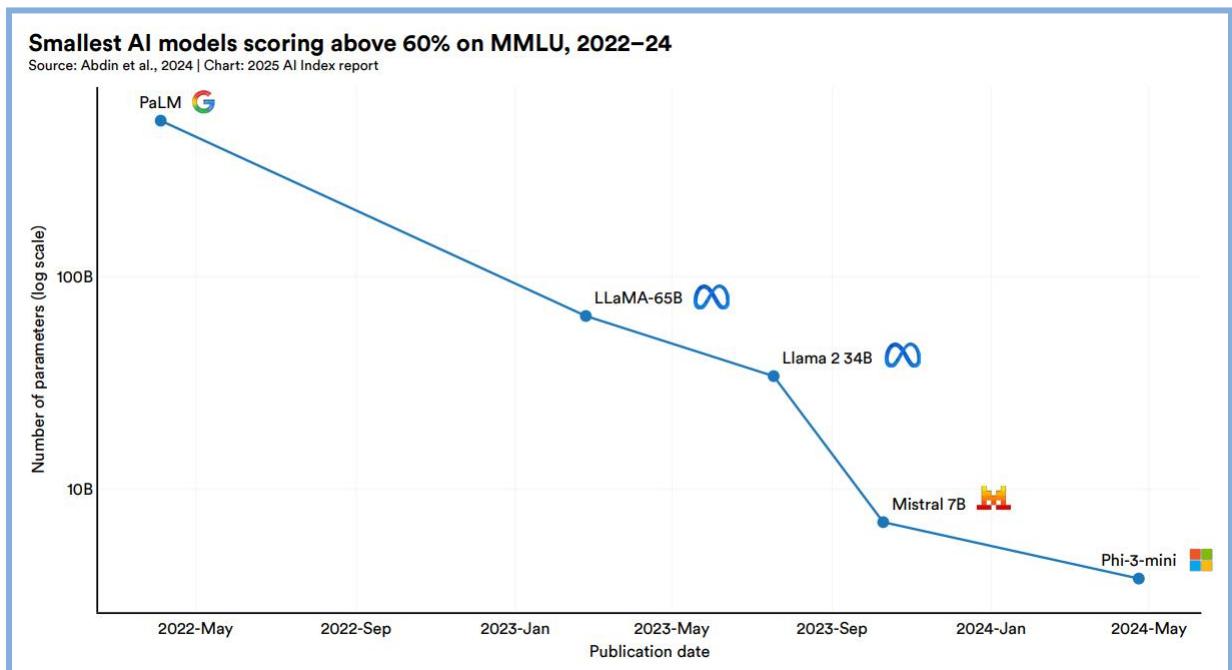
HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

28



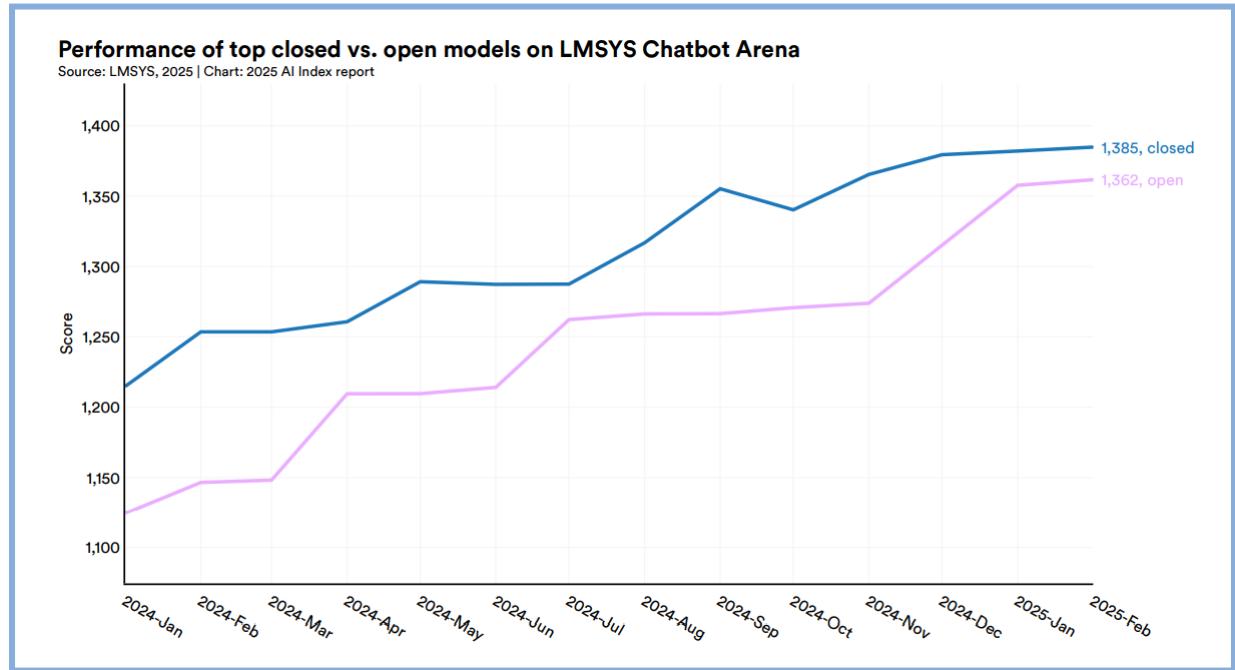
HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

29



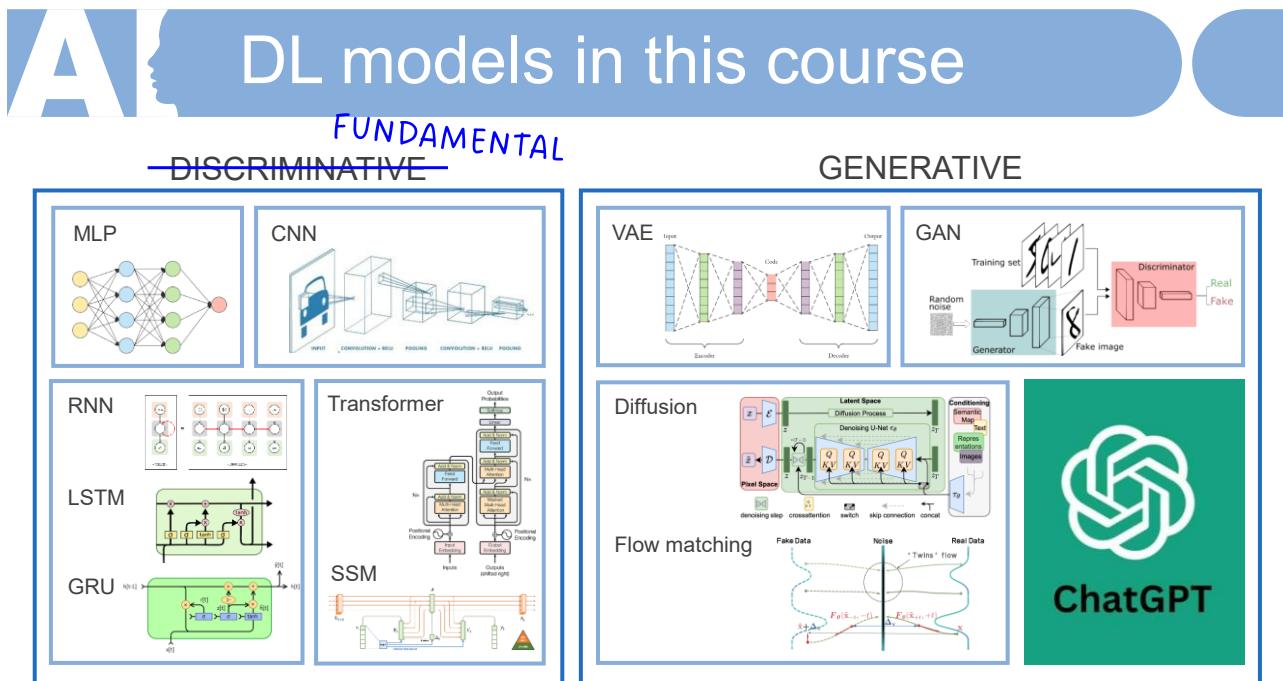
HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

