The background of the slide is a dense field of 3D-rendered numbers in various shades of blue. The numbers are of different sizes and are scattered across the frame, creating a sense of depth and complexity. Some numbers are larger and more prominent, while others are smaller and recede into the background. The overall effect is a vibrant, data-driven aesthetic.

LLMs and Advanced GPU Course

Chap 1- Introduction

Electrical Engineering department of
Amirkabir University of technology

Dr. Mohammadreza Pourfard

February 2026

صفحات درس و استاد



Ssail.aut.ac.ir



اتاق: طبقه سوم ساختمان ابوریحان، دانشکده مهندسی برق، آزمایشگاه سیستمهای هوشمند دیجیتال، پزشکی شخصی، هوش مصنوعی و GPU

@Electronic_daneshbonyan



pourfardm@gmail.com



+982164543373



http://ssail.aut.ac.ir



- [Home](#)
- [Research Groups](#)
- [Research Projects](#)
- [People](#)
- [News](#)
- [About Us](#)
- [Admission](#)

Research Groups

In today's rapidly evolving technological landscape, numerous research groups are at the forefront of innovation, each focusing on transformative areas that have the potential to reshape industries and enhance our daily lives. Among these, the Internet of Things (IoT) research group is dedicated to exploring the interconnectedness of devices and sensors, striving to create smart environments that foster efficiency and convenience. Meanwhile, our Artificial Intelligence (AI) team is delving into advanced algorithms and machine learning techniques to develop systems that can learn, adapt, and make informed decisions, pushing the boundaries of automation and intelligence. Additionally, the Big Data research group is



Laboratory page



Profile Contents

Home

Courses

- About
- Theses
- Research Interests
- Research Groups
- Research Projects
- Current Students
- Employment Records
- Contact
- Publications
- News

About

Adjunct Professor of Amirkabir University of Technology

Theses

B.S Degree:

Simulating Multi-Level Cache Memory

M.S Degree:

Detection and Tracking of a Human in different positions of its body

PhD Degree:

Texture Analysis and Separation for Characterization of Material's Structure through

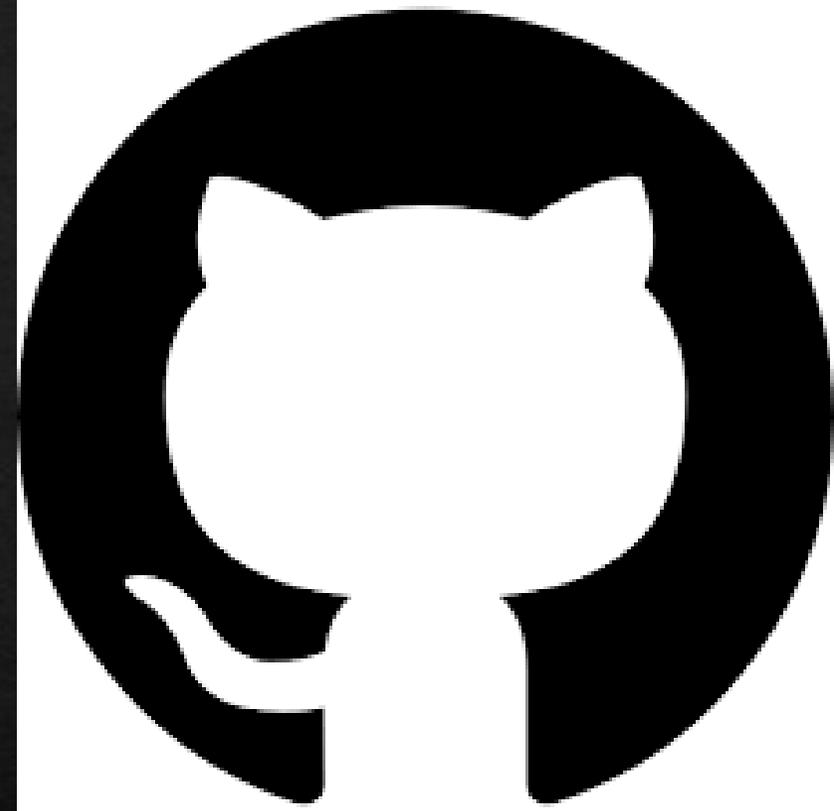
Research Interests

- Genetics data processing,

کانالها و گروه‌های مرتبط با درس

ADVANCED CODING

- AI -



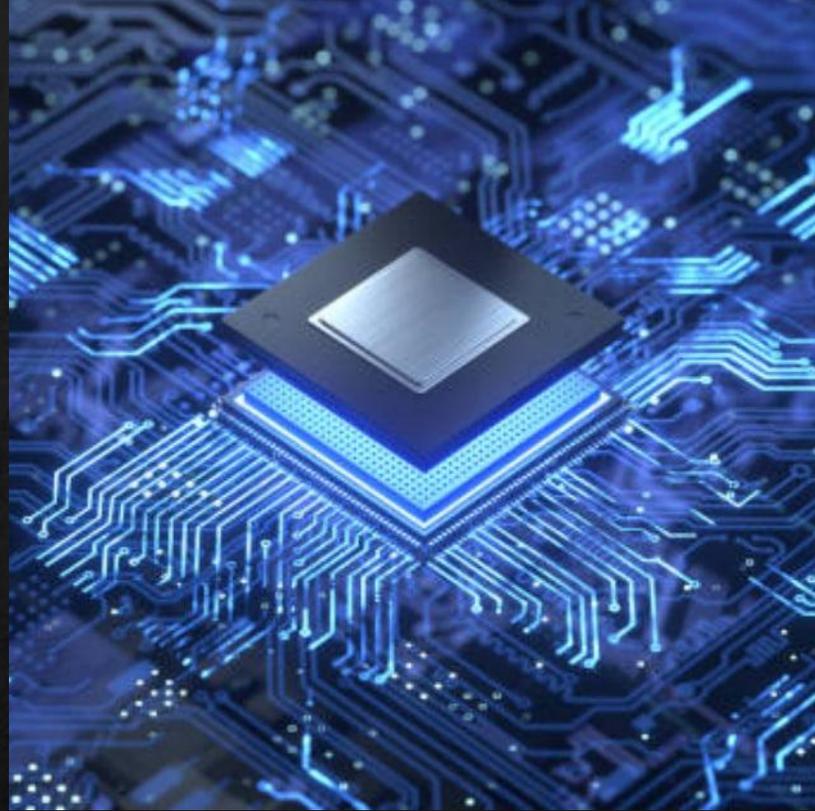
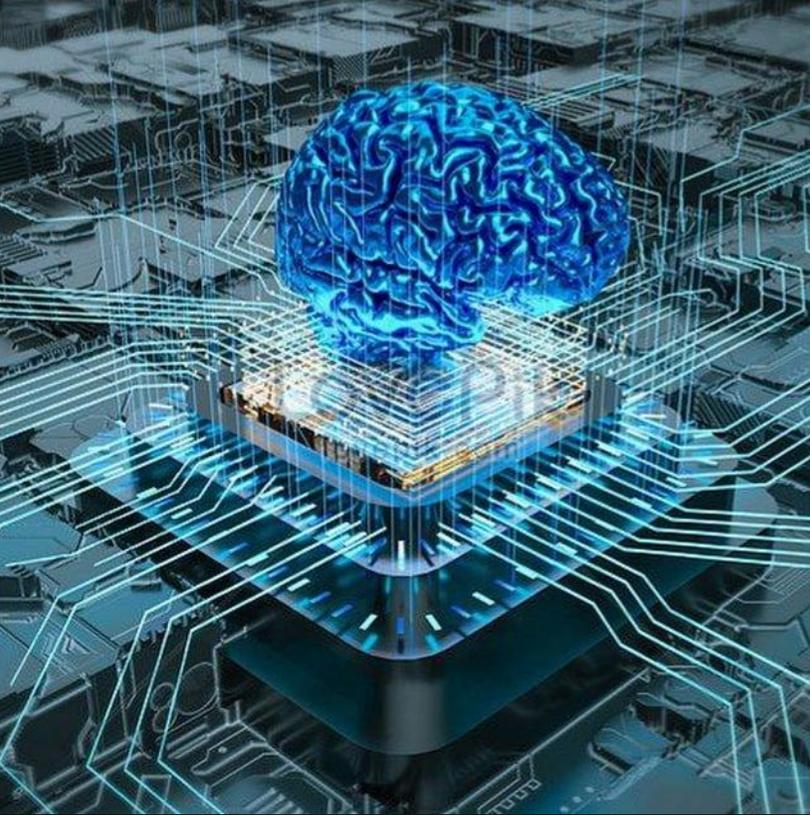
Telegram Channel:
https://t.me/Advanced_programming_algorithm_c



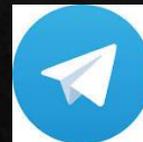
Telegram Group:
https://t.me/Advanced_Programming_Algorithm



GitHub Page:
<https://github.com/Pourfardm/>



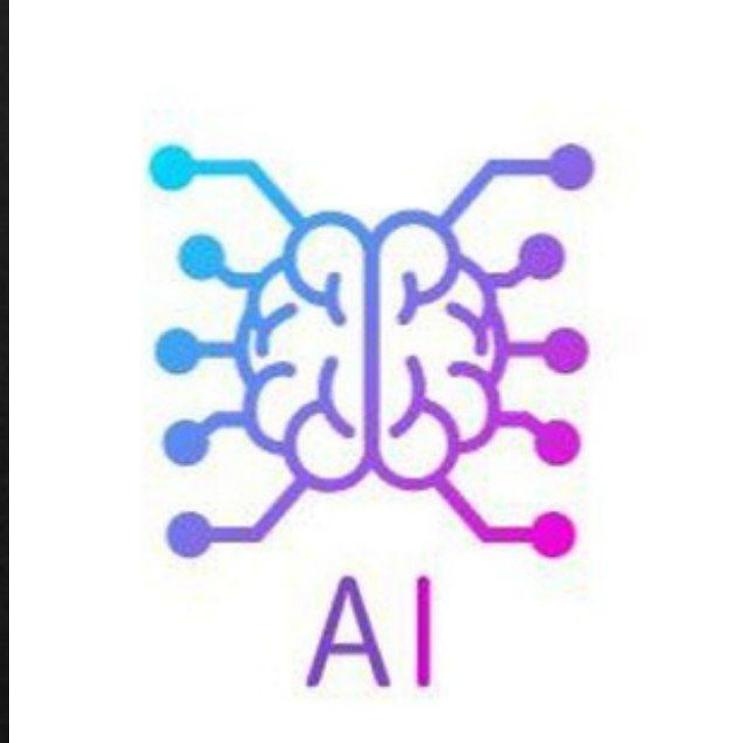
Telegram Channel:
[https://t.me/fpga_digital_logic_d
esign](https://t.me/fpga_digital_logic_design)



Telegram Group:
[https://t.me/Advanced_digital_a
nalog_design](https://t.me/Advanced_digital_analog_design)



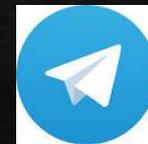
YouTube Channel:
[https://www.youtube.com/chan
nel/UCa8ZFUxp37Vy-
KEbx7SbbHQ](https://www.youtube.com/channel/UCa8ZFUxp37Vy-KEbx7SbbHQ)



Telegram group:
https://t.me/AI_Bigdata_amirkabir_group



Telegram Channel:
https://t.me/Deeplearning_BigData_AI



Private Telegram Group:
For class student



Telegram Channel:

https://t.me/AI_Neuroscience_Genomic_Data

YouTube Channel:

<https://www.youtube.com/channel/UCa8ZFUxp37Vy-KEbx7SbbHQ>

مروی بر نسل دانشگاهها در جهان
اهمیت وروده دانشگاه نسل سوم و
ارایه درس تناسب با دانشگاه نسل سوم

نسل های مختلف دانشگاهها در جهان

عنوان	نسل اول	نسل دوم	نسل سوم	نسل چهارم
هدف	کاشفیت از حقیقت عالم	توسعه فناوری تحقیق و پژوهش	تولید ثروت و ارزش افزوده	ارزش افزوده قابل توجه ایجاد GDP قابل توجه نسبت به GDP کشور
تاکید	آموزش	آموزش پژوهش محور	پژوهش نیاز محور	توجه به نیازهای کلیدی و مهم کشور
شاخص ارزیابی	کیفیت تدریس ترجمه کتاب	مقالات ISI (امتیاز بالا بدون سقف در آیین نامه ارتقای اساتید)	حجم قرارداد دانشگاه با صنعت درآمد زایی دانشگاه ثبت پتنت با ارجاعات بالا تالیف کتاب	نسبت دانشگاه با شرکتهای بزرگ خصوصی با سهم قابل توجه در اقتصاد کشور
مشکلات	عدم آشنایی با مرز علم	پراکنده کاری شناسایی موضوعات مقاله خور تغییر مرتب موضوعات عدم توجه به نیازهای کشور	عدم اثرگذاری قابل توجه در اقتصاد ملی	مخالفت کشورهای غیرهمسو با توسعه صادرات کشورهای رقیب

دانشگاهها در ایران عمدتاً نسل هستند.