30



HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

31

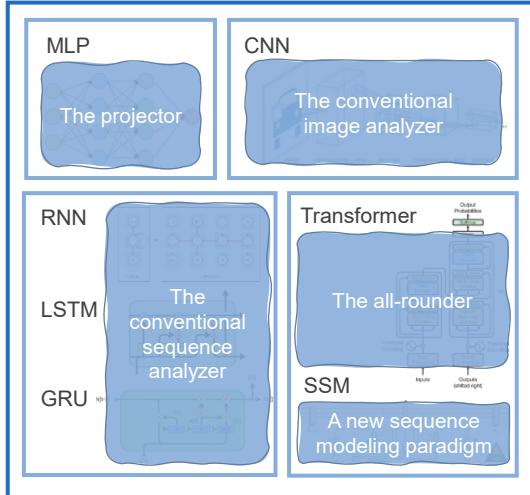


32

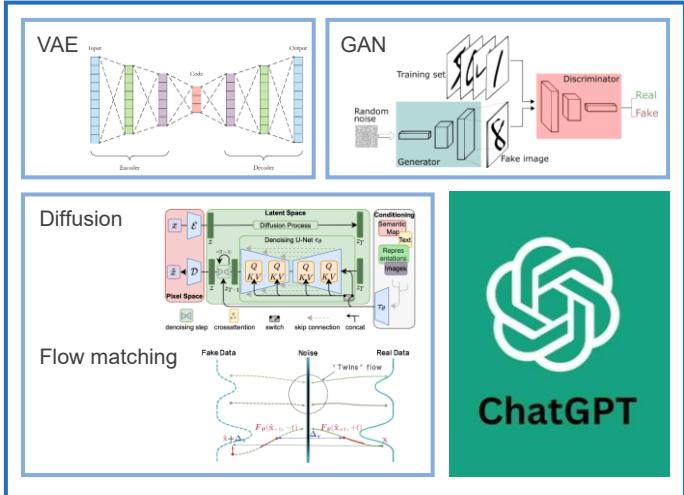
AI DL models in this course

FUNDAMENTAL

~~DISCRIMINATIVE~~



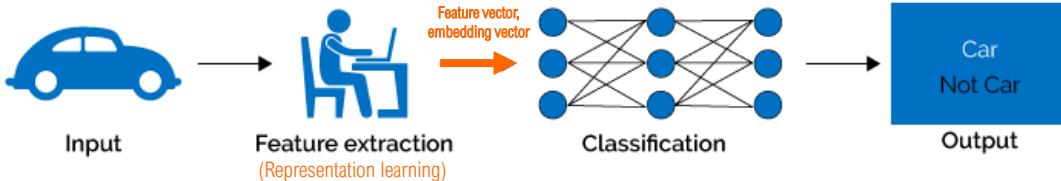
GENERATIVE



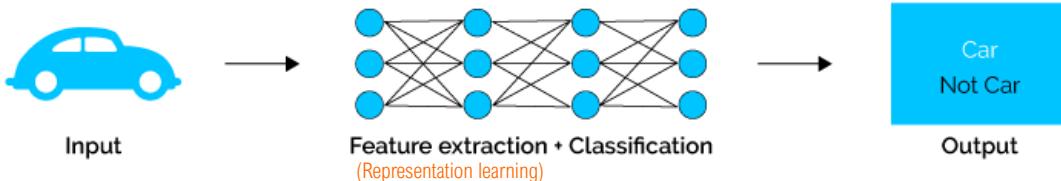
33

AI Paradigm shift: ML > DL

Machine Learning

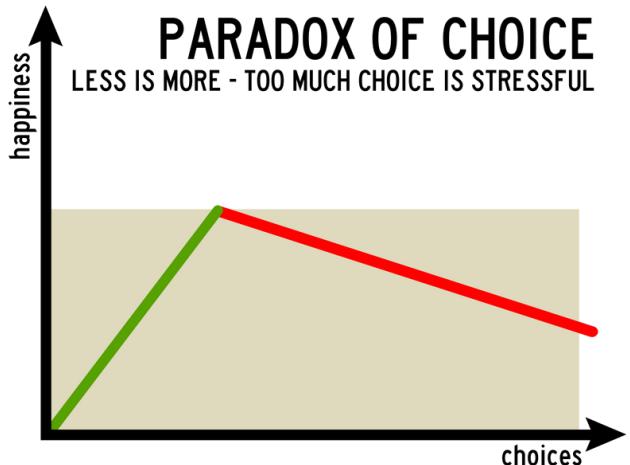
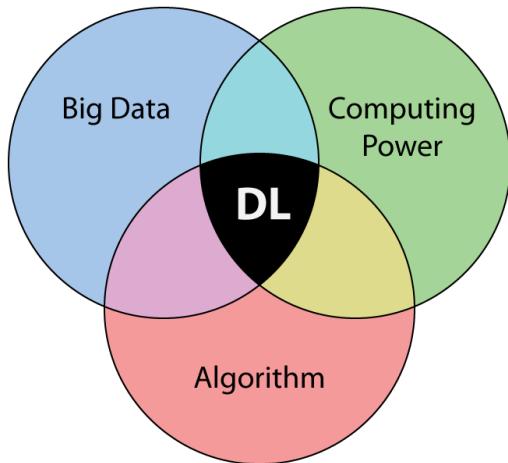


Deep Learning

Image credit: <https://www.softwaretestinghelp.com/data-mining-vs-machine-learning-vs-ai/>

34

AI WHY Deep Learning



35

AI WHY NOT Deep Learning

- Deep learning usually requires a huge amount of input data. For supervised learning, data collection and annotation are sometimes very laborious.

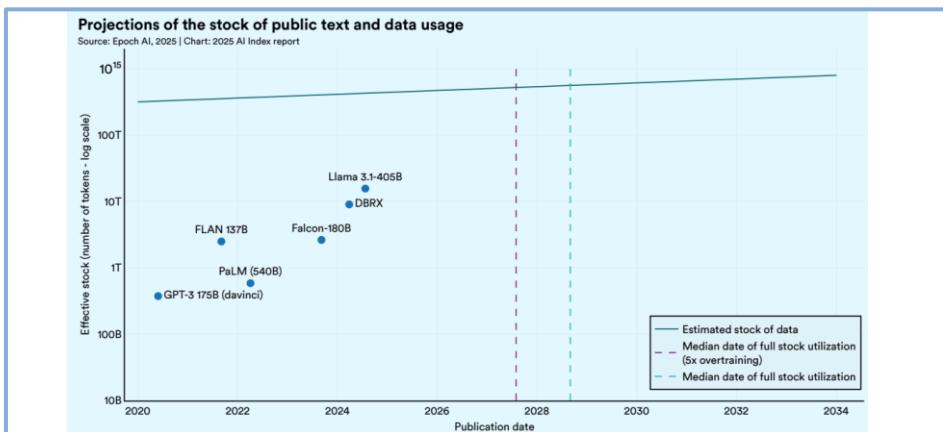


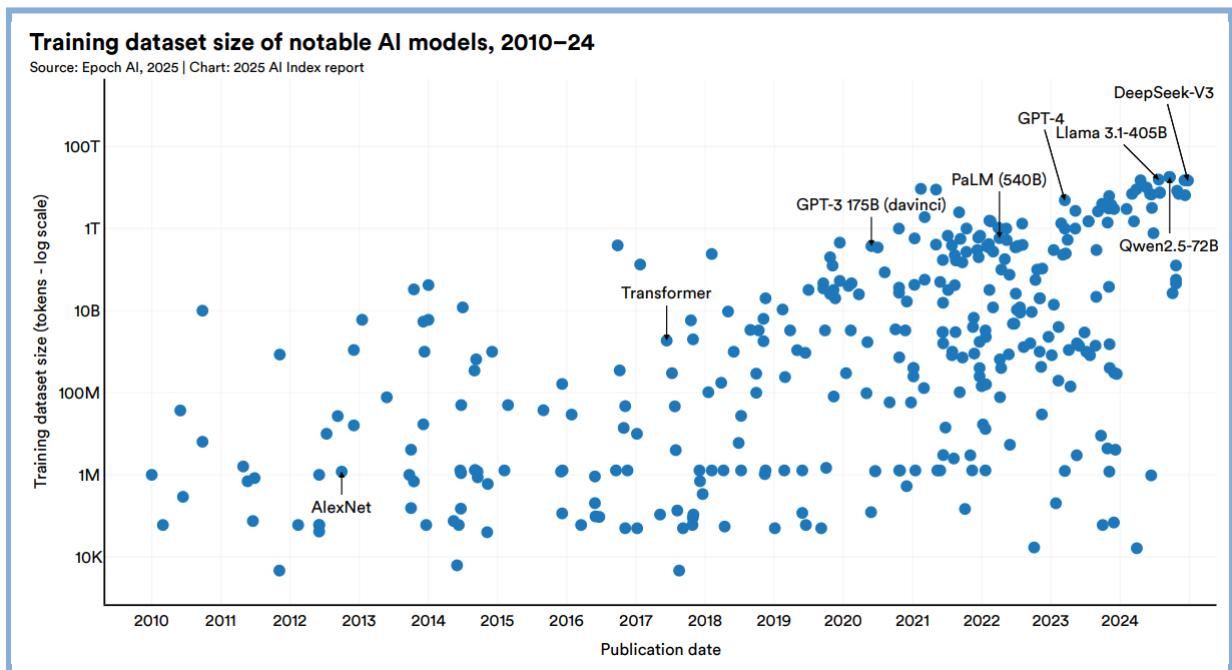
Image credit: HAI AI Index Report, April 2025, https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

36

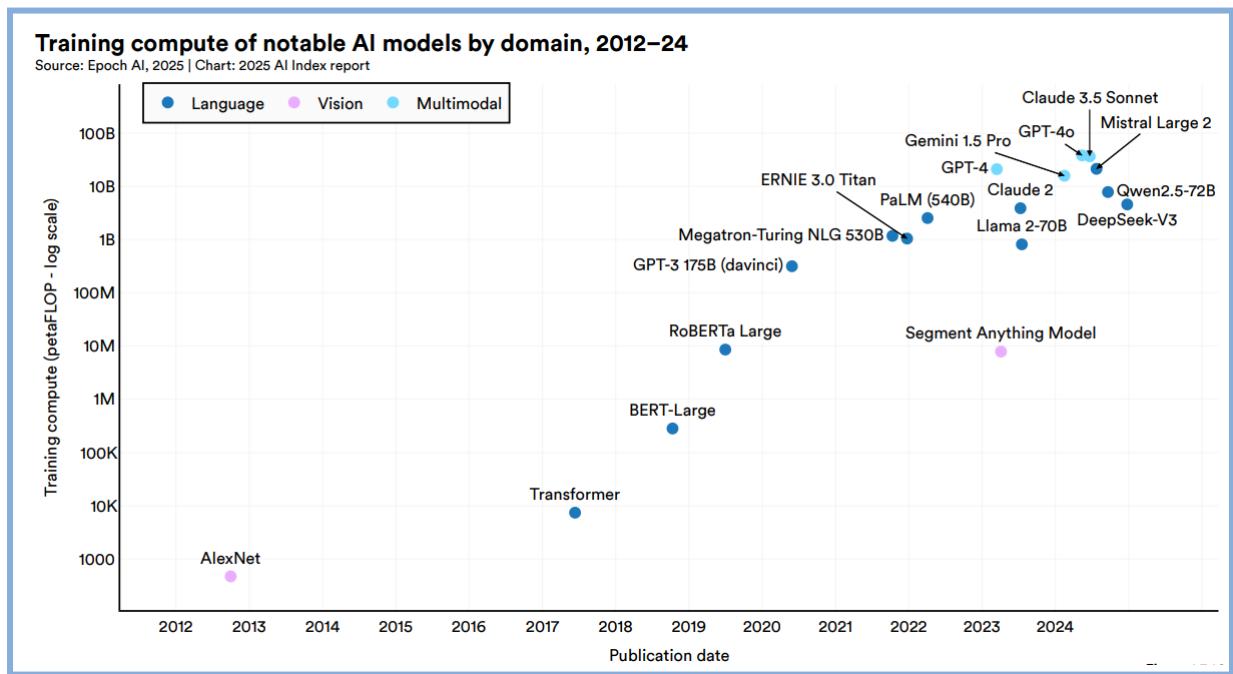
AI WHY NOT Deep Learning

- Training deep networks consumes huge computational resources (GPU memory, training time, and inference time).
 - The largest **StyleGAN2** (Nvidia 2019) was trained for 69d 23h on one Tesla V100 GPU.
 - At theoretical 28 TFLOPS for Tesla V100 and lowest 3-year reserved cloud pricing, **GPT-3** (OpenAI 2020) was still trained for 355 GPU-years and cost \$4.6M.
 - **Vision Transformer** (Google 2020) was trained for 2,500 TPUs v3-core-days.
 - **CLIP** (OpenAI 2021) was trained for 2 weeks on 256 GPUs.
 - **WangchanBERTa** (VISTEC x DEPA 2021) was trained for 134 days on eight Tesla V100 GPUs (one DGX-1 server).
 - **Stable Diffusion** (2022) was trained with 256 A100 GPUs on Amazon Web Services for a total of 150,000 GPU hours, at \$600,000.
 - **GPT-4** (OpenAI 2023) used an estimated \$78 million worth of computing to train.
 - **Gemini Ultra** (Google 2023) cost \$191 million to compute.
 - In 2024, OpenAI is likely to spend inference costs around \$4 billion on processing power supplied by Microsoft (at a special discount rate), and expects to spend \$3 billion on training models and data.
 - **DeepSeek-R1** (DeepSeek 2025) was officially reported as using only \$294,000 training cost with 512 NVIDIA H800 GPUs, excluding ~\$6 million for the base model.

37

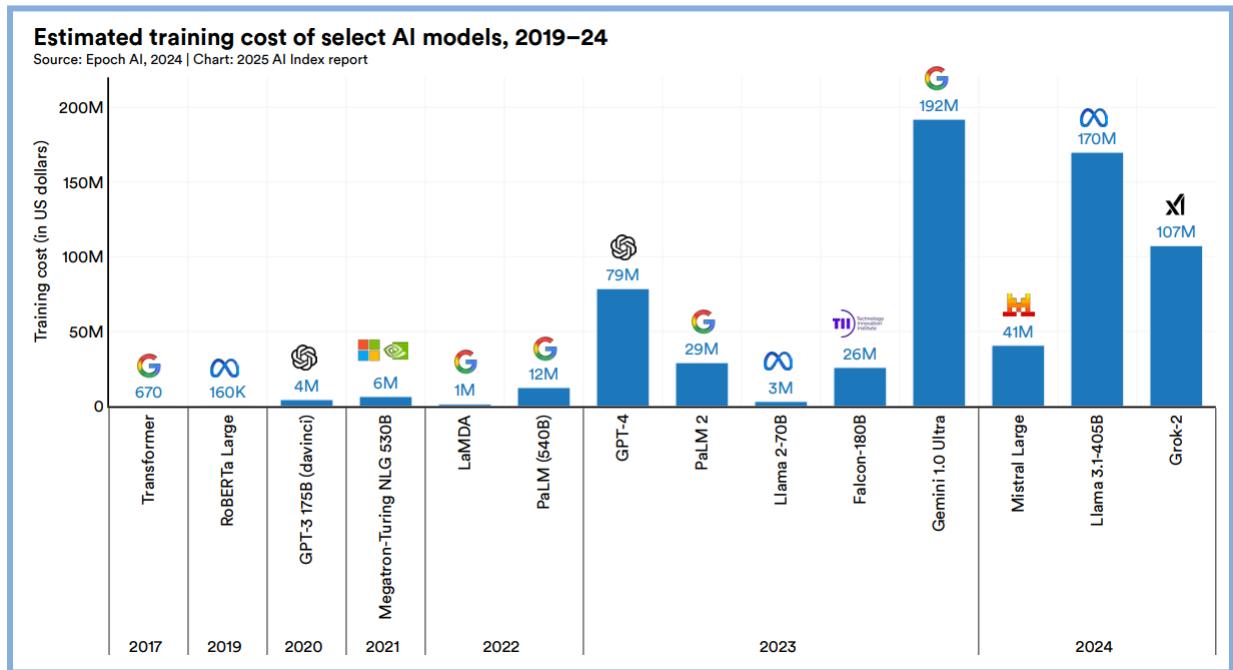
HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

38



HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

39

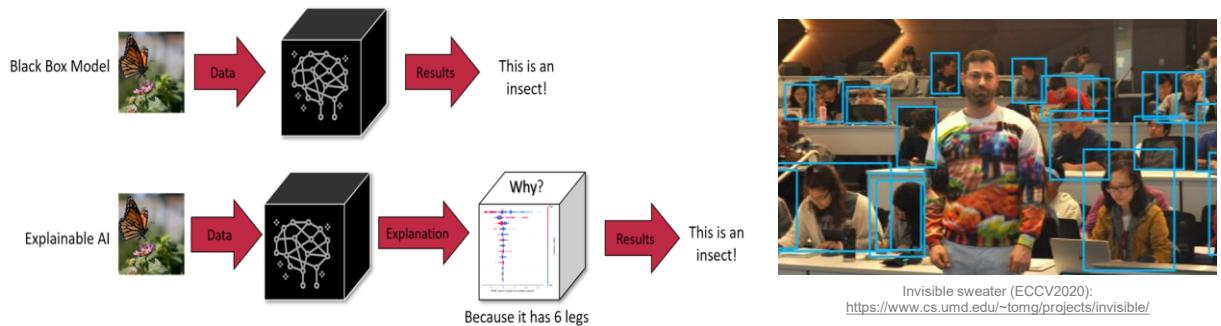


HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

40

AI WHY NOT Deep Learning

- Deep learning suffers from **the infamous black-box issue**.
 - Hence, it is necessary to use as many indicators as possible to evaluate DLs
- Deep learning can fail and can be attacked (**Adversarial Attack**).



41



Hardware Acceleration

Prepare necessary hardware to accelerate our deep learning project

42

REUTERS World ▾ Business ▾ Markets ▾ Sustainability ▾ Legal ▾ Breakingviews ▾ Technology ▾ Ir

Technology

Nvidia briefly joins \$1 trillion valuation club

By Akash Sriram ▾ and Samritha A ▾
May 31, 2023 7:33 AM GMT+7 · Updated 3 days ago

Reuters 31MAY2023 <https://www.reuters.com/technology/nvidia-sets-eye-1-trillion-market-value-2023-05-30/>

Market Cap data as of May 30, 2023 | Source: [richards.com/companymarketcap.com](#)

Visual Capitalist 30MAY2023 <https://www.visualcapitalist.com/nvidia-joins-the-trillion-dollar-club/>

Reuters World ▾ Business ▾ Markets ▾ Sustainability ▾ Legal ▾ Breakingviews ▾ Technology ▾ Investigations

Technology

Nvidia overtakes Apple as No. 2 most valuable company

By Noel Randewich
June 6, 2024 1:18 PM GMT+7 · Updated an hour ago

Nvidia's stock market value overtakes Apple

Sources: LSEG
Created by Thomson Reuters

Reuters 6JUN2024 <https://www.reuters.com/technology/nvidia-verge-overtaking-apple-no-2-most-valuable-company-2024-06-05/>

Reuters World ▾ Business ▾ Markets ▾ Sustainability ▾ Legal ▾ Breakingviews ▾ Technology ▾ Ir

Technology

Nvidia becomes world's most valuable company

By Reuters
June 19, 2024 1:17 AM GMT+7 · Updated 2 months ago

Reuters 19JUN2024 <https://www.reuters.com/technology/view-nvidia-becomes-worlds-most-valuable-company-2024-06-18/>

43

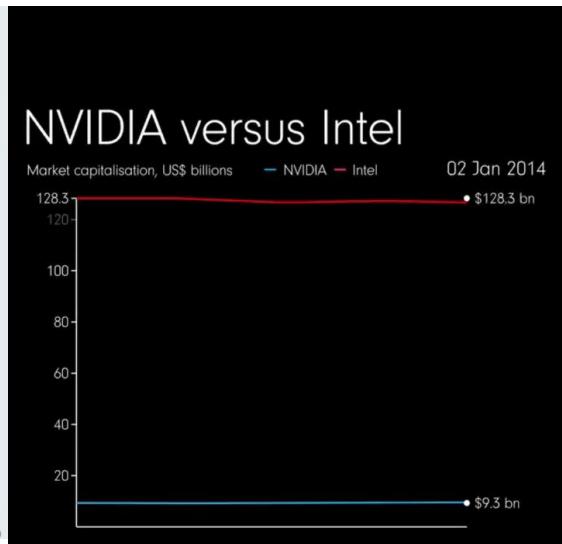
Nvidia H100 GPU Shipments by Customer

Estimated 2023 H100 shipments by end customer.

Omdia estimates Nvidia sold ~500k H100 and A100 GPUs in Q3, and lead time for H100-based servers is up to 52 weeks.

Source: Omdia Research

yahoo/finance Search for news, symbols or companies



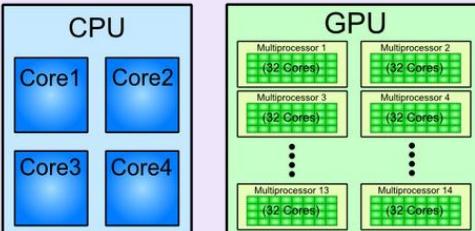
NVIDIA Crushes Rivals: Secures Unprecedented 90% of GPU Market in Q3 2024

Nauman khan

December 12, 2024 • 1 min read

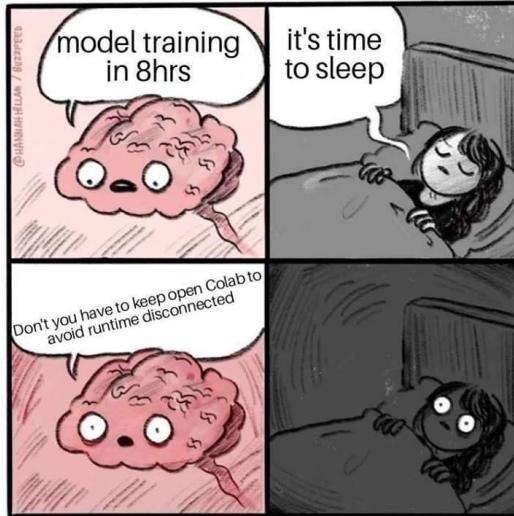
44

AI Local Machine vs. Cloud

  	<p>Free Google</p> <p>colab Pro+ version for \$49.99/month since 8/2021 Pro version for \$9.99/month since 2/2020 Free version for public since 2018</p> <p>kaggle Free GPU since 2018 (approximately)</p>  
Local Machine	Cloud