نکته مهم این درس

به دلیل رویکرد دانشگاه نسل سوم

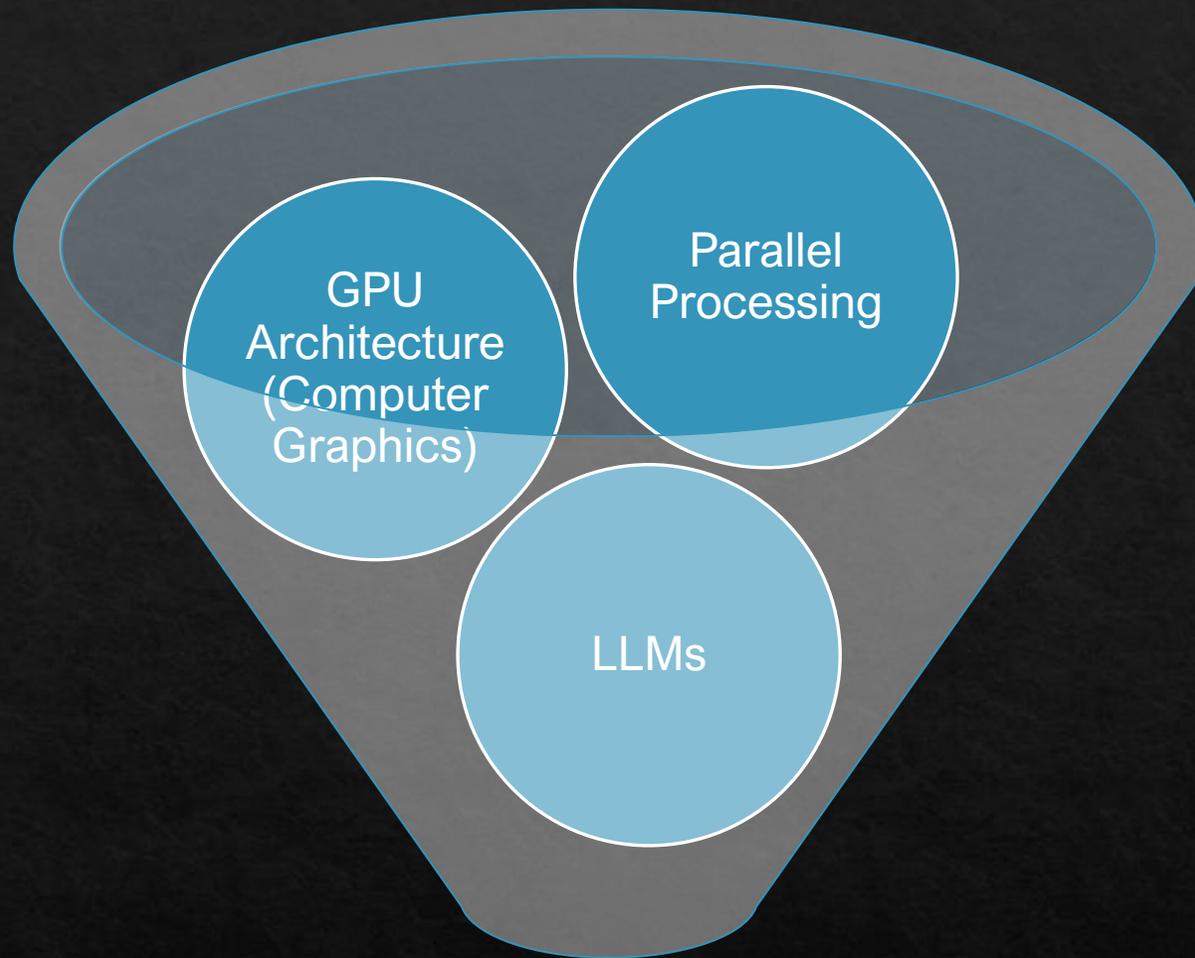
مركز این درس فقط روی الگوریتم نیست

اهداف درس،

اهداف کلیدی درس

- ❖ آشنایی با برنامه نویسی موازی و زبان برنامه نویسی CUDA
- ❖ آشنایی با معماری داخلی آخرین نسل GPU کمپانی Nvidia (Blackwell 2024)
- ❖ آشنایی با آخرین مدل‌های هوش مصنوعی در حوزه LLM
- ❖ فهم ارتباط آخرین دستاوردهای پردازنده گرافیکی با آخرین مدل‌های هوش مصنوعی
- ❖ آشنایی با کاربردهای مفید هوش مصنوعی و LLM به منظور راه اندازی کسب و کار
- ❖ آشنایی با ابزارهای مورد استفاده در شرکت‌های دانش بنیان برای هوش مصنوعی
- ❖ بررسی مدل‌های LLM به صورت عملی بر روی داده‌های مختلف (شامل داده‌های پزشکی)

محتوی درس



This course

جامعہ مخاطب درس



بارم بندی

بارم بندی درس (۲۴ نمره)

◆ نمرات اصلی

◆ پایان ترم ۸ نمره

◆ میان ترم -نداریم

◆ تکالیف ۶ نمره

◆ پروژه پایانی ۴ نمره

◆ ارایه مقاله ۲ نمره

◆ نمرات امتیازی

◆ تکالیف اختیاری ۲ نمره

◆ ارزیابی استاد از دانشجو ۱ نمره اضافی (به صورت استثنا)

◆ تولید محتوی ۱ نمره اضافی

سیلابس درس

بخش اول) بخش پردازش موازی و معماری GPU های نسل آخر NVidia (۳۳ درصد)

◇ الف) سیستم های پردازش موازی

◇ آشنایی با معماری سیستم های چند هسته ای

◇ آشنایی با برنامه نویسی چند نخ، مدل ها و زبانهای برنامه نویسی مرتبط با آن

◇ آشنایی با مفاهیم پردازش برداری، SIMD، SSE، AVX و نحوه استفاده از آن

◇ پیاده سازی الگوریتم ها به صورت چند نخ و برداری با استفاده از زبان های برنامه نویسی چند هسته ای (OpenMP)

◇ آشنایی با روشهای متداول همگام سازی نخ، قفل، مانع

بخش اول) بخش پردازش موازی و معماری GPU های نسل آخر NVidia (۳۳ درصد)

- ◇ آشنایی با معماری پردازنده‌های گرافیکی، سلسله مراتب حافظه در GPU
- آشنایی با تاریخچه معماری های GPU کمپانی Nvidia شامل Fahrenheit, Kelvin Turing, Hopper, Blackwell, Blackwell Ultra
- آشنایی با بلوک های چیپ گرافیکی CB202 Chipset
- آشنایی با واحدهای GPC, Memory Controller, Cache, AMP & Giga Thread Engine, NVENC / NVDEC, Optical Flow Engine, PCI Express 5.0 Host Interface.
- بررسی ماژول Mixed FP32/INT32 module
- آشنایی با معماری و قابلیت های هسته های پردازشی Ray tracing core, Cude core, Tensor core
- آشنایی با واحد های Warp Schedulers & Dispatch Units
- آشنایی با واحد Texture Units: perform texture fetches and filtering
- آشنایی با بخش Load/Store Units (LD/ST): handle memory access (global/local memory 10-reads/writes)
- آشنایی با بخش Register File: private storage per thread
- آشنایی با حافظه های سریع Shared Memory / L1 Cache: memory accessible by all threads in the SM. و پیاده سازی سریع الگوریتم ها به کمک این حافظه ها
- آشنایی با واحد پردازش هوش مصنوعی AI Management Processor (AI kernels, Multi-GPU Scaling, tensor parallelism, Training massive models (LLMs), data parallelism) و نحوه استفاده از آنها در پیاده سازی مدل‌های LLM

بخش اول) بخش پردازش موازی و معماری GPU های نسل آخر NVidia (۳۳ درصد)

آشنایی با معماری CUDA و GPU Driver

User-Space Driver of GPU (Resource & State Management (Logical View), آشنایی با بخش
Build Command Buffers, Interface with Kernel Driver, API Front-End, State Manager,
Compiler Stack, Resource Manager, Command Buffer Builder, Caching & Pipeline
Database)

Kernel-Space Driver of GPU (Context Manager, Memory Manager, Command Processor Interface, Scheduler / Dispatcher, Synchronization Manager, Interrupt Handler, Power & Thermal Control, Virtualization Layer) آشنایی با بخش

بخش اول) بخش پردازش موازی و معماری GPU های نسل آخر NVidia (۳۳ درصد)

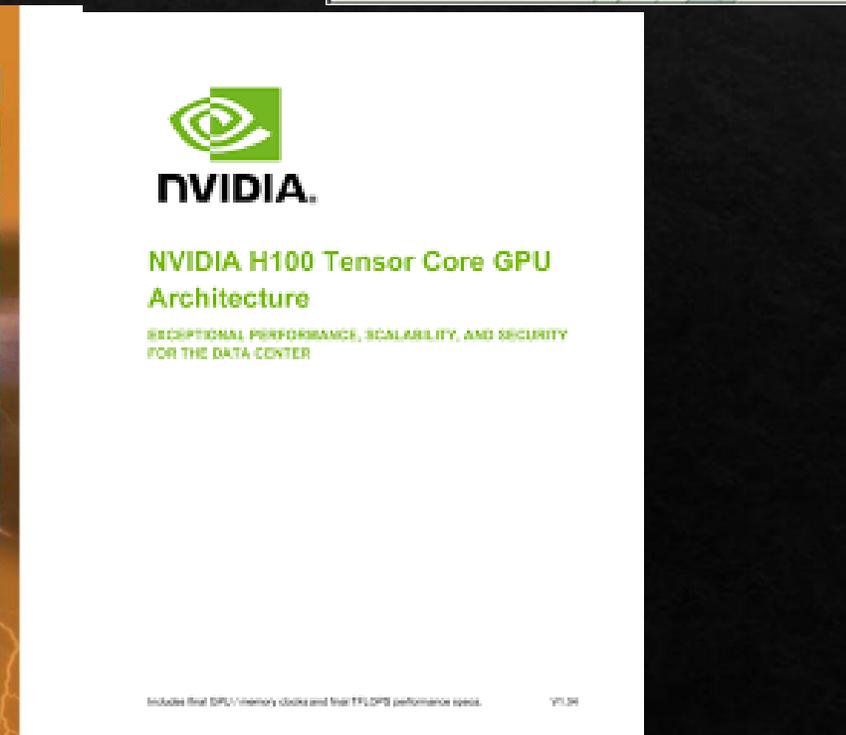
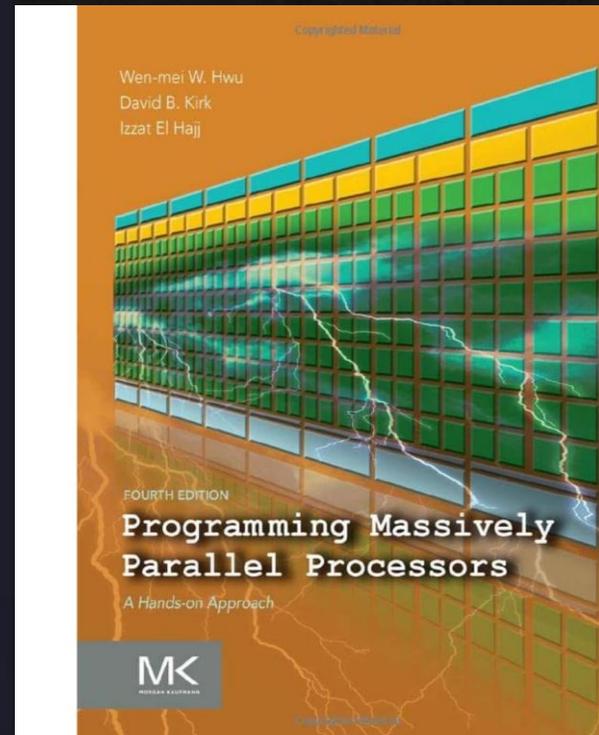
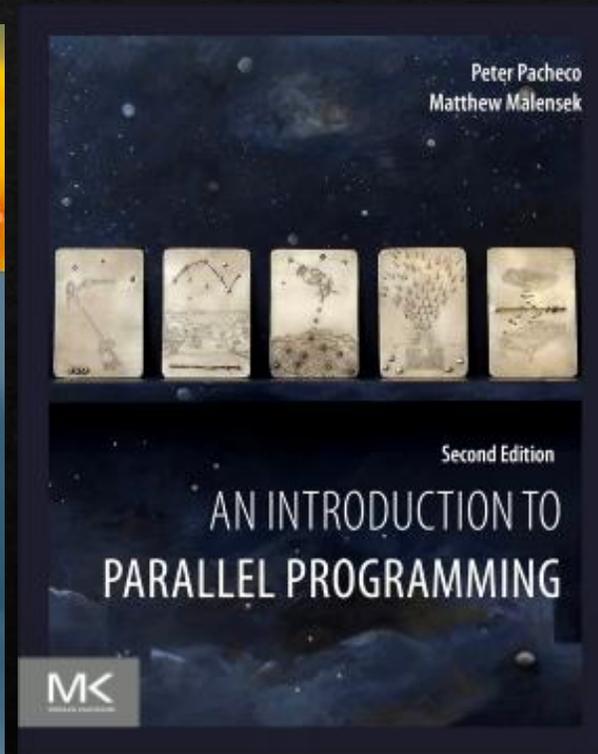
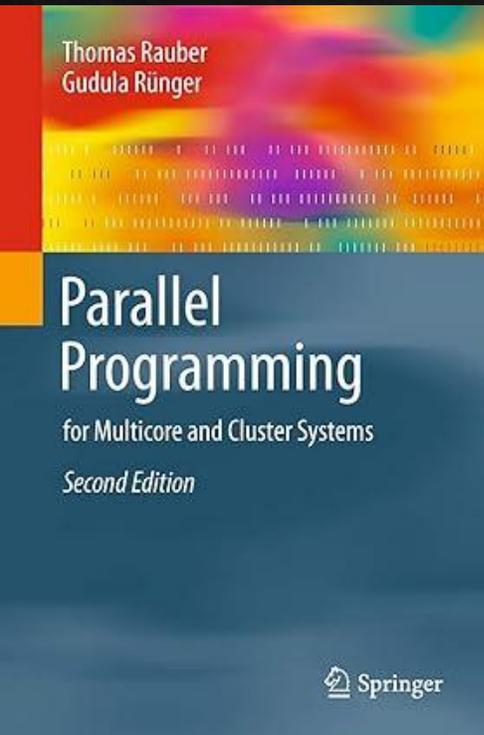
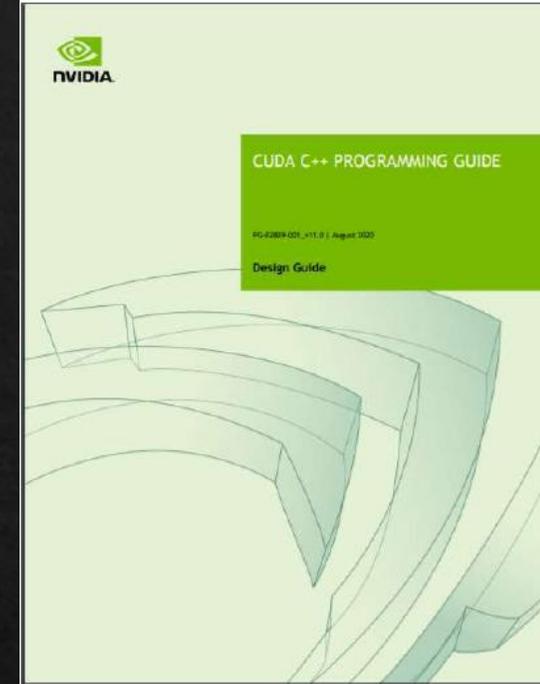
◇ آشنایی با برنامه‌نویسی موازی GPU و زبان برنامه‌نویسی CUDA