45

REVIEW: Google Colab



Version 1:

```
function ClickConnect(){
  console.log("Connect Clicked - Start");
  document.querySelector("#top-toolbar > colab-connect-button").shadowRoot.querySelector("button").click();
  console.log("Connect Clicked - End");
};  
setInterval(ClickConnect, 60000)
```

Version 2: If you would like to be able to stop the function, here is the new code:

```
var startClickConnect = function startClickConnect(){
  var clickConnect = function clickConnect(){
    console.log("Connect Clicked - Start");
    document.querySelector("#top-toolbar > colab-connect-button").shadowRoot.querySelector("button").click();
    console.log("Connect Clicked - End");
  };
  setInterval(clickConnect, 60000);
}
```

```
var stopClickConnectHandler = function stopClickConnect() {
  console.log("Connect Clicked Stopped - Start");
  clearInterval(intervalId);
  console.log("Connect Clicked Stopped - End");
};

return stopClickConnectHandler;
};

var stopClickConnect = startClickConnect;
```

In order to stop, call:

```
stopClickConnect();
```



Credit: <https://stackoverflow.com/questions/57113226/how-can-i-prevent-google-colab-from-disconnecting>

46

REVIEW: Google Colab Pro+

The image shows three cards side-by-side, each representing a different plan:

- Pay As You Go:** Shows 0 compute units available and a price of THB343.47 for 100 Compute Units. It includes a note that units expire after 90 days and a link to purchase more.
- Recommended Colab Pro:** Shows 100 compute units per month at a price of THB343.47 per month. It lists benefits: Faster GPUs (upgrade to more powerful premium GPUs), More memory (access higher memory machines), and Terminal (ability to use a terminal with the connected VM).
- Colab Pro+:** Shows 500 compute units per month at a price of THB1,677.76 per month. It lists benefits: Faster GPUs (priority access to more powerful premium GPUs), More memory (access higher memory machines), Background execution (upgrade notebooks to keep executing for up to 24 hours even if browser closed), and Terminal (ability to use a terminal with the connected VM).

Update on AUG2023:

- Nvidia A100 GPU is available for both Colab Pro and Pro+.
- The better the GPU, the more compute units it will consume.
- When running out of compute units, we can simply buy more compute units to continue running.



- Good:**

- Convenient as we can train the model while being far away from our computer or even when closing the computer
- At least 10+ hours of continuous running without interruption

- Bad:**

- After 2 continuous weeks of heavy GPU usage (e.g., 4 concurrent tabs), it banned us from further GPU usage with a suggestion to move to Google Cloud ML.

AI

Review from a real user (using Google Colab Pro+ during January-February 2022)

47



- TPU** is an AI accelerator hardware specially designed for neural networks and is proprietary to Google.
- TPU** has been used to power Google's data centers internally since 2015 and in 2018 made available for 3rd party use.
- According to Google (in 2017), **TPU**-v1 was about 80x faster than CPU and 30x faster than GPU at neural network inference.
- Currently, **TPU** is mostly available as a cloud service or a smaller version of the chip (Edge TPU).



TPU v1
Launched in 2015
Inference only



TPU v2
Launched in 2017
Inference and training

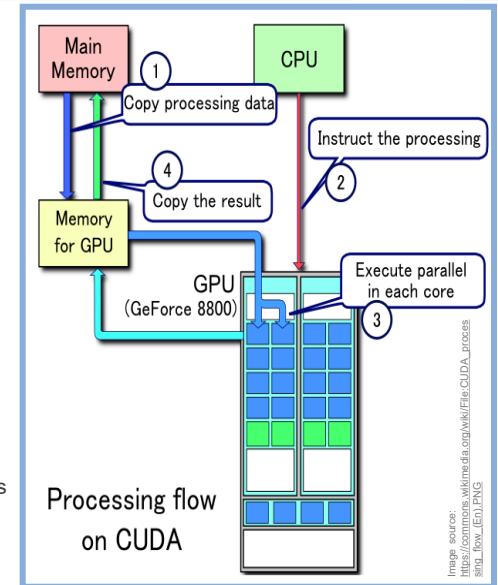
5APR2017: <https://cloud.google.com/blog/products/gcp/quantifying-the-performance-of-the-tpu-our-first-machine-learning-chip>

48

AI CPU vs. GPU vs. TPU

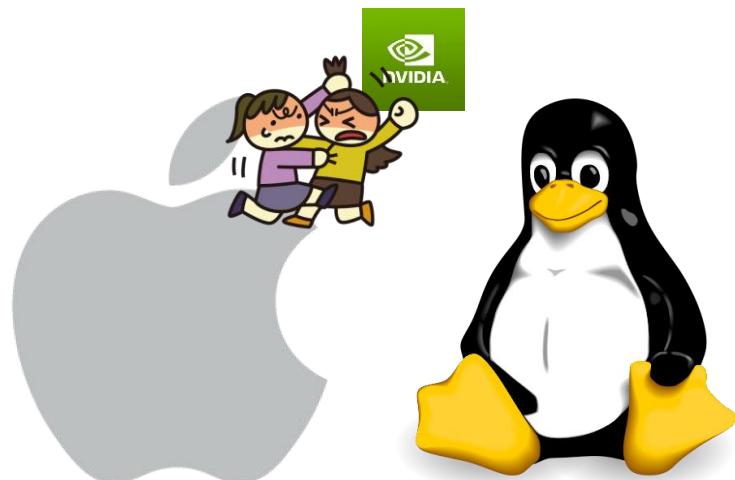
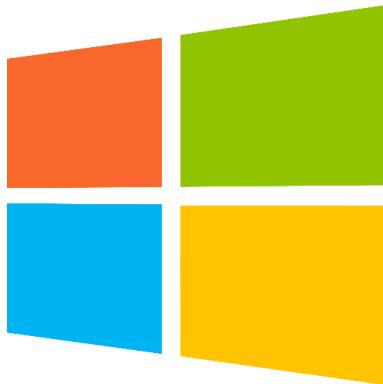
- Overall, **TPU** has the highest training throughput (the number of examples trained per second).
- **Deep FC networks:**
 - For a big batch size, it is **TPU**, GPU and CPU respectively.
 - For a big model, **GPU** is the best.
 - For a big batch size and model, **GPU** is better than CPU because of better parallelization.
 - For very big models, **CPU** is better than GPU, and GPU is better than TPU. This is because CPU has the highest memory per core.
- **Deep CNNs:**
 - **TPU** is better than GPU and is the best for training CNN, particularly very big CNNs. This is because TPU is specifically optimized for spatial reuse characteristics of CNN.
- **RNNs:**
 - For speed aspect, **TPU** is a lot better than GPU.
 - Right now, TPU is good at dense computation (MatMul) whereas GPU is better in non-MatMul operations. In the future, if TPU is further optimized for non-MatMul operations, it will do even better than GPU.

[31JUL2019] Harvard University, <https://arxiv.org/abs/1907.10701>



49

AI Deep Learning's Computer



Why not Mac OS? [23NOV2019] Apple and Nvidia are over <https://gizmodo.com/apple-and-nvidia-are-over-1840015246>
How much GPU's RAM for training DL models? [18FEB2020] Choosing the Best GPU for Deep Learning in 2020 <https://lambdalabs.com/blog/choosing-a-gpu-for-deep-learning/>
Why Mac OS? [15NOV2020] How is the Apple M1 going to affect Machine Learning? <https://medium.com/disruptive-nerd/how-is-the-apple-m1-going-to-affect-machine-learning-2d9da1beef86>

50



How to choose a computer for deep learning?

51

AI Neural Processing Unit (NPU)



Home AI Data Center Driving Gaming Pro Graphics Robotics Healthcare Startups AI Podcast NVIDIA Life

NVIDIA RTX 500 and 1000 Professional Ada Generation Laptop GPUs Drive AI-Enhanced Workflows From Anywhere

Thin and light designs deliver advanced AI, compute and graphics horsepower for professionals on the go.

February 26, 2024 by [John Della Bona](#)

The next generation of mobile workstations with Ada Generation GPUs, including the RTX 500 and 1000 GPUs, will include both a neural processing unit (NPU), a component of the CPU, and an NVIDIA RTX GPU, which includes Tensor Cores for AI processing. The NPU helps offload light AI tasks, while the GPU provides up to an additional 682 TOPS of AI performance for more demanding day-to-day AI workflows.

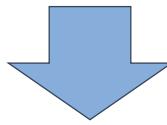
26FEB2024 <https://blogs.nvidia.com/blog/rtx-ada-ai-workflows/>

52

AI Neural Processing Unit (NPU)

- **Problems of GPU:**

- **Still rely on CPU:** Although GPU has the advantage in parallel computing capability, it does not work alone and needs the CPU's co-processing. The construction of neural network models and data streams are still carried out on the CPU.
- **The problem of high-power consumption and large size:** The higher the performance, the larger the GPU, the higher the power consumption, and the more expensive it is, which will not be available for some small devices and mobile devices.



- **We need a small size, low power consumption, high computational performance, and high computational efficiency of a dedicated chip.** This is the birth of **NPU**.

Credit: (28DEC2021) <https://www.utmel.com/blog/categories/integrated%20circuit/neural-processing-unit-npu-explained>

53

AI Neural Processing Unit (NPU)

- **NPUs are becoming increasingly common in PC, laptop, and mobile domains. It usually refers to all specialized AI processors.**
 - For example: AI chip in Apple iPhone, Exynos 9 (9820) the AI-enabled NPU in Samsung mobile
- **NPUs are designed to complement the functions of CPUs and GPUs, ensuring that no single processor gets overwhelmed, maintaining smooth operation across the system.**
 - CPUs handle a broad range of tasks.
 - GPUs excel in rendering detailed graphics.
 - **NPUs** specialize in executing AI-driven tasks swiftly.