◇ ارائه مثال‌هایی از پیاده‌سازی کاربردهای معمول در GPU

◇ برنامه‌نویسی با CUDA

◇ تغییر در کرنل‌های CUDA برای یک منظور خاص همانند FFT با تعداد نقاط دلخواه

مراجع این بخش



بخش دوم) شبکه های LLM (۳۴ درصد)

◇ بررسی شبکه های LLM و آخرین پیشرفت ها در این حوزه

◇ آشنایی با مبانی Large Language Modeling

◇ آشنایی با Tokenization و اثرات پنهان آن

◇ آشنایی با مدل های پایه GPT، BERT، T5

◇ آشنایی با Reasoning در LLM ها

◇ آشنایی با انواع Hallucination و عدم قطعیت

◇ آشنایی با انواع روشهای RAG و اتصال به دانش خارجی (همانند Naïve RAG, Advanced RAG, Modular RAG)

◇ ارزیابی و بحران Benchmark

◇ بررسی معماری آخرین نسل Llama, GPT, Deep Seek و سایر مدل های بروز جهانی

بخش دوم) شبکه های LLM (۳۴ درصد)

◇ مراجع درس در این بخش، مقالات بروز جهانی از سال ۲۰۲۲ به بعد خواهد بود.

Gao, Yunfan, et al. "Retrieval-augmented generation for large language models: A survey." *arXiv preprint arXiv:2312.10997* 2.1 (2023). ◇

Hoffmann, Jordan, et al. "Training compute-optimal large language models." *arXiv preprint arXiv:2203.15556* (2022). ◇

Wei, Jason, et al. "Emergent abilities of large language models." *arXiv preprint arXiv:2206.07682* (2022). ◇

Ji, Ziwei, et al. "Survey of hallucination in natural language generation." *ACM computing surveys* 55.12 (2023): 1-38. ◇

◇ اسلاید های درس دانشگاهی CMU – 11-667: Large Language Models

◇ اسلاید های درس دانشگاهی Stanford – CS324 / CME 295

◇ اسلاید های درس دانشگاهی Princeton – COS 597G

بخش سوم) کاربرد LLM و GPU بر روی پردازش داده های مختلف (شامل داده های پزشکی) (۳۳ درصد)

- ◆ پیاده سازی روشهای LLM به کمک GPU بر روی انواع مختلفی از داده ها
- ◆ نحوه استفاده از مفهوم RAG در استفاده از پایگاه های داده متغیر با زمان ژنتیک در مدل های LLM
- ◆ نحوه استفاده از واحدهای AI Processor , AI Kernel در پردازش داده
- ◆ پیاده سازی برخی از الگوریتم های سنگین بر روی GPU های نسل آخر کمپانی Nvidia با استفاده از حافظه سریع L1 Cache memory.
- ◆ آشنایی با فرمت تصاویر پزشکی همانند MRI, CT-scan, Pet-scan, Radiology, Mammography و نقش هر یک از این داده ها در فهم مسایل پزشکی
- ◆ آشنایی با فرمت داده های سیگنال های پزشکی همانند نوار مغز، نواز قلب، نوار عضله و عصب و نقش هر یک از این داده ها در فهم مسایل پزشکی
- ◆ آشنایی با مفاهیم پایه ژنتیک همانند DNA, RNA, Protein و Gut Microbiota
- ◆ آشنایی با سطوح تحلیل داده های پزشکی شامل ژنومیکس، اپی ژنومیکس، ترنسکریپتومیکس، پرتئومیکس، متابولومیکس، فنومیکس
- ◆ نحوه Fine-tuning مدل های LLM برای تحلیل داده های مختلف Multi-omics
- ◆ نحوه طراحی لایه های Encoder , Decoder در مدل های LLM بر روی تصاویر، سیگنال و ژنتیک
- ◆ آشنایی با پایگاههای داده ژنتیک در پزشکی

بخش سوم) کاربرد LLM و GPU بر روی پردازش داده های مختلف (شامل داده های پزشکی) (۳۳ درصد)

◇ مراجع درس در این بخش، مقالات بروز جهانی از سال ۲۰۲۴ به بعد خواهد بود.

- 1- Sarumi, Oluwafemi A., and Dominik Heider. "Large language models and their applications in bioinformatics." *Computational and Structural Biotechnology Journal* 23 (2024): 3498-3505. ◇
- 2- Van, Minh-Hao, Prateek Verma, and Xintao Wu. "On large visual language models for medical imaging analysis: An empirical study." *2024 IEEE/ACM Conference on Connected Health: Applications, Systems and Engineering Technologies (CHASE)*. IEEE, 2024. ◇
- 3- Kao, Jyun-Ping, and Huan-Tang Kao. "Large Language Models in radiology: A technical and clinical perspective." *European Journal of Radiology Artificial Intelligence* (2025): 100021. ◇

موضوعات تکالیف

- ۱- پیاده سازی چند الگوریتم به صورت موازی بر روی بستر GPU با استفاده از CUDA
- ۲- استفاده از Cache L1 از GPU برای پردازش سریع داده های ژنتیک
- ۳- انجام Fine tuning یک شبکه LLM متن باز
- ۴- پیاده سازی یک سیستم RAG
- ۵- طراحی لایه های انکودری و دیکدري LLM برای روی داده های پزشکی

استفاده از ۳۷ گیگابایت رم از ۴۰ گیگابایت رم A100 GPU فقط برای پردازش اطلاعات یک فرد

NVIDIA GH200

Built for the new era
of AI supercomputing

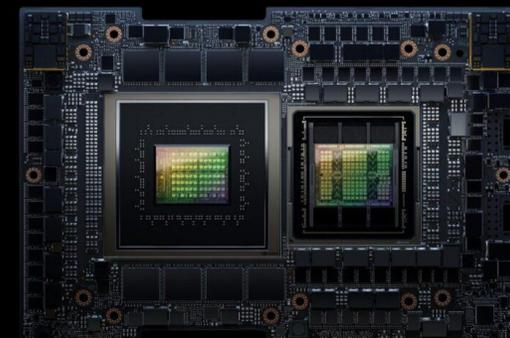
CPU to GPU Bandwidth
900 GB/s
NVLink-C2C

Memory Bandwidth
4.9 TB/s
HBM3e per GPU

Energy Efficiency
1.9X
Performance vs H100

QFT Quantum Simulation
90X
Performance vs dual x86 CPU

RAG LLM Inference
100X
Performance vs dual x86 CPU



624 GB High-Speed Memory | 4.9 TB/s | 4 PF AI Perf | 72 Arm Cores



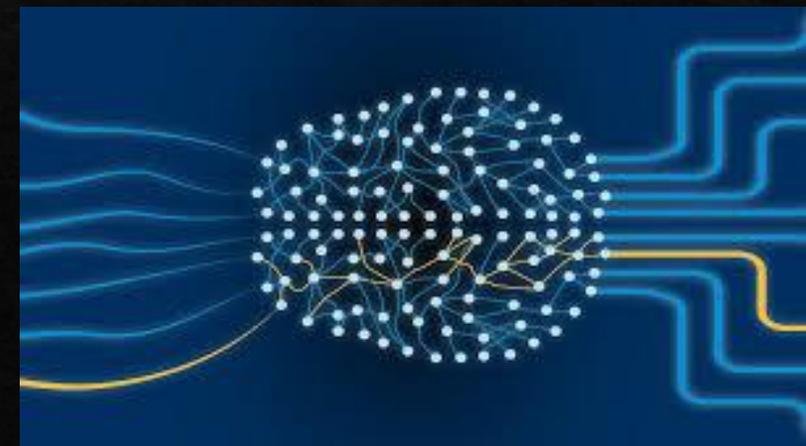
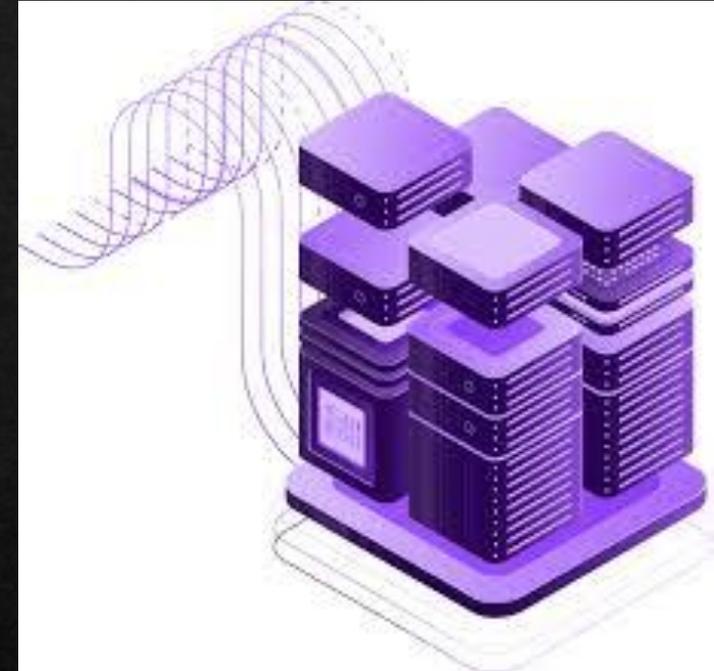
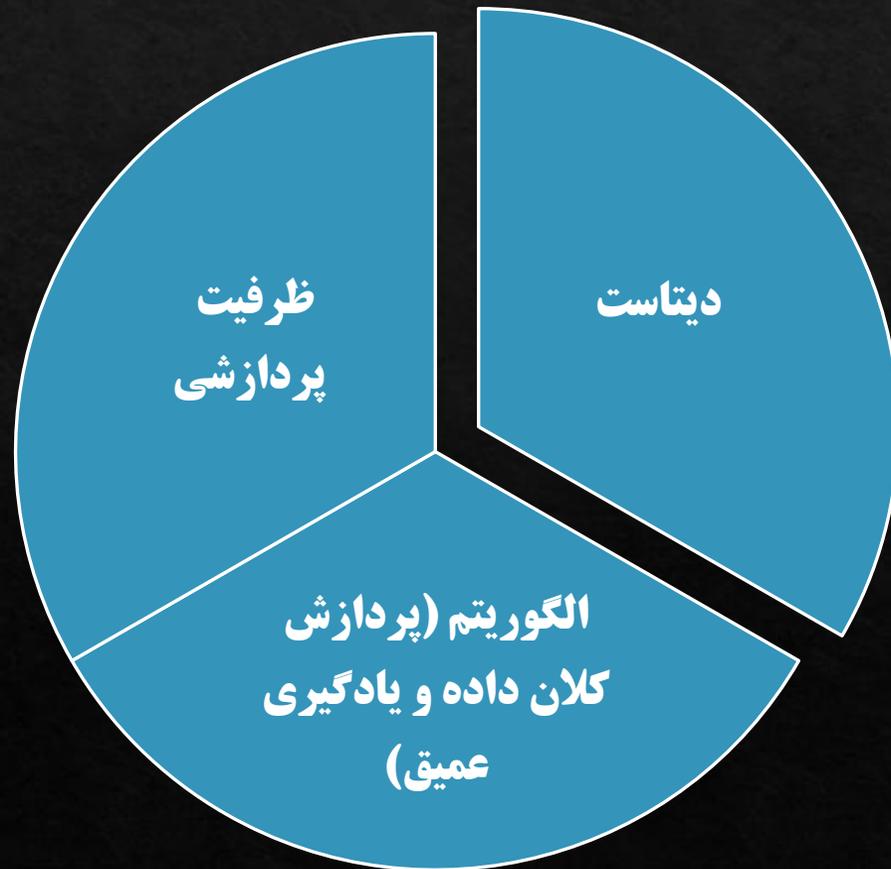
مهارت‌هایی که در درس پیدا خواهید کرد