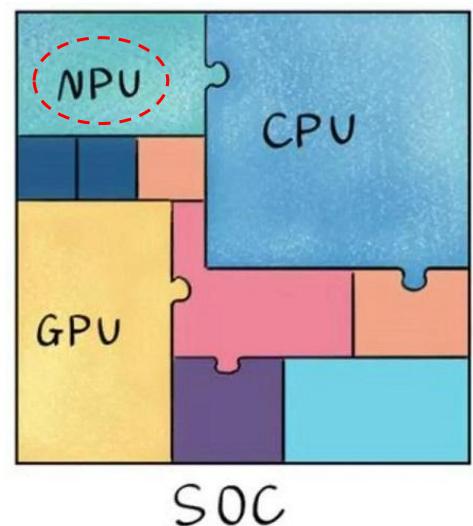


Image source: (28DEC2021) <https://www.utmel.com/blog/categories/integrated%20circuit/neural-processing-unit-npu-explained>

54

AI Neural Processing Unit (NPU)

- **CPU vs. GPU vs. NPU:**

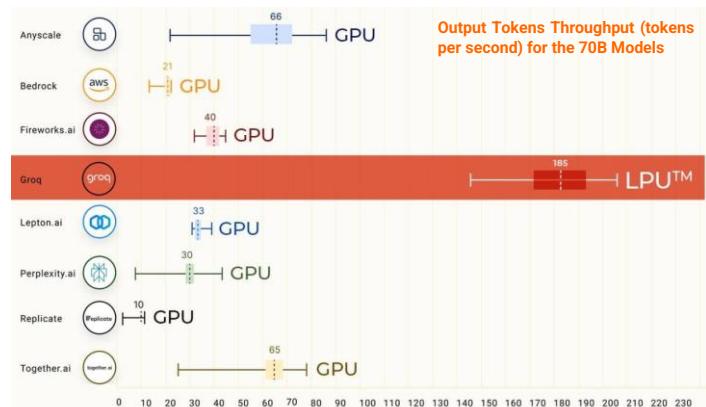
- Unlike traditional CPUs and GPUs, an **NPU**, at its core, is a specialized processor explicitly designed/optimized for executing AI/ML/DL algorithms, particularly for ANN mathematical computation and DL acceleration.
- This degree of specialization allows **NPUs** to deliver significantly higher performance for AI workloads compared to CPUs and even GPUs in certain scenarios.
 - **NPUs** excel in processing vast amounts of data in parallel, making them ideal for tasks like image recognition, natural language processing, and other AI-related functions.
 - For example, if you were to have an **NPU** within a GPU, the **NPU** may be responsible for a specific task like object detection or image acceleration.
- The concept of **GPNPU** (GPU-NPU hybrid) has emerged, aiming to combine the strengths of GPUs and **NPUs**. GPNPUs leverage the parallel processing capabilities of GPUs while integrating **NPU** architecture for accelerating AI-centric tasks.

Credit: (28DEC2021) <https://www.utmel.com/blog/categories/integrated%20circuit/neural-processing-unit-npu-explained>

55

AI Language Processing Unit (LPU)

- In February 2024, a largely unknown company, **Groq** (not to be confused with Elon's Grok AI) demonstrated unprecedented speed running open-source LLMs such as Llama-2 (70 billion parameters) at more than 100 tokens per second, and Mixtral at nearly 500 tokens per second per user on a Groq's **Language Processing Unit (LPU)**.
- **LPU** is a special kind of computer brain designed to be highly specialized for language-intensive tasks, boasting optimized hardware and software. In short, it uses faster processing times and lower power consumption.
- **LPU** employs sequential processing, meticulously handling tasks step-by-step, mirroring the natural flow of language.



Sources:
 26FEB2024 <https://www.linkedin.com/pulse/new-ai-compute-paradigm-language-processing-unit-lpu-dheerend/>
 27FEB2024 <https://www.analyticsvidhya.com/blog/2024/02/what-is-the-difference-between-lpu-and-gpu/>

56

Forbes

Cerebras Update: The Wafer Scale Engine 3 Is A Door Opener

Karl Freund Contributor
Founder and Principal Analyst, Cambrian-AI Research LLC

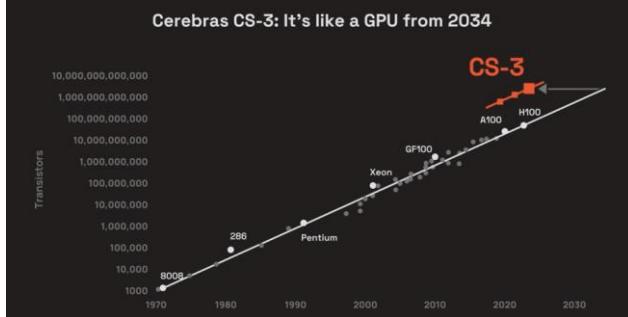
Follow



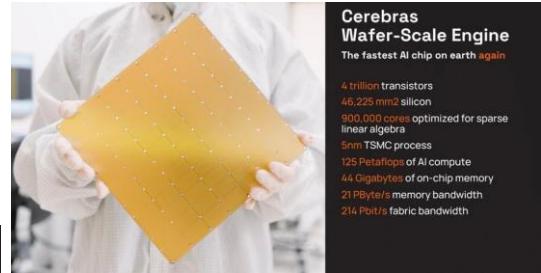
0

Mar 25, 2024, 03:16pm EDT

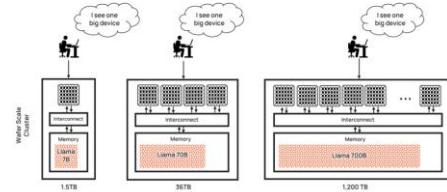
Updated Mar 25, 2024, 06:19pm EDT



AI Supercomputer



And Your Model Always Fits
1B or 1T Parameters



AI

25MAR2024: <https://www.forbes.com/sites/karlfreund/2024/03/25/cerebras-update-the-wafer-scale-engine-3-is-a-door-opener/>

57

AI startup Cerebras debuts 'world's fastest inference' service - with a twist

The AI computer maker claims its inference service is dramatically faster and makes new kinds of 'agentic' AI possible.

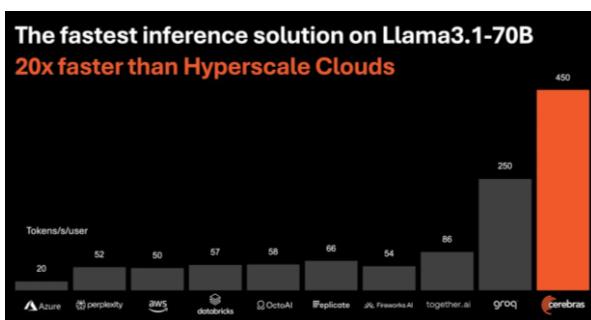


Written by Tiernan Ray, Senior Contributing Writer

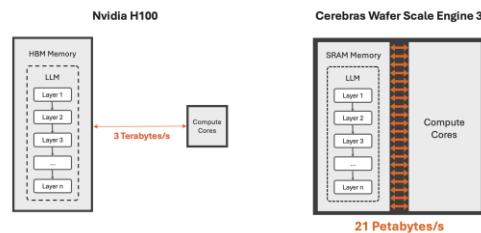
Aug. 27, 2024 at 6:03 p.m. PT

Cerebras vs Nvidia H100
Llama3.1-70B

	Tokens/sec/user	Cost per M tokens ¹
Cerebras	450	\$0.60
Hyperscale H100 Cloud	20	\$2.90 ²
X-factor	22x	1/5



Storing the entire model in on-chip SRAM increases memory bandwidth by 7,000x



AI

27AUG2024: <https://cerebras.ai/blog/introducing-cerebras-inference-ai-at-instant-speed>, <https://www.zdnet.com/article/ai-startup-cerebras-debuts-worlds-fastest-inference-with-a-twist/>

58

CNBC MARKETS BUSINESS INVESTING TECH POLITICS VIDEO INVESTING CLUB JOIN PRO JOIN LIVESTREAM

TECH

Nvidia buying AI chip startup Groq's assets for about \$20 billion in its largest deal on record

PUBLISHED WED, DEC 24 2025 3:54 PM EST | UPDATED FRI, DEC 26 2025 9:14 AM EST

David Faber
@DAVIDFABER

SHARE

KEY POINTS

- Nvidia is making its largest purchase ever, acquiring assets from 9-year-old chip startup Groq for about \$20 billion.
- The company was founded by creators of Google's tensor processing unit, or TPU, which competes with Nvidia for artificial intelligence workloads.
- Groq, which was valued at \$6.9 billion in a financing round in September, framed the deal as a "non-exclusive licensing agreement," with its CEO and other senior leaders joining Nvidia.

**A**24DEC2025: <https://www.cnbc.com/2025/12/24/nvidia-buying-ai-chip-startup-groq-for-about-20-billion-biggest-deal.html>

59

**NOTEBOOKCHECK**

Reviews News Videos Benchmarks / Tech Buyers Guide Magazine Library

The GPU war is over: Nadella says AI's new ceiling is the power grid

Microsoft CEO Satya Nadella confirms the GPU shortage is over. The new AI bottleneck is power: Microsoft has chips sitting in inventory it can't plug in. Why the race for data center grid capacity and Small Modular Reactors (SMRs) is now the real AI power crunch.

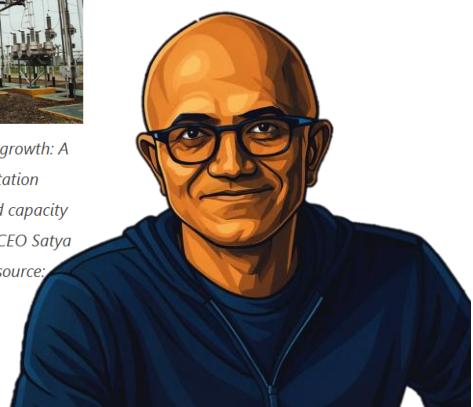
Darryl Linington, Published 12/07/2025 [ES](#) [PT](#) ... [AI](#) [Software](#)

Microsoft CEO Satya Nadella says the AI boom has hit a new kind of ceiling... and it's not GPUs.

Speaking on the [BG2 podcast](#) alongside OpenAI CEO Sam Altman, Nadella said Microsoft is "not chip supply constrained" anymore. The real problem, he explained, is finding enough powered, fully built-out data centres — the "warm shells" close to grid capacity — actually to switch all those accelerators on.



The real bottleneck for AI growth: A typical electrical substation representing the local grid capacity challenges that Microsoft CEO Satya Nadella cited. (Image source: Freepik.com)

**A**1NOV2025: <https://www.youtube.com/watch?v=GnI833wXRz0>, <https://www.notebookcheck.net/The-GPU-war-is-over-Nadella-says-AI-s-new-ceiling-is-the-power-grid.1180290.0.html>

60



DL Frameworks

Which framework to use in our deep learning project

61

Deep Learning Frameworks



The diagram features a large iceberg floating in a blue ocean under a blue sky with white clouds. The visible part of the iceberg above the water surface is labeled with logos for Keras (red 'K'), Pytorch Lightning (purple lightning bolt), TensorFlow (orange 'G' and yellow 'T'), JAX (green 'G' and purple 'A'), and Flax (blue 'C'). Below the waterline, the submerged portion of the iceberg is also labeled with these same logos. This visual metaphor represents that many deep learning frameworks have significant underlying complexity and capabilities that are not immediately apparent from their surface-level interfaces.