مهارت هایی که در درس کسب خواهید کرد

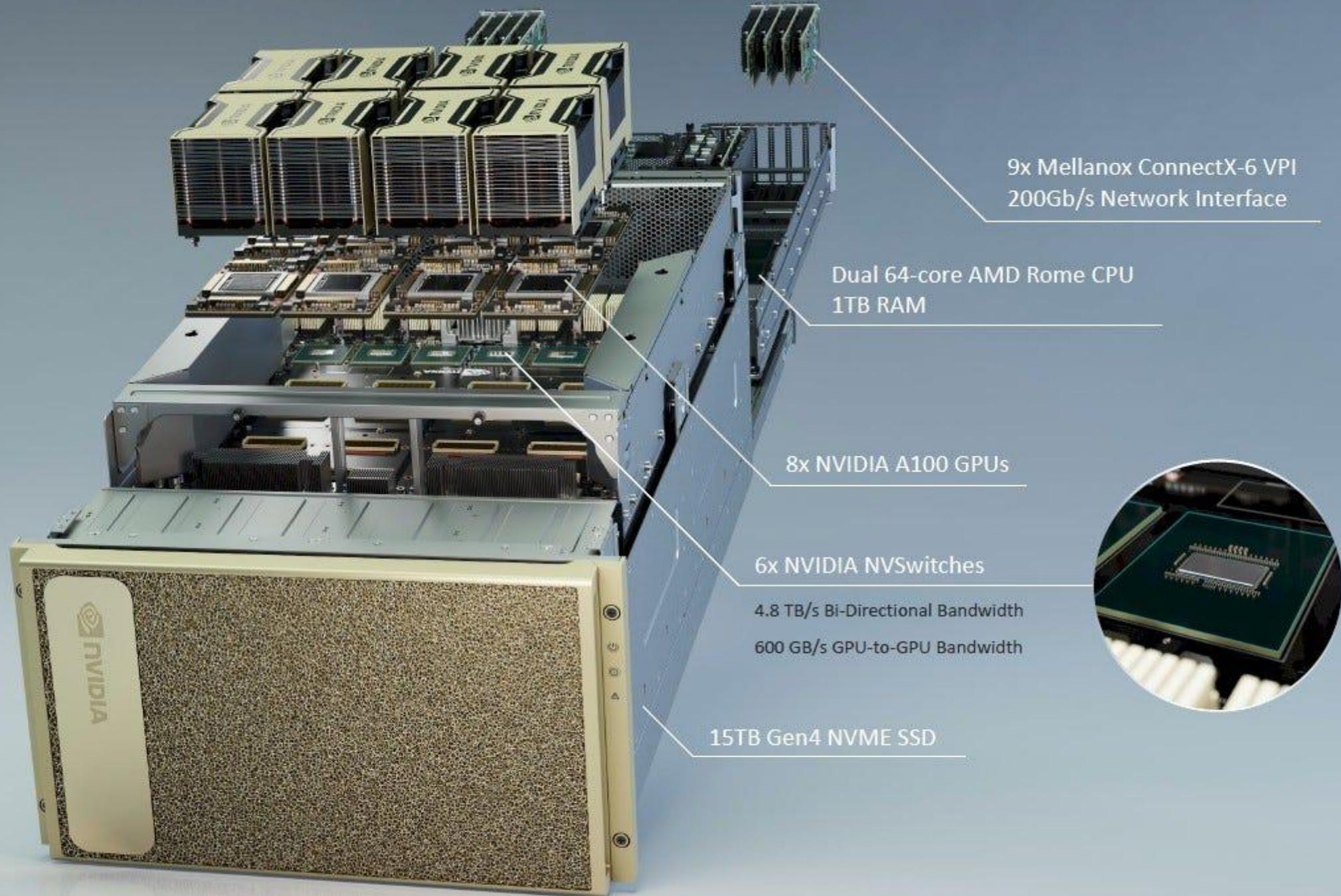
- ◆ برنامه نویسی با Python و یادگیری حین انجام پروژه
- ◆ کار با پلتفرم برنامه نویسی موازی CUDA و یادگیری حین انجام پروژه
- ◆ انجام چند پروژه واقعی در حوزه LLM
- ◆ تقویت تفکر موازی و تبدیل الگوریتم های سری به موازی و استفاده بهینه از GPU بجای CPU
- ◆ آشنایی با معماری GPU های مختلف به منظور انتخاب GPU بهینه برای حل یک مساله
- ◆ کار بر روی یک نمونه داده واقعی بزرگ
- ◆ نحوه کار با GPU های نسل آخر Nvidia از طریق کرایه سرور
- ◆ کار در بستر Linux و یادگیری حین انجام پروژه
- ◆ آشنایی با مقالات جدید حوزه LLM

چراپی تعریف درس
(یکی از جدیدترین دروس دانشگاه
و در مرز علم جهانی)

سه راس مهم هوش مصنوعی







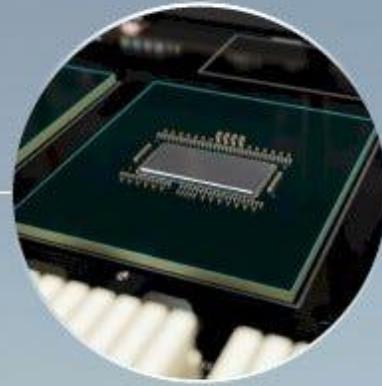
9x Mellanox ConnectX-6 VPI
200Gb/s Network Interface

Dual 64-core AMD Rome CPU
1TB RAM

8x NVIDIA A100 GPUs

6x NVIDIA NVSwitches
4.8 TB/s Bi-Directional Bandwidth
600 GB/s GPU-to-GPU Bandwidth

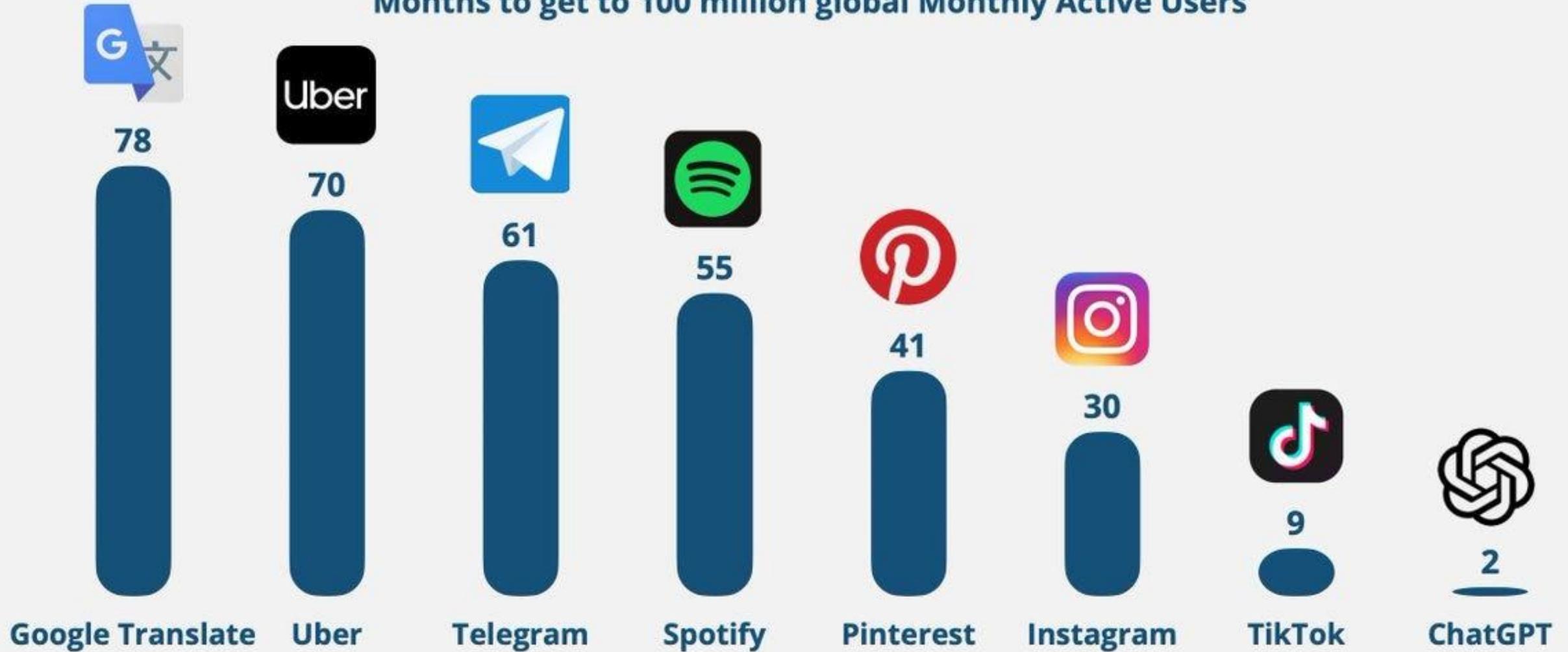
15TB Gen4 NVME SSD



اهمیت موضوع درس

Time to Reach 100M Users

Months to get to 100 million global Monthly Active Users



Source: UBS / Yahoo Finance

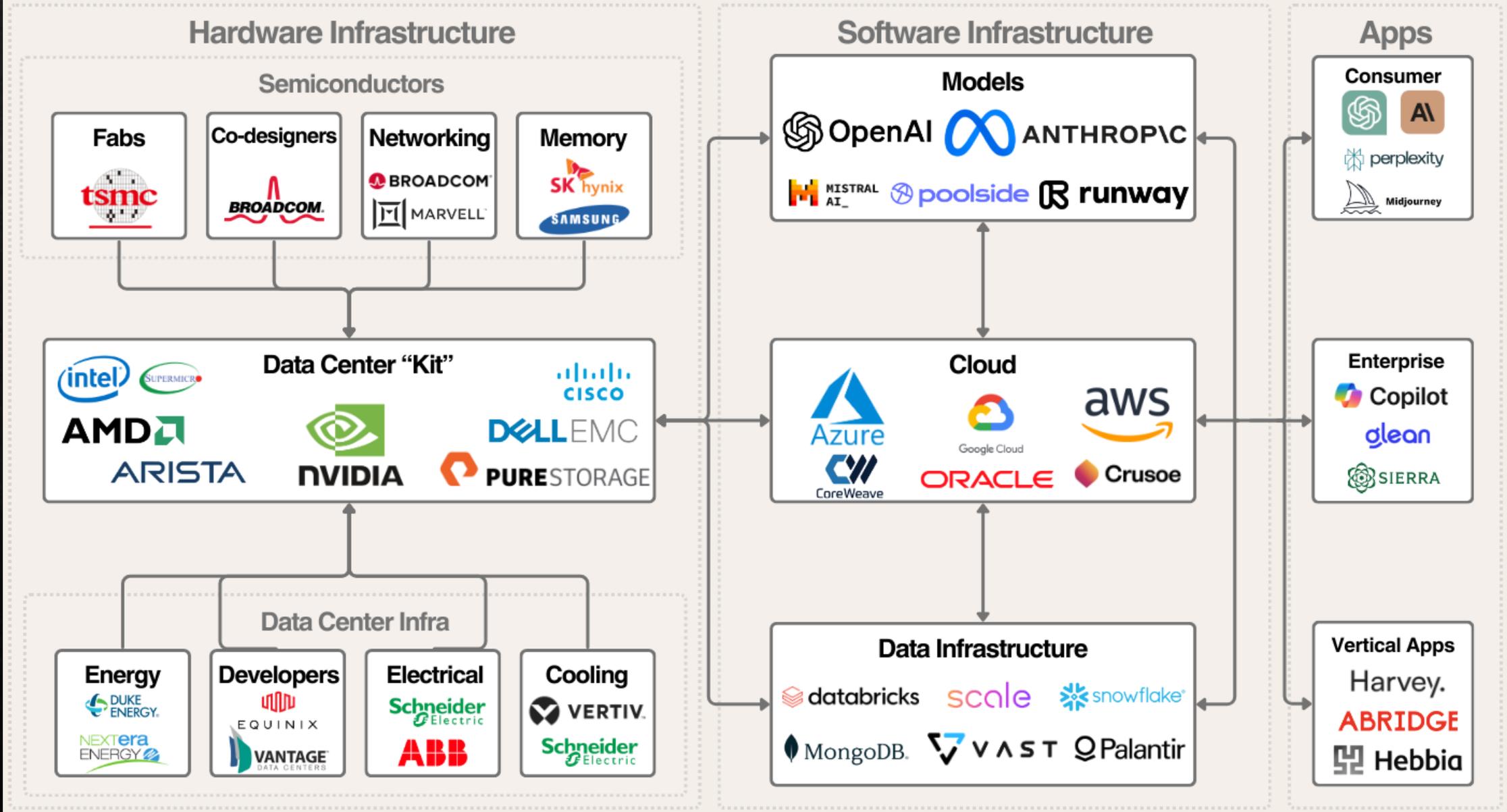
 @EconomyApp

 APP ECONOMY INSIGHTS

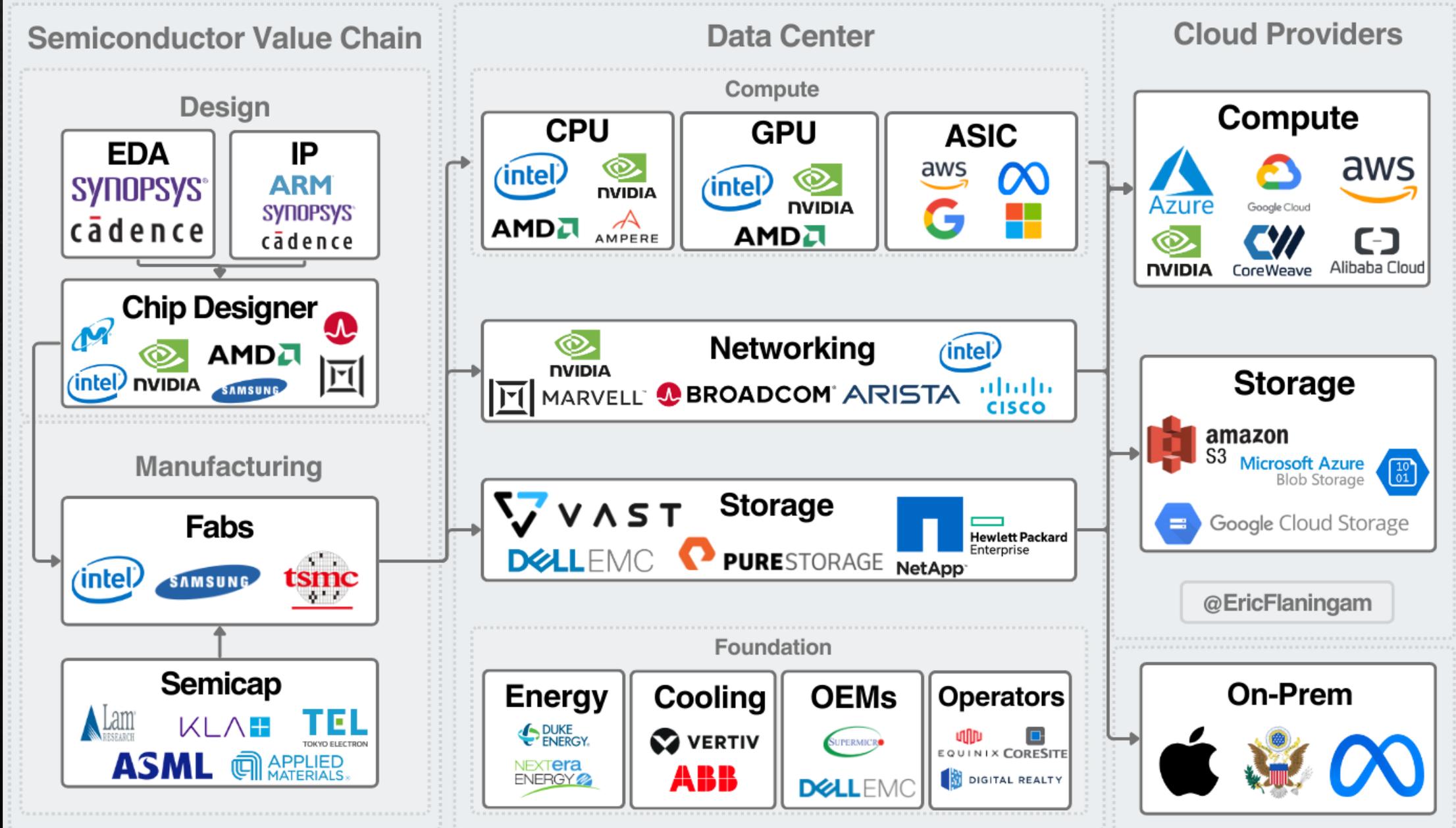
تحولات هوش مصنوعی در علوم مختلف در جهان

- ◇ حذف پزشکی سنتی و جایگزینی آن با پزشکی شخصی (Personalized medicine)
- ◇ حذف داروسازی سنتی و ایجاد داروی ویژه هر فرد با استفاده از پردازش داده ژنتیک
- ◇ از بین رفتن دموکراسی
- ◇ حذف رانندگی سنتی و ایجاد خودرو هوشمند و جاده هوشمند و نقش مهم رگولاتوری داده در آینده
- ◇ حذف وکالت در آینده
- ◇ نقش مهم مهندسان داده در آینده

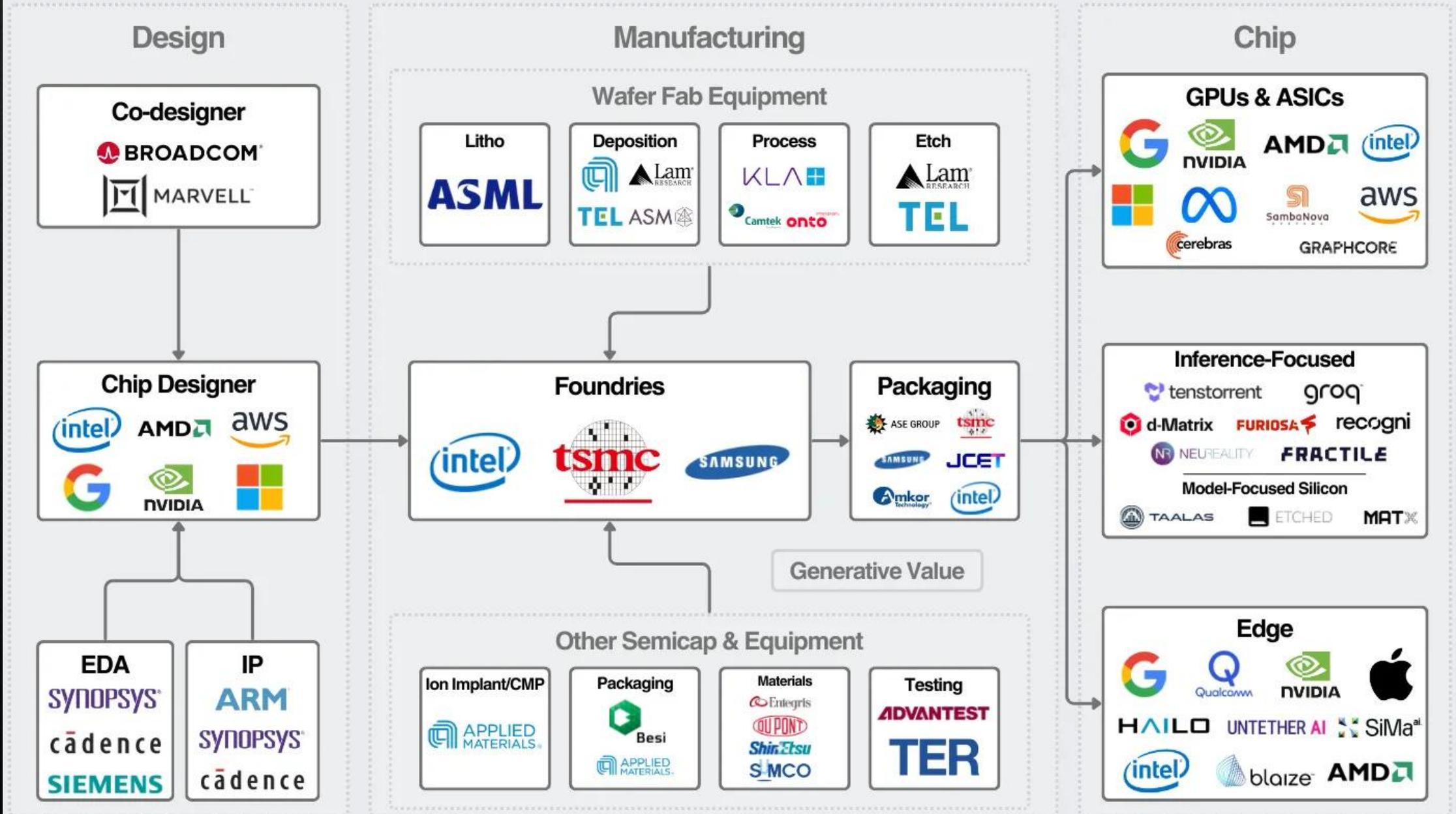
The AI Value Chain



Data Center Value Chain



AI Semiconductor Value Chain

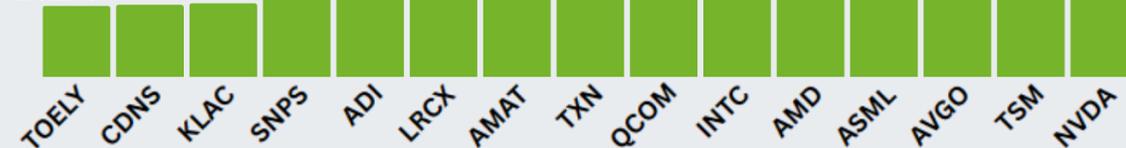


*Not an exhaustive list of companies/segments. Companies may fall into multiple segments. Several inference-focused chips can be used for training.

Semiconductor Companies Ranked by Market Cap



TOKYO ELECTRON



Nvidia's market cap of \$1.2T is 2.5x TSMC's

For reference, TEL's market cap is \$70B.

Rank	Name	Market Cap	Price	Today	Price (30 days)	Country
1	 NVIDIA NVDA	\$4.621 T	\$189.82	▲ 1.02%		 USA
2	 Apple AAPL	\$3.888 T	\$264.58	▲ 1.54%		 USA
3	 Alphabet (Google) GOOG	\$3.809 T	\$314.90	▲ 3.74%		 USA
4	 Microsoft MSFT	\$2.952 T	\$397.23	▼ 0.31%		 USA
5	 Amazon AMZN	\$2.255 T	\$210.11	▲ 2.56%		 USA
6	 TSMC TSM	\$1.921 T	\$370.54	▲ 2.82%		 Taiwan
7	 Meta Platforms (Facebook) META	\$1.658 T	\$655.66	▲ 1.69%		 USA
8	 Saudi Aramco 2222.SR	\$1.657 T	\$6.85	▲ 0.39%		 S. Arabia
9	 Broadcom AVGO	\$1.577 T	\$332.65	▼ 0.40%		 USA
10	 Tesla TSLA	\$1.545 T	\$411.82	▲ 0.03%		 USA