Image credit: <https://softwaremill.com/ml-engineer-comparison-of-pytorch-tensorflow-jax-and-flax/>

- Most deep learning frameworks can be used for two different things:
 - Replacing Numpy-like operations with GPU-accelerated operations
 - Building deep neural networks
- **TensorFlow** (before 2.0) used **static computation graphs**. However, **PyTorch** uses **dynamic computation graph**, allowing greater flexibility in building complex architectures.

A graph is created on the fly

```
from torch.autograd import Variable
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

W_h h W_x x

62

Static vs. Dynamic Computation Graphs

- Most **static declaration paradigms** use a two-step process:
 1. **Define** a computational architecture that is represented as a computational graph
 - For example, take a 64x64 image and perform two layers of convolution, predict the class of the image out of 100 classes, and optionally calculate the loss with respect to the correct label
 2. **Execute** the computation
 - For example, a user repeatedly populates the 64x64 matrix with data and the library executes the previously declared computation graph. At test time, the predictions can be used. At training time, the loss can be calculated and back-propagated through the graph to compute the gradients required for parameter updates
- Despite a number of advantages, there are many scenarios where it is difficult for simple static declaration tools to handle.
 - For example, variably sized inputs, variably structured inputs, non-trivial inference algorithms, variably structured outputs
- In **the dynamic declaration paradigm**, a user defines the computation graph programmatically. There are no separate steps for definition and execution: the necessary computation graph is created, on the fly, as the loss calculation is executed, and a new graph is created for each training instance.
- Recommended reading materials that comparatively explain the pros and cons of each paradigm:
 - [15JAN2017] DyNet: The Dynamic Neural Network Toolkit <https://arxiv.org/abs/1701.03980>
 - [11MAR2019] <https://medium.com/analytics-vidhya/dynamic-vs-static-computation-graph-2579d1934ecf>

AI

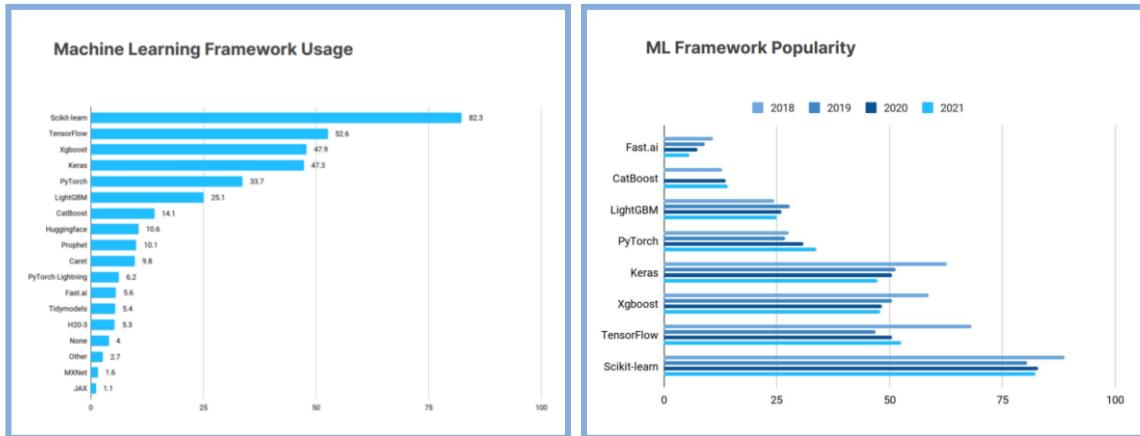
63



64

AI Google TensorFlow vs. Meta PyTorch

- Previously, TensorFlow took the lead when considering the number of users. It has been a go-to framework for deployment-oriented applications.



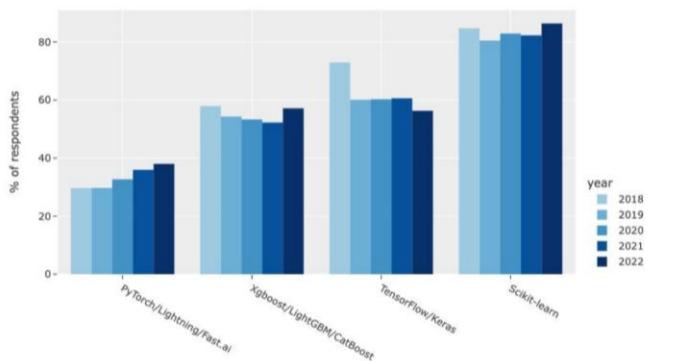
[14OCT2021] Kaggle: <https://www.kaggle.com/kaggle-survey-2021>

65

AI Google TensorFlow vs. Meta PyTorch

Kaggle DS & ML Survey 2022

Scikit-learn is the most popular ML framework while PyTorch has been growing steadily year-over-year



Google Cloud

[11-13OCT2022] Kaggle:
<https://www.kaggle.com/kaggle-survey-2022>

66

AI Google TensorFlow vs. Meta PyTorch

- PyTorch quickly became the preferred choice within research/developer communities.

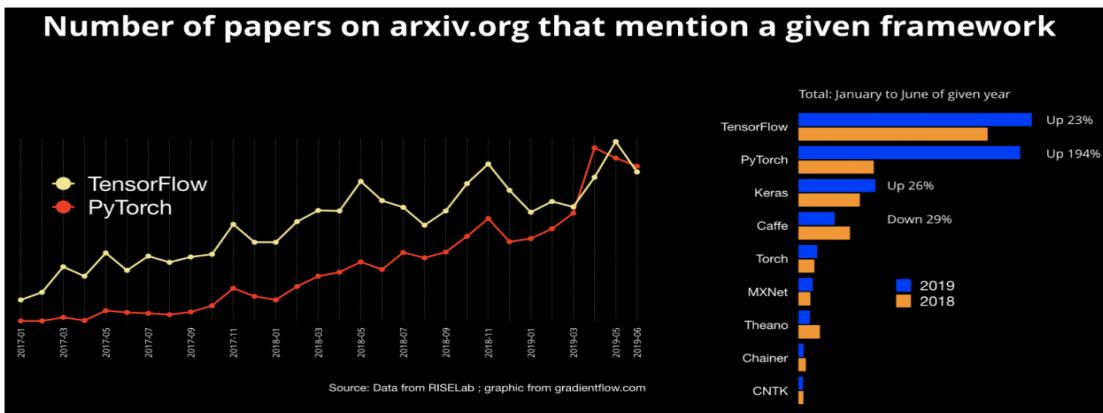


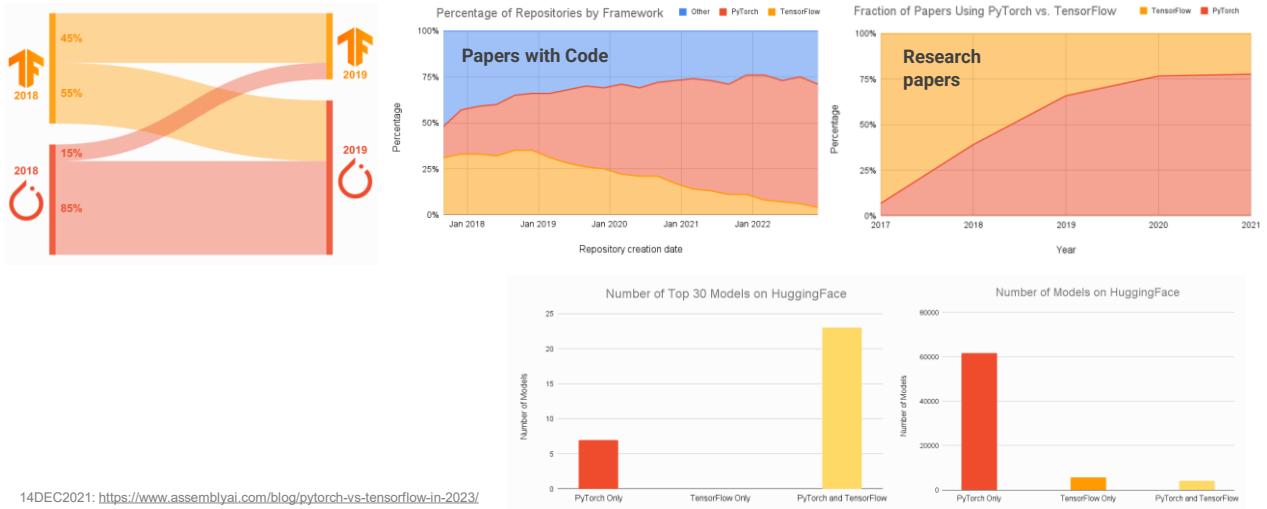
Figure 1. Number of papers posted on arXiv.org that mention each framework. Source: Data from RISELab and graphic by Ben Lorica.

[JULY2019] <https://www.oreilly.com/content/one-simple-graphic-researchers-love-pytorch-and-tensorflow/>

67

AI Google TensorFlow vs. Meta PyTorch

- PyTorch has steadily overshadowed TensorFlow.

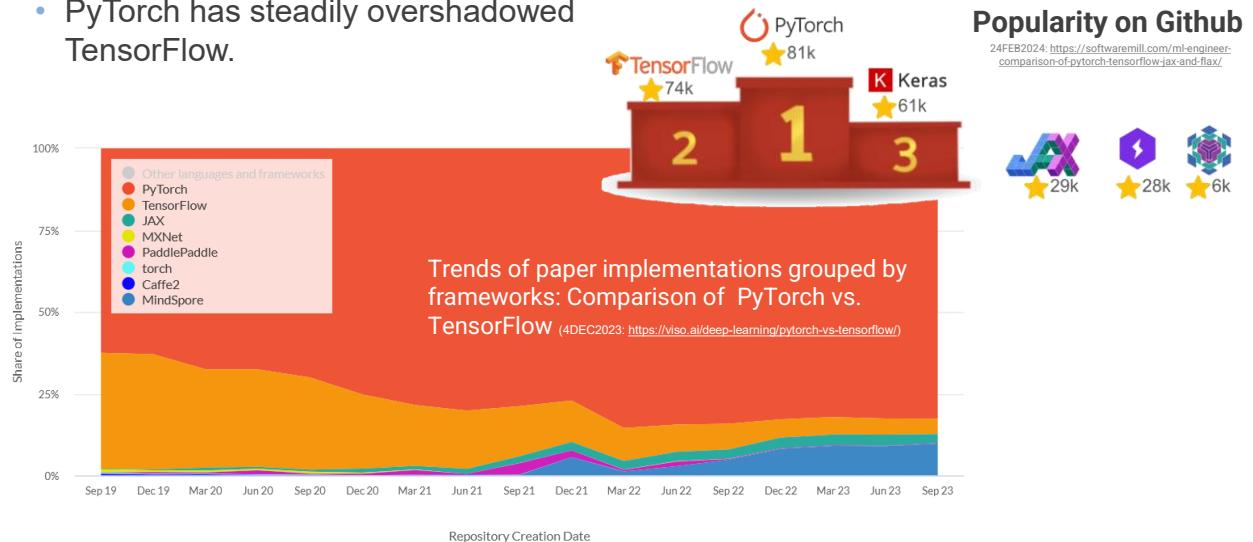


14DEC2021: <https://www.assemblyai.com/blog/pytorch-vs-tensorflow-in-2023/>

68

AI Google TensorFlow vs. Meta PyTorch

- PyTorch has steadily overshadowed TensorFlow.



69

AI Deep Learning Frameworks



since 2015 (with several backends including TF)
Since 6/2019, officially included in TF 2.0
since 2023, Keras 3.0 includes full APIs for TF, JAX, and Torch backends



since February 2024 (stable release)



since 2015
since 6/2019 (TF 2.0), include keras's high-level APIs