NVIDIA Revenue Breakdown

In \$ million

Q2 FY25

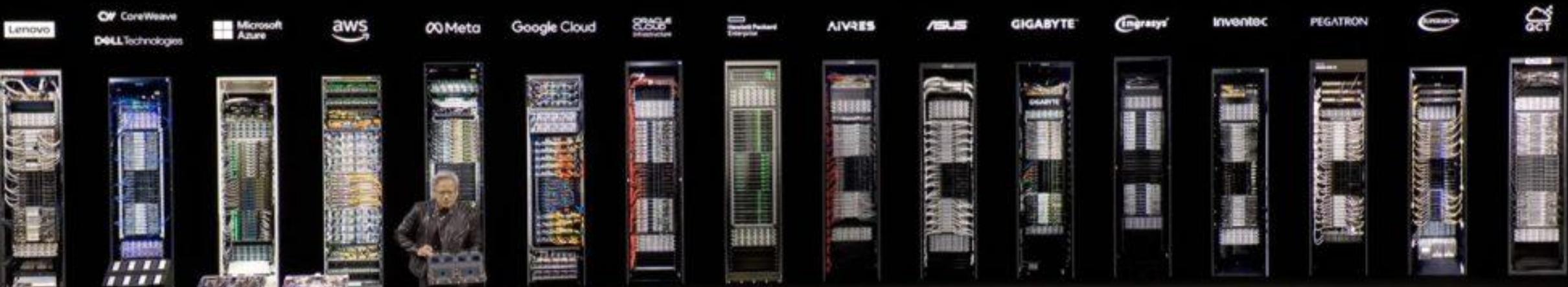
Ending July 2024



HOW THEY MAKE MONEY

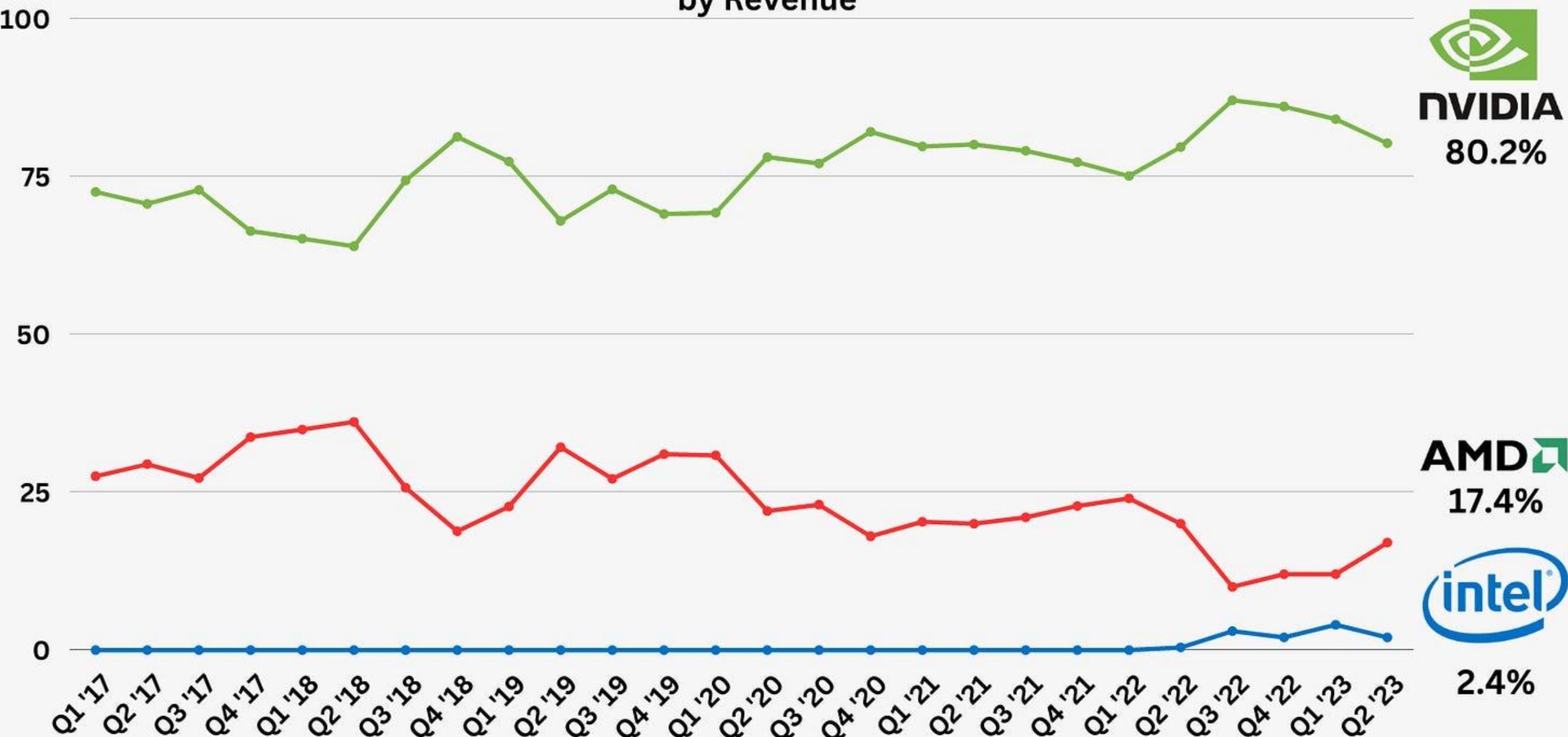


Grace Blackwell in Full Production



Desktop GPU Market Share

by Revenue




NVIDIA
80.2%


AMD
17.4%

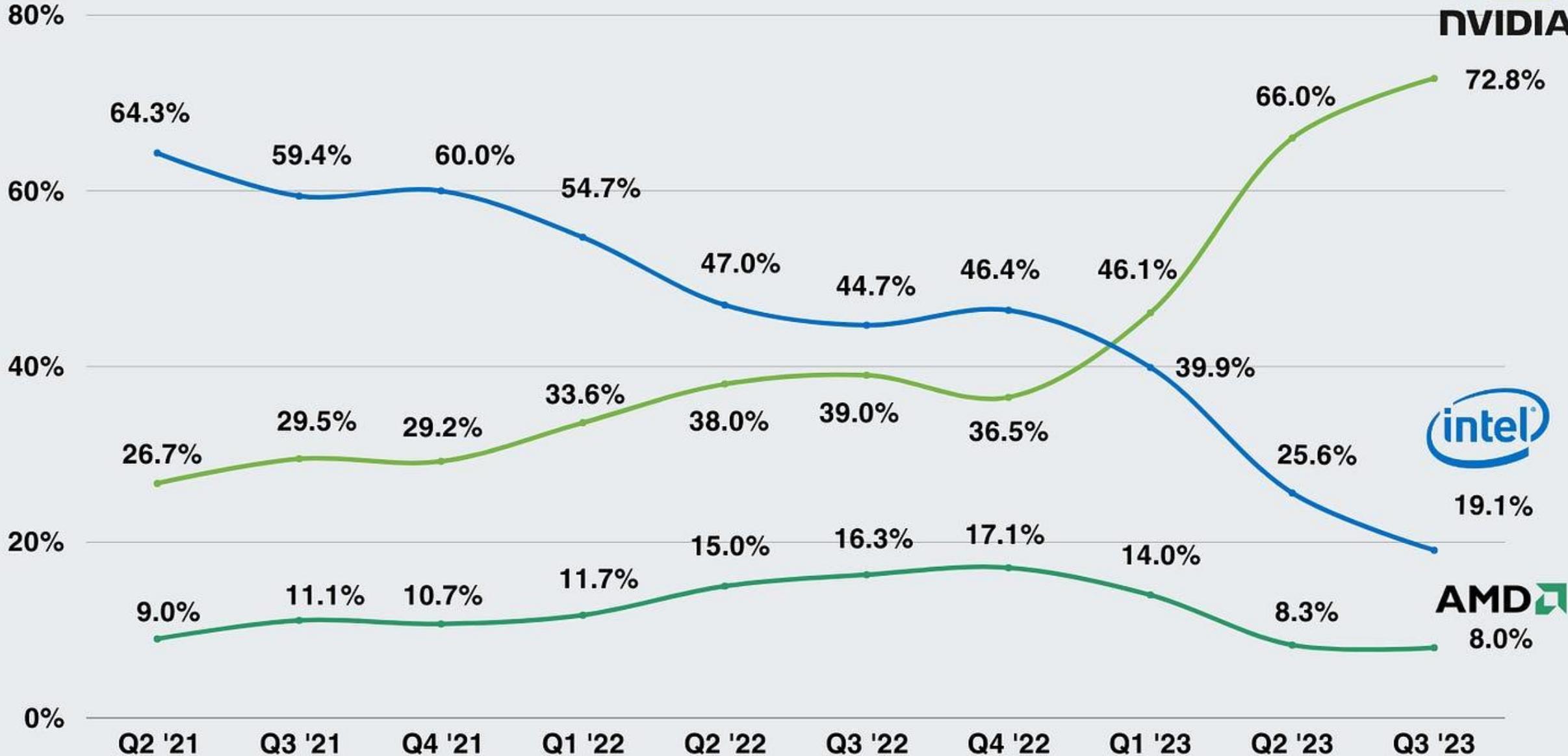

intel
2.4%

Data Center Market Share

3 largest chip providers

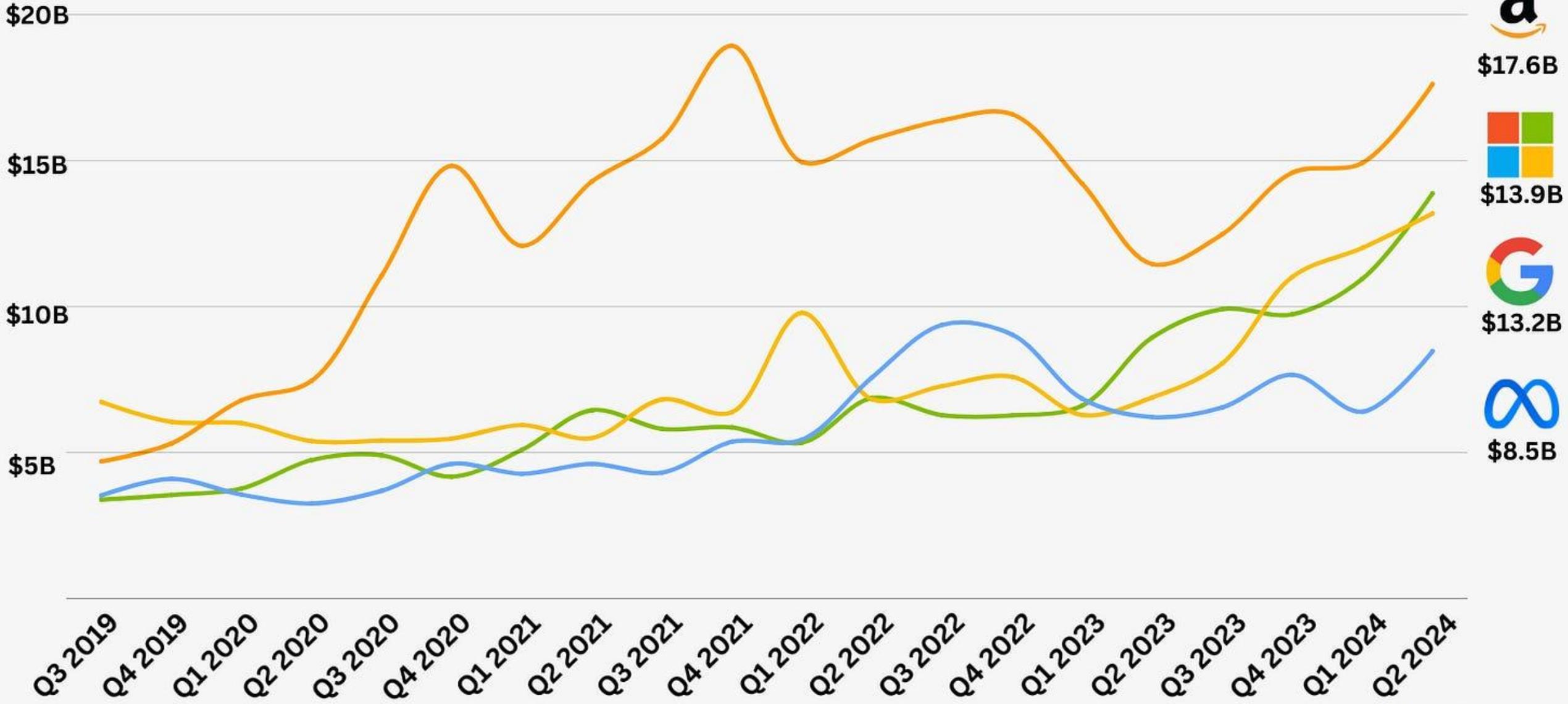


NVIDIA



Hyperscaler Quarterly Capex

MSFT GOOGL META AMZN



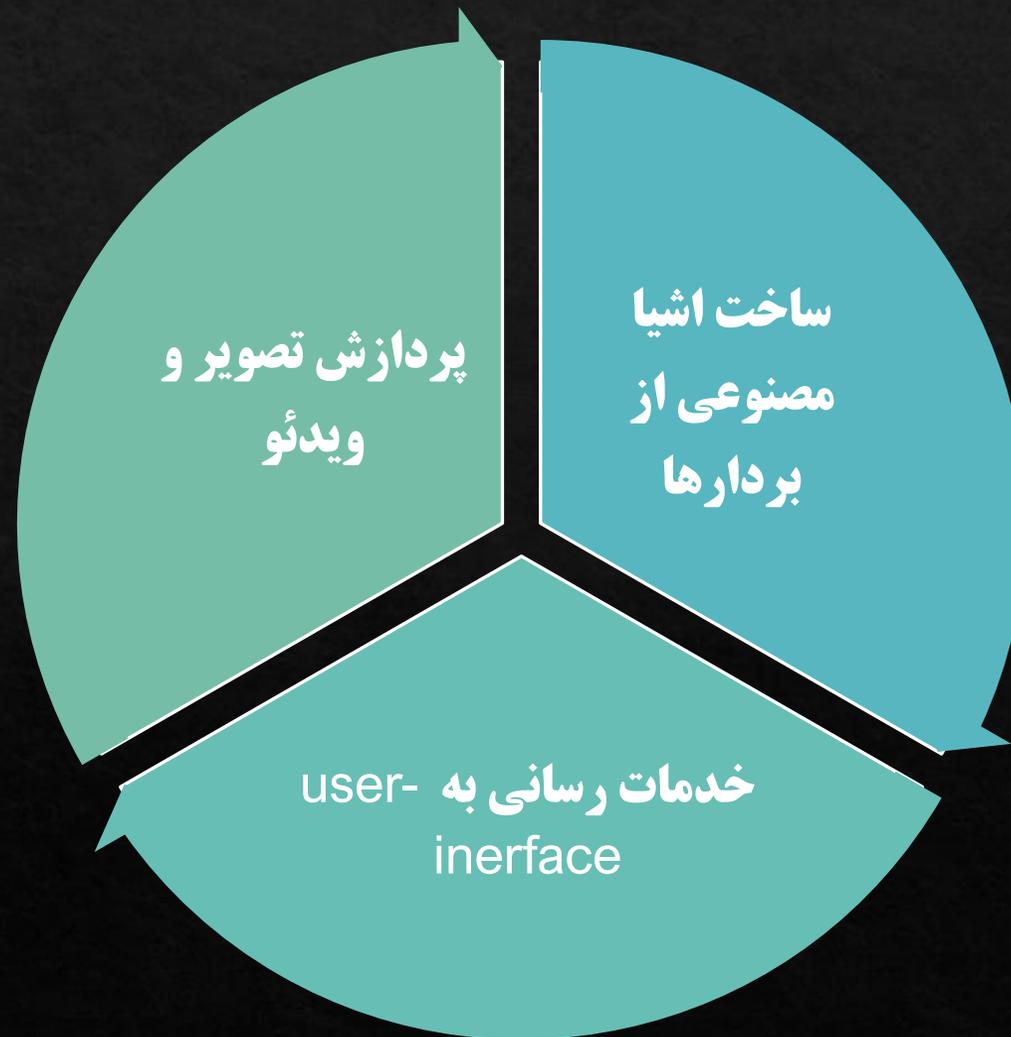
Amazon \$17.6B
Microsoft \$13.9B
Google \$13.2B
Meta \$8.5B

تاریخچه سرداننده های کراچی

و

اتفاقى مهم در این حوزه!

کارهای GPU های اولیه





The Rise of Nvidia

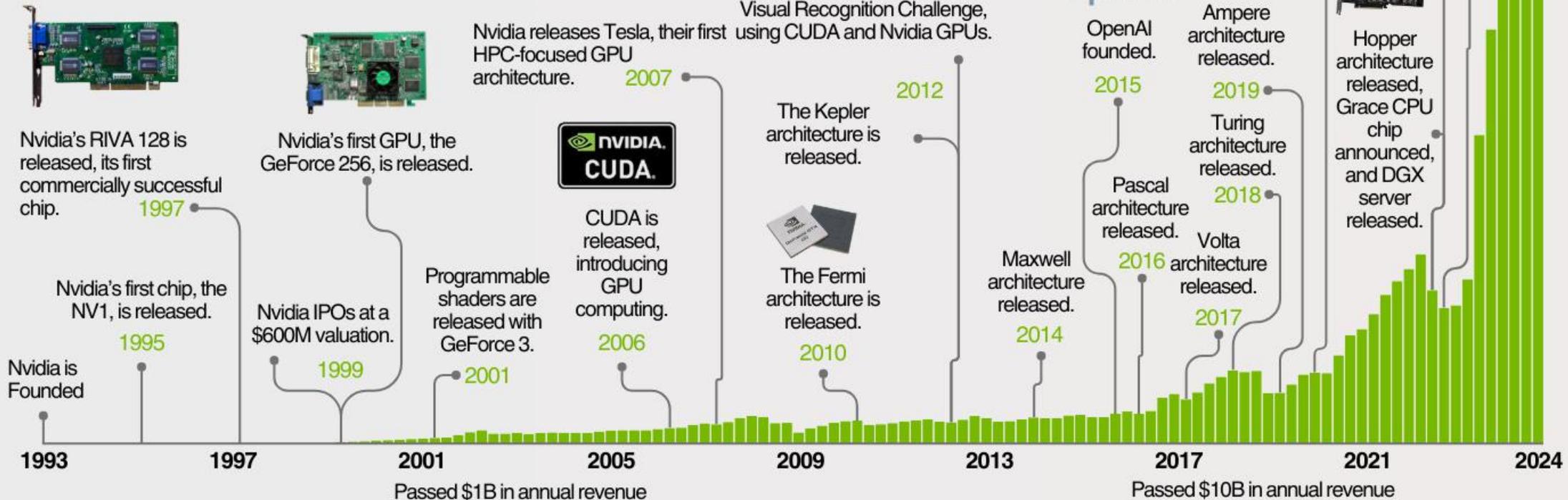


Nvidia surpasses \$35B in quarterly revenue, up 5x in just 6 quarters.