```
Simply change from:  
from keras.XXX import YYY  
To:  
from tensorflow.keras.XXX import YYY
```

70



\$ python example.py
Using TensorFlow backend



\$ python example.py
Using PyTorch backend



\$ python example.py
Using JAX backend

```

import keras

model = keras.Sequential([
    keras.layers.Input(shape=(num_features,)),
    keras.layers.Dense(512, activation="relu"),
    keras.layers.Dense(512, activation="relu"),
    keras.layers.Dense(num_classes, activation="softmax"),
])
model.summary()

model.compile(
    optimizer=keras.optimizers.AdamW(learning_rate=1e-3),
    loss=keras.losses.CategoricalCrossentropy(),
    metrics=[
        keras.metrics.CategoricalAccuracy(),
        keras.metrics.AUC(),
    ],
)

history = model.fit(
    x_train, y_train, batch_size=64, epochs=8, validation_split=0.2
)
evaluation_scores = model.evaluate(x_val, y_val, return_dict=True)
predictions = model.predict(x_test)

```



Image credit: "Introducing Keras 3.0," https://keras.io/keras_3/

71













```

import keras
from keras import ops

class TokenAndPositionEmbedding(keras.Layer):
    def __init__(self, max_length, vocab_size, embed_dim):
        super().__init__()
        self.token_embed = self.add_weight(
            shape=(vocab_size, embed_dim),
            initializer="random_uniform",
            trainable=True,
        )
        self.position_embed = self.add_weight(
            shape=(max_length, embed_dim),
            initializer="random_uniform",
            trainable=True,
        )

    def call(self, token_ids):
        # Embed positions
        length = token_ids.shape[-1]
        position_embeddings = np.arange(length, dtype="int32")
        position_vectors = ops.takeself(position_embeddings, axis=0, keepdims=True)
        token_ids = ops.cast(token_ids, dtype="int32")
        token_vectors = ops.takeself(token_ids, token_embed, token_ids, axis=0)
        # Sum
        embed = token_vectors + position_vectors
        # Normalize embeddings
        embed = ops.norm(embed * square_norm(embed), axis=1, keepdims=True)
        return embed / ops.sqrt(ops.maximum(power_sum, 1e-7))

```

```

model = get_keras_core_model()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
loss_fn = torch.nn.CrossEntropyLoss()

def train_step(inputs, targets):
    inputs, targets = inputs.type(torch.FloatTensor), targets.type(torch.LongTensor)
    logits = model(inputs, training=True)
    loss = loss_fn(logits, targets)

    # Compute gradients.
    model.zero_grad()
    loss.backward()
    optimizer.step()

    # Update weights.
    optimizer.step()

    return loss

# Iterate over epochs.
for epoch in range(num_epochs):
    # Iterate over batches of the dataset.
    for step, (inputs, targets) in enumerate(dataset):
        loss = train_step(inputs, targets)
        print(f"loss: {loss.detach().item():.4f}")

```

```

import torch
class TokenAndPositionEmbedding(keras.Layer):
    ...

    def call(self, token_ids):
        length = token_ids.shape[-1]
        position_embeddings = np.arange(length, dtype="int32")
        position_vectors = ops.takeself(position_embeddings, axis=0, keepdims=True)
        token_ids = token_ids.int()
        token_embed = self.get_weights("token_embed")
        token_ids = token_ids.type(token_embed.dtype)
        token_ids = token_ids + 1
        token_ids = token_ids - 1
        token_ids = token_ids * token_embed
        token_ids = token_ids.sum(dim=-1, keepdim=True)
        embed = token_vectors + position_vectors
        # Normalize
        embed = jax.numpy.jnp.sum(embed, axis=1, keepdims=True)
        embed = jax.numpy.jnp.sqrt(jnp.sum(embed**2, axis=1, keepdims=True))
        embed = embed / (jax.numpy.jnp.sqrt(jnp.sum(embed**2, axis=1, keepdims=True)))
        return embed / (jax.numpy.jnp.sqrt(jnp.sum(embed**2, axis=1, keepdims=True)))

```

```

model = get_keras_core_model()
optimizer = keras.optimizers.Adam(learning_rate=1e-3)
loss_fn = keras.losses.CategoricalCrossentropy(from_logits=True)

# ALL variables must be built before training starts.
optimizer.build(model.trainable_variables)

def compute_loss_and_updates(trainable_vars, non_trainable_vars, data):
    # Stateless function to compute the loss and non-trainable var updates.
    X, y = data
    pred, non_trainable_vars = model.stateless.call(trainable_vars, non_trainable_vars, x)
    loss = loss_fn(pred, y)
    return loss, non_trainable_vars

# Compute gradients for the trainable vars.
grad_fn = jax.vjp_and_grad(lambda x: compute_loss_and_updates(x, non_trainable_vars))

@jax.jit
def train_step(state, data):
    # Stateless function that calls the grad_fn and computes trainable var updates.
    trainable_vars, non_trainable_vars, optimizer_vars = state
    print(f"trainable_vars: {trainable_vars}, non_trainable_vars: {non_trainable_vars}, optimizer_vars: {optimizer_vars}")
    (loss, state) = grad_fn(data)
    (loss, state) = jax.lax.stateless_update(stateless_fn, apply_fn, optimizer_vars, grads, trainable_vars)
    return loss, (trainable_vars, non_trainable_vars, optimizer_vars)