NVIDIA®

The key moments leading to its rise from a market cap of \$600M to \$3.5T.

Nvidia's Quarterly Revenue (since IPO)



Nvidia acquires Mellanox, the leader in InfiniBand networking.



ChatGPT released.

DGX Cloud released.

OpenAI

OpenAI founded.

Ampere architecture released.

Hopper architecture released, Grace CPU chip announced, and DGX server released.

Turing architecture released.

Volta architecture released.

Pascal architecture released.

Maxwell architecture released.



The Fermi architecture is released.

The Kepler architecture is released.



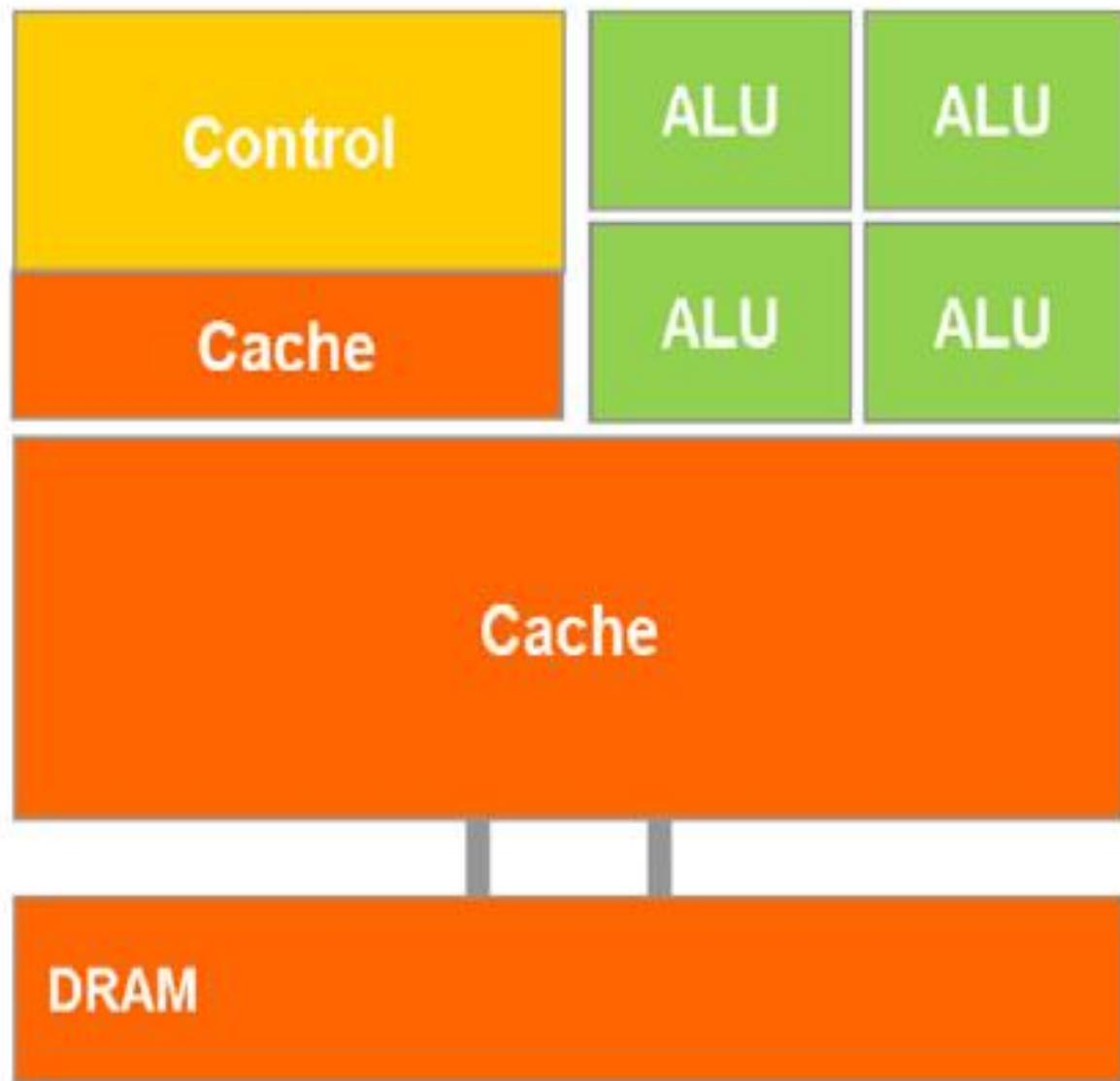
CUDA is released, introducing GPU computing.



Nvidia's first GPU, the GeForce 256, is released.

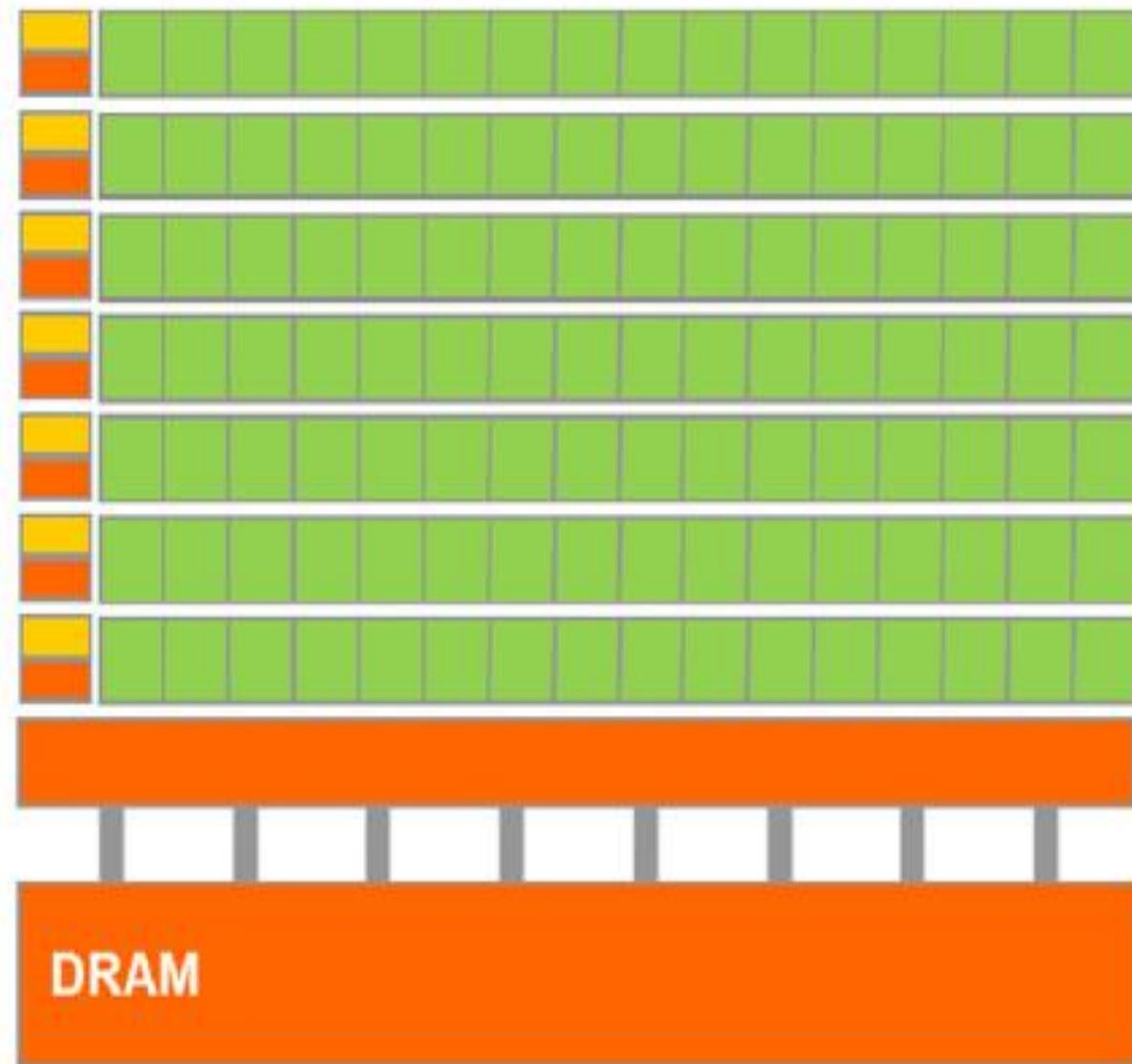


Nvidia's RIVA 128 is released, its first commercially successful chip.



CPU

(latency-oriented design)

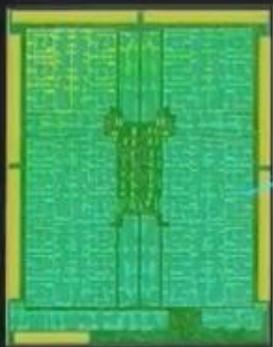


GPU

(throughput-oriented design)