# Prepare model state.
state = model.trainable_variables, model.non_trainable_variables, optimizer.variables

# Iterate over epochs.
for epoch in range(num_epochs):
    # Iterate over batches of the dataset.
    for step, (inputs, targets) in enumerate(dataset):
        # Each train_step call is entirely stateless (no side effects).
        loss, state = train_step(state, data)
        print(f"loss: {loss.item():.4f}")

```

Writing a custom training loop for a Keras model:
PyTorch,
TensorFlow,
JAX

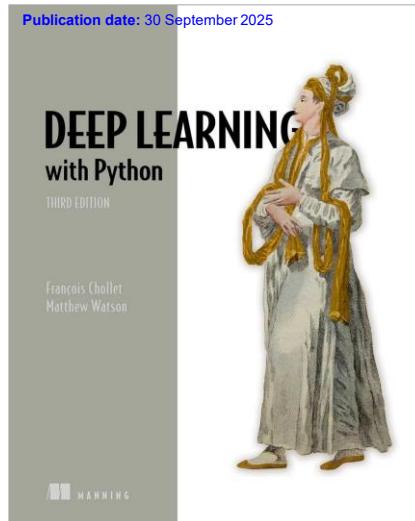
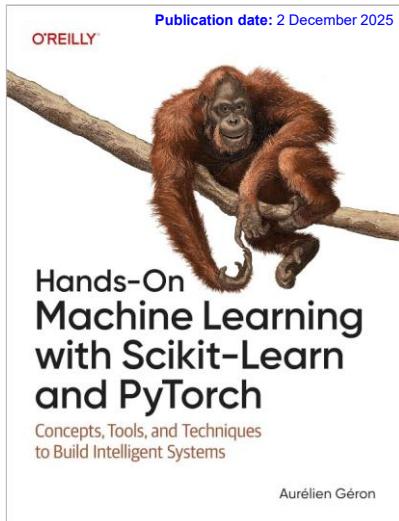
...
or use
your framework of choice
for backend-specific
components



Image credit: "Introducing Keras 3.0," https://keras.io/keras_3/

72

AI Google TensorFlow vs. Meta PyTorch



73

AI DL frameworks in this course



74



75

Tensor in Deep Learning

- **Tensor** can be thought as a general term for nd -array where n can be 0, 1, 2, 3,

Scalar Vector Matrix Tensor



How many indexes required?	Computer Science	Mathematics	
0	Number	Scalar	Rank-0 tensor (0 dimension)
1	Array	Vector	Rank-1 tensor (1 dimension)
2	2D-array	Matrix	Rank-2 tensor (2 dimensions)
n	nd-array	nd-tensor	Rank-n tensor (n dimensions)

76



FLOPS vs. FLOPs

- One **floating-point operation** refers to one multiplication, one division, one addition, and so on.
- **FLOPS (FLoating-point OPerations per Second)** is a unit of speed that often uses to measure computing hardware performance.
 - For example, NVIDIA RTX 2080 has a theoretical performance in FP32 mode of 10.07 TFLOPS. This means the theoretical maximum number of floating-point operations that the hardware might be capable of (if we are extremely lucky).
 - For modern CPUs, **FLOPS** is calculated from repeated uses of a “fused multiply then add” instruction, so that one instruction counts as two floating point operations.
- **FLOPs (FLoating-point OPerations)** is a unit of amount that often uses to describe how many operations are required to run a single instance of a given model.

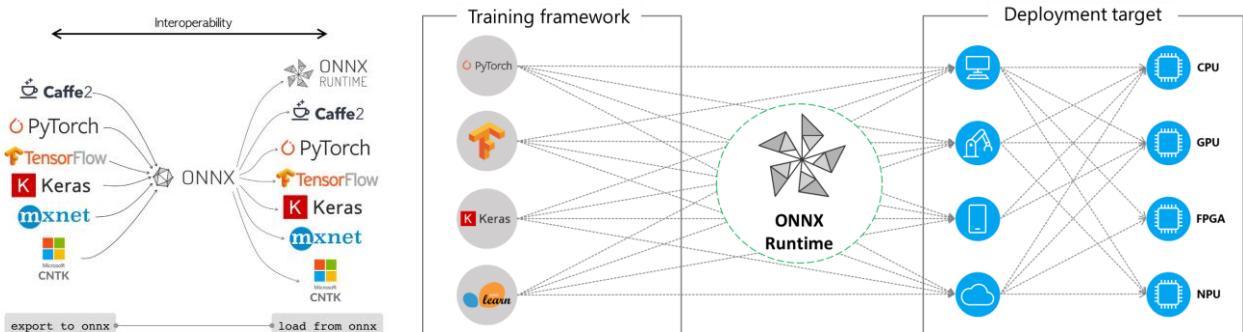
Model	Input Size	Param Size	Flops
AlexNet	227 x 227	233 MB	727 MFLOPs
CaffeNet	224 x 224	233 MB	724 MFLOPs
VGG-VD-16	224 x 224	528 MB	16 GFLOPs
VGG-VD-19	224 x 224	548 MB	20 GFLOPs
GoogleNet	224 x 224	51 MB	2 GFLOPs
ResNet-34	224 x 224	83 MB	4 GFLOPs
ResNet-152	224 x 224	230 MB	11 GFLOPs
SENet	224 x 224	440 MB	21 GFLOPs

77



ONNX , ONNX Runtime

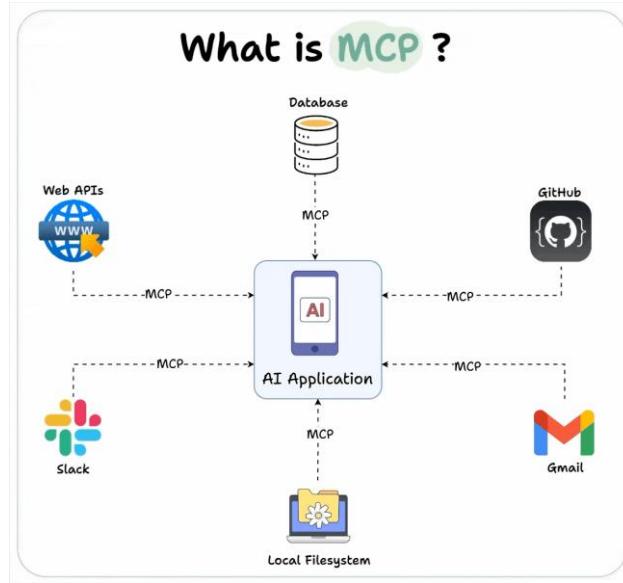
- Open Neural Network Exchange (**ONNX**)



Images from: <https://www.aurigait.com/blog/onnx-onnx-runtime-and-tensorrt/> , <https://onnxruntime.ai/docs/execution-providers/>

78

AI Model Context Protocol

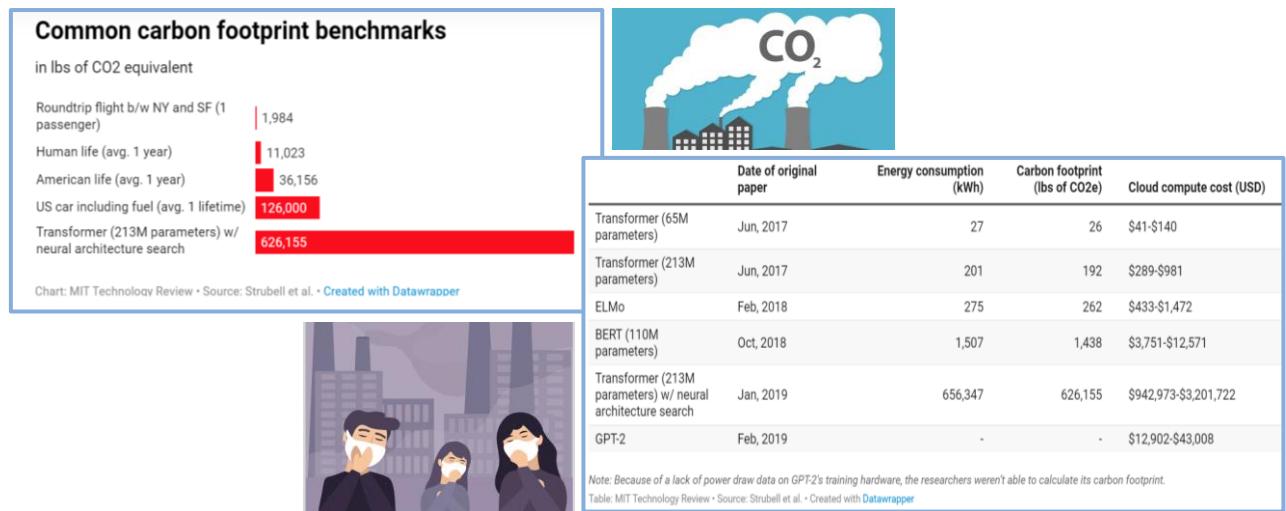


- **MCP** is like a universal connector for AI applications, allowing different AI tools and models to interact seamlessly with various data sources.
- **MCP** is an open protocol that standardizes how applications provide context to LLMs.
- **MCP** was originally built to improve Claude's ability to interact with external systems. Anthropic decided to open-source MCP in early 2024 to encourage industry-wide adoption.
- Read more <https://www.ibm.com/think/topics/model-context-protocol>

Image credit: (29JUL2025) <https://abhishek-iit.medium.com/model-context-protocol-mcp-assets-8cae94599875>

79

AI Carbon Emissions Problem



(6JUN2019) MIT Technology Review: <https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>

80

AI Carbon Emissions Problem

- In the DL era, the computational resources needed to produce a best-in-class AI model has on average doubled every 3.4 months; this translates to a 300,000x increase between 2012 and 2018.
- A study in 2019 estimated that:
 - Training a single DL model (a particularly energy-intensive model) can generate up to 626,155 pounds of CO₂ emissions—roughly equal to the total lifetime carbon footprint of five cars. As a point of comparison, the average American generates 36,156 pounds of CO₂ emission in a year.
 - Training an average-sized model (much smaller than GPT) and examining not just the energy required to train the final version, but the total number of trial runs that went into producing the final version:
 - Over the course of six months, 4,789 different versions of the model were trained, requiring 9,998 total days' worth of GPU time (more than 27 years).
 - Taking all these runs into account, the researchers estimated that building this model generated over 78,000 pounds of CO₂ emission in total—more than the average American adult will produce in two years.

(17JUN2020): <https://www.forbes.com/sites/robtoews/2020/06/17/deep-learning-climate-change-problem/?sh=7120e1176b43>

81

AI Carbon Emissions Problem

- Deploying AI models to take action in real-world settings—a process known as inference—comes even more energy than training does. NVIDIA estimates that 80-90% of the cost of a neural network is in inference rather than training.
- Machine Learning Emissions Calculator <https://mlco2.github.io/impact/#compute>

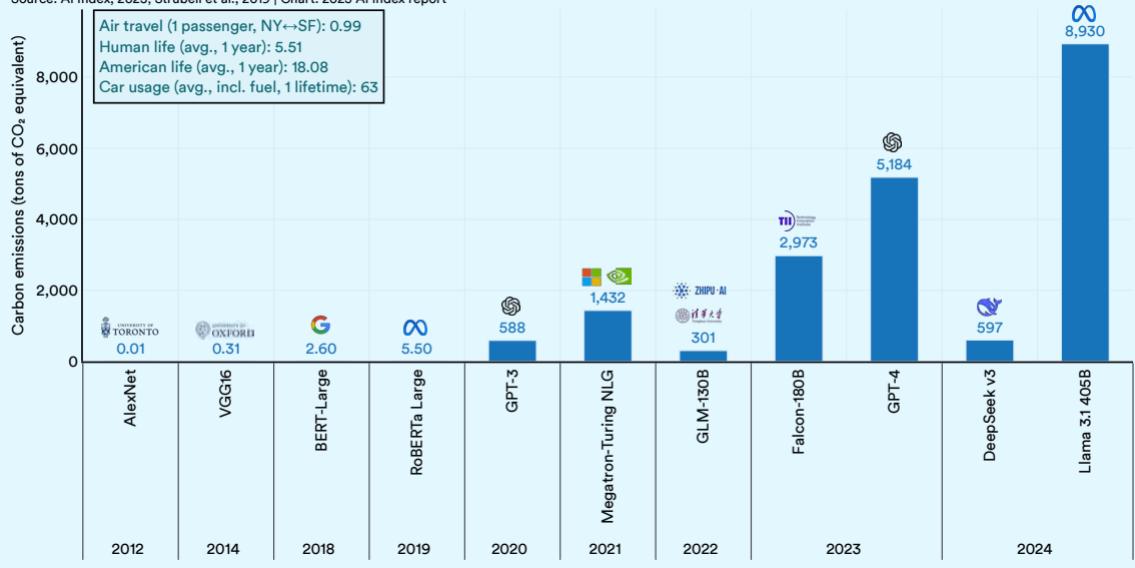


(17JUN2020): <https://www.forbes.com/sites/robtoews/2020/06/17/deep-learning-climate-change-problem/?sh=7120e1176b43>

82

Estimated carbon emissions from training select AI models and real-life activities, 2012–24

Source: AI Index, 2025; Strubell et al., 2019 | Chart: 2025 AI Index report



HAI AI Index Report 2025 (April 2025), https://hai.stanford.edu/assets/files/hai_ai_index_report_2025.pdf

83



84

Next class, a ready-to-use GPU computer

- Cloud
 - Google Colab
- Local installation
 - NVIDIA driver, CUDA toolkit, cuDNN
 - <https://pytorch.org/get-started/locally/>
- NVIDIA's Docker containers
 - <https://docs.nvidia.com/deeplearning/frameworks/user-guide/index.html>



AI

85

Thank You



86