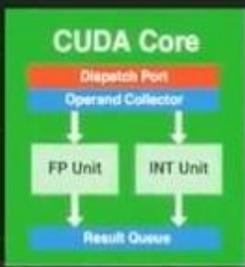
GPU

CPU

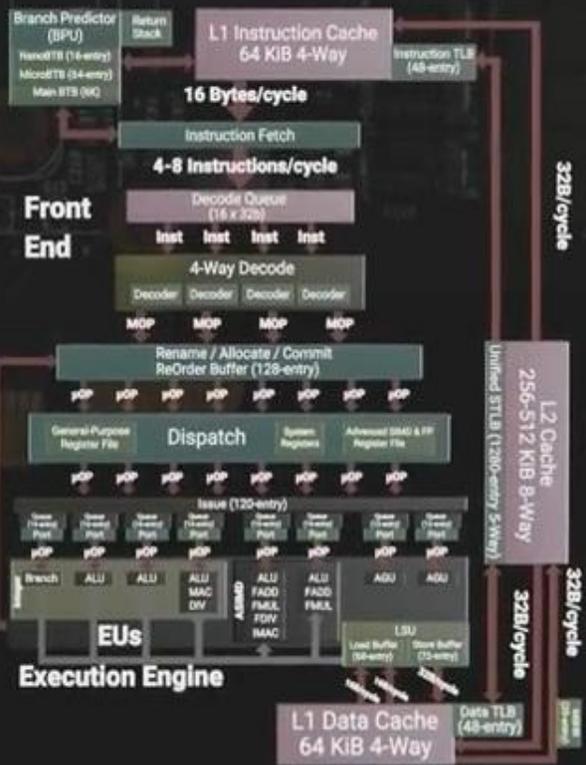


3854 Cores

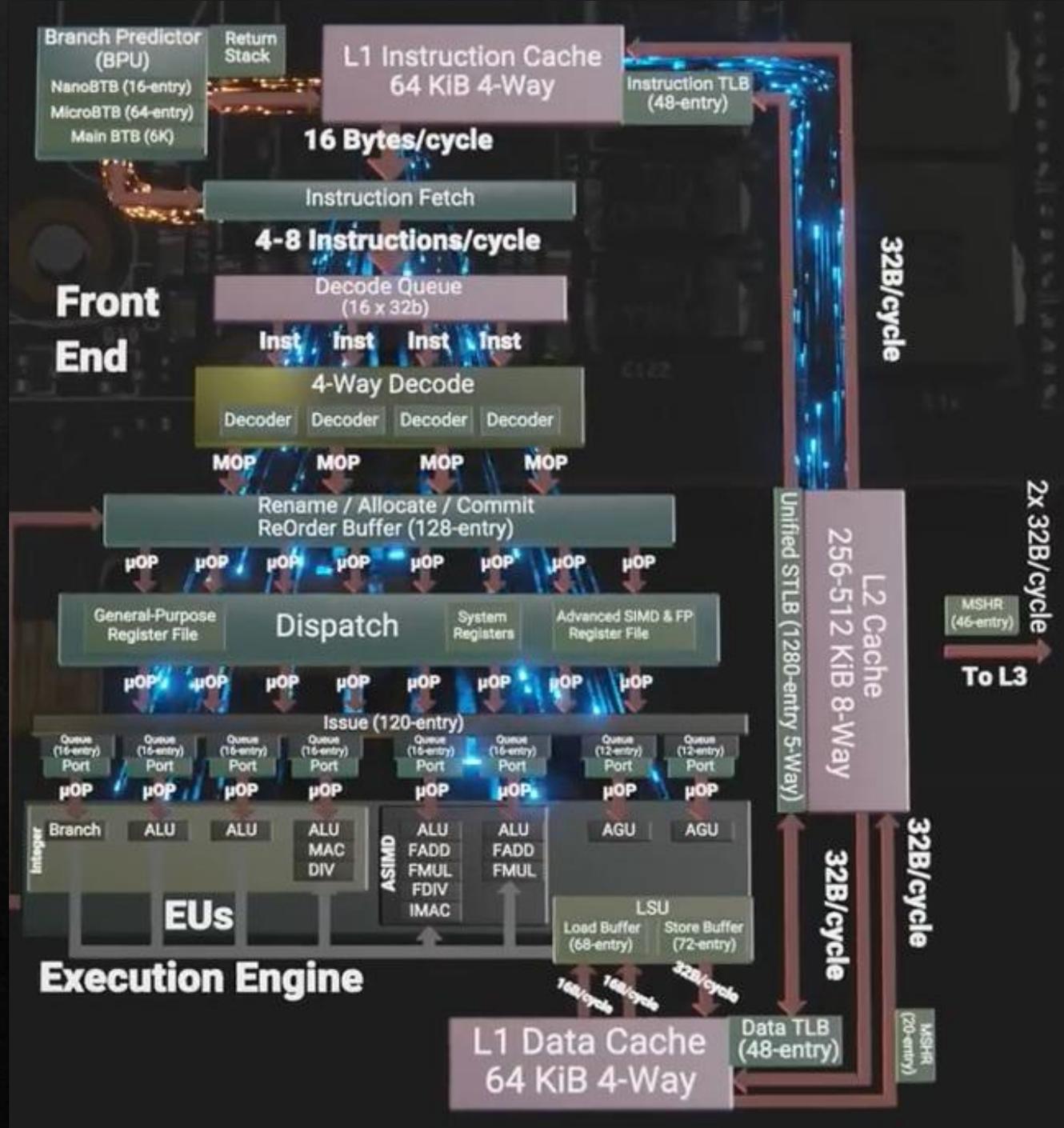
10 Cores



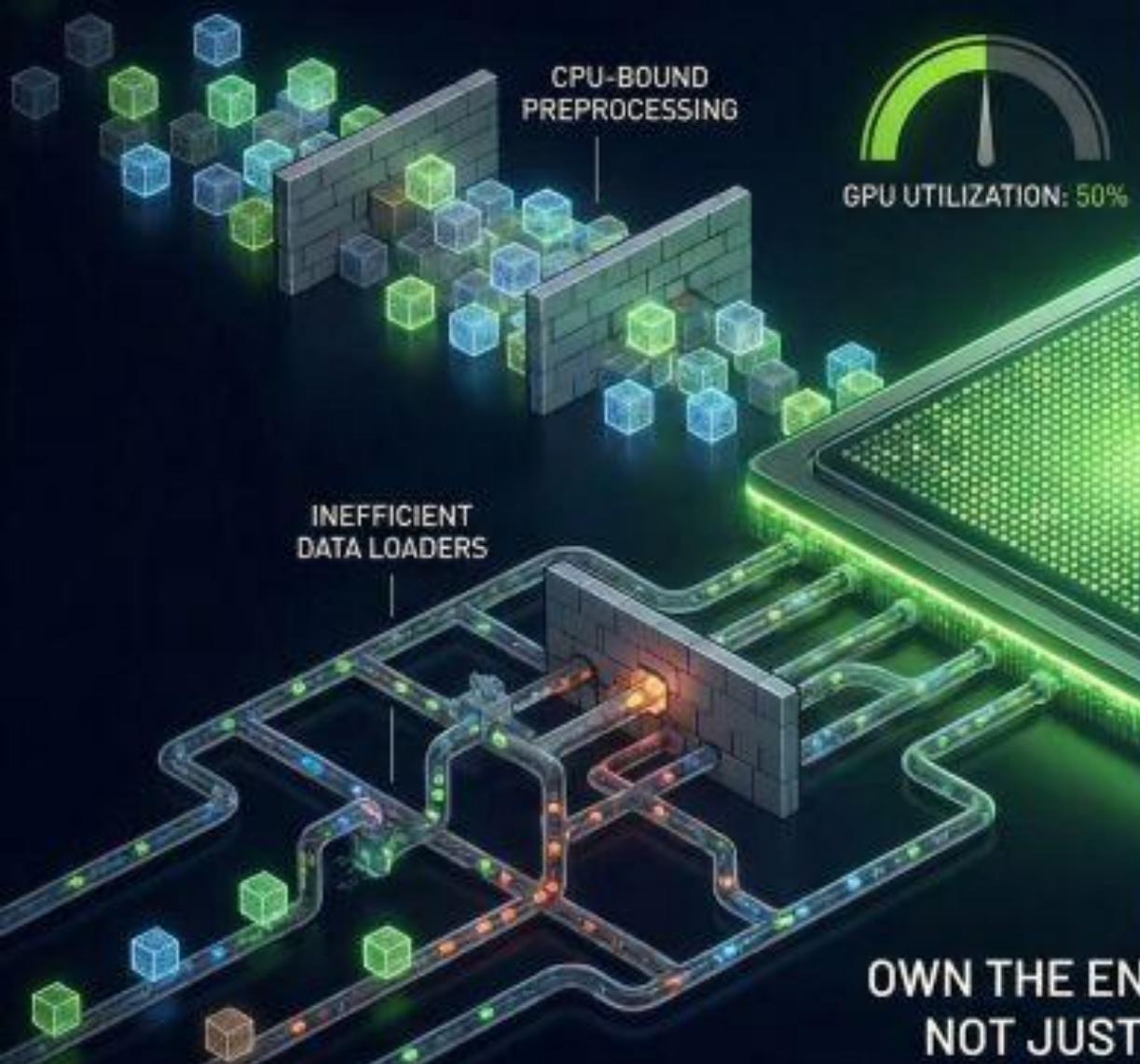
nvidia NVIDIA CORPORATION
Santa Clara, California, USA



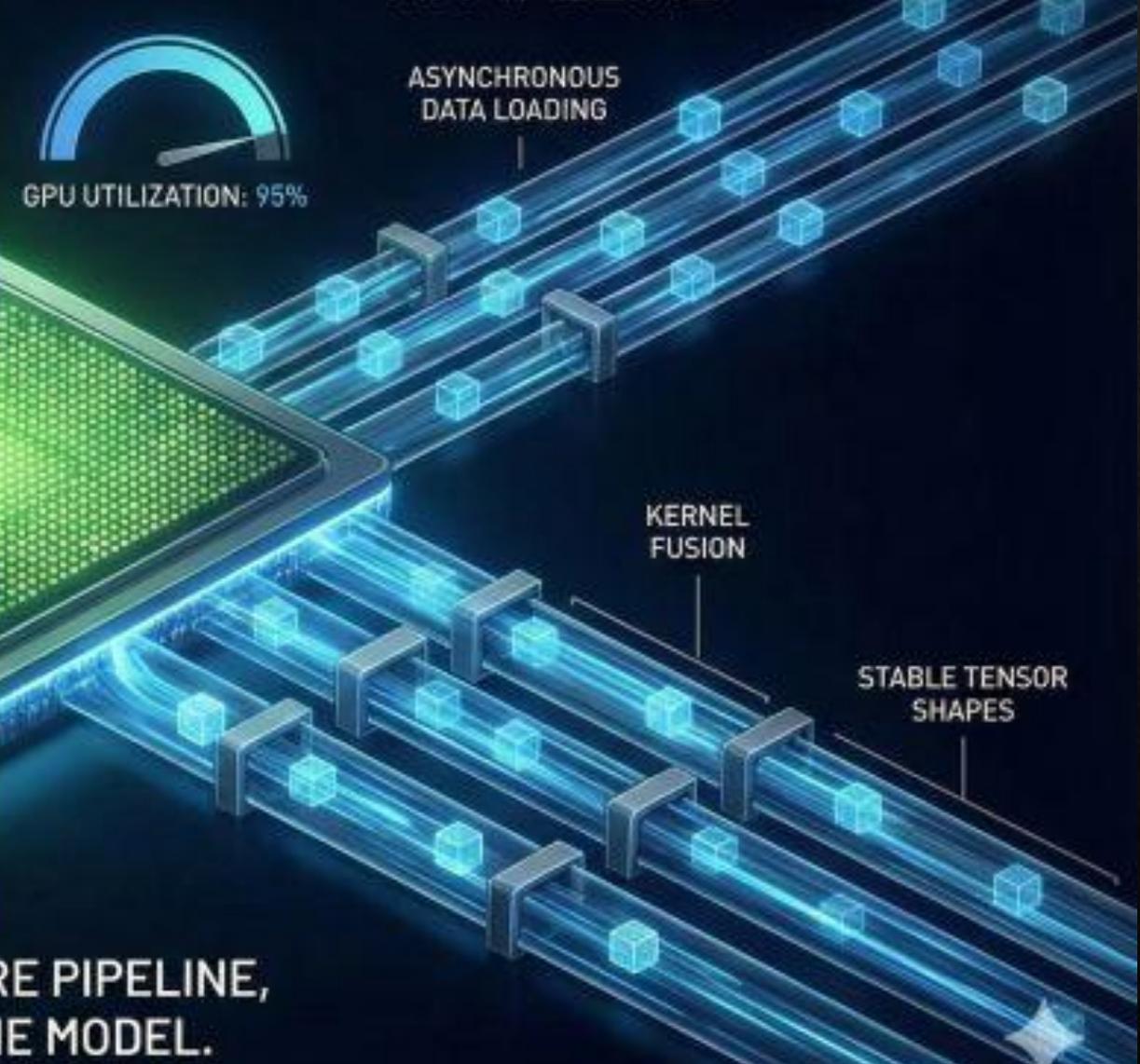
Streaming Multiprocessor (SM)



WHY GPUs ARE FAST (BUT IDLE)



HIGH-PERFORMANCE AI PIPELINE



OWN THE ENTIRE PIPELINE,
NOT JUST THE MODEL.

NVIDIA's Evolution: From Graphics to AI Leadership

Introduction of
CUDA programming
model

2006

Launch of Ada
Lovelace
architecture

2022

Introduction of
Blackwell
architecture

2023



نسل های برنده های

پردازنده های گرافیکی (GPU)



GX200

GB200

GH200

GH100

GA100

GV100

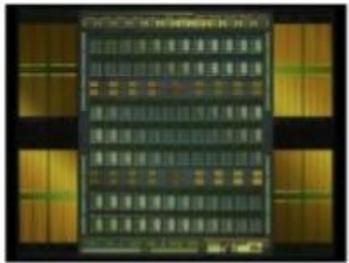
GP100

	Nvidia B200	Nvidia H100 SXM5	AMD Radeon Instinct MI300X
Power Target	1000W	700W	750W
Maximum Clock Speed	1.965 GHz	1.98 GHz	2.1 GHz
Core Count	148 SMs	132 SMs	304 Compute Units (CUs)
Last Level Cache	126 MB	50 MB	126 MB
Die Arrangement	2 Dies	Monolithic	8x Compute Dies 4x Base Dies
VRAM	288 GB HBM3E 8 TB/s	80 GB HBM3 3.3 TB/s	192 GB HBM3 5.3 TB/s
Process	TSMC 4NP	TSMC 4N	TSMC 5nm (Compute Dies) TSMC 6nm (Base Dies)

	Supported CUDA Core Precisions									Supported Tensor Core Precisions								
	FP8	FP16	FP32	FP64	INT1	INT4	INT8	TF32	BF16	FP8	FP16	FP32	FP64	INT1	INT4	INT8	TF32	BF16
NVIDIA Tesla P4	No	No	Yes	Yes	No	No	Yes	No	No	No	No	No	No	No	No	No	No	No
NVIDIA P100	No	Yes	Yes	Yes	No	No	No	No	No	No	No	No	No	No	No	No	No	No
NVIDIA Volta	No	Yes	Yes	Yes	No	No	Yes	No	No	No	Yes	No						
NVIDIA Turing	No	Yes	Yes	Yes	No	No	Yes	No	No	No	Yes	No	No	Yes	Yes	Yes	No	No
NVIDIA A100	No	Yes	Yes	Yes	No	No	Yes	No	Yes	No	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes
NVIDIA H100	No	Yes	Yes	Yes	No	No	Yes	No	Yes	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes

نسل ہامی معماری

پردازنده ہامی گرافیکی (GPU)



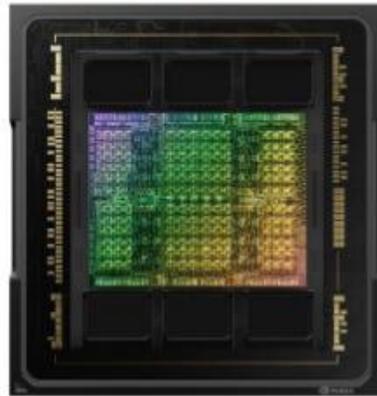
Volta

>21 billion transistors
815mm²
TSMC 12nm FFN



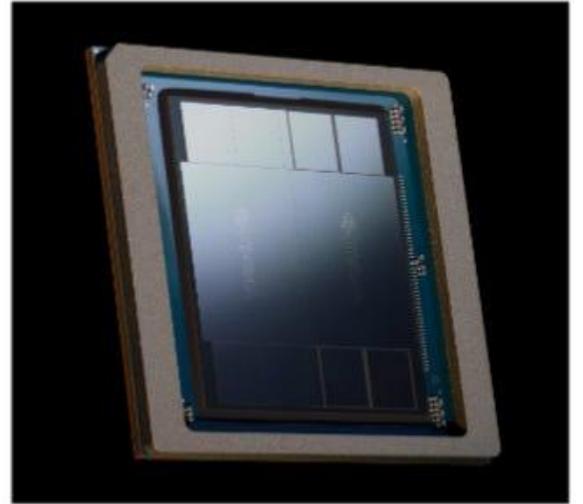
Ampere

>54 billion transistors
826 mm²
TSMC N7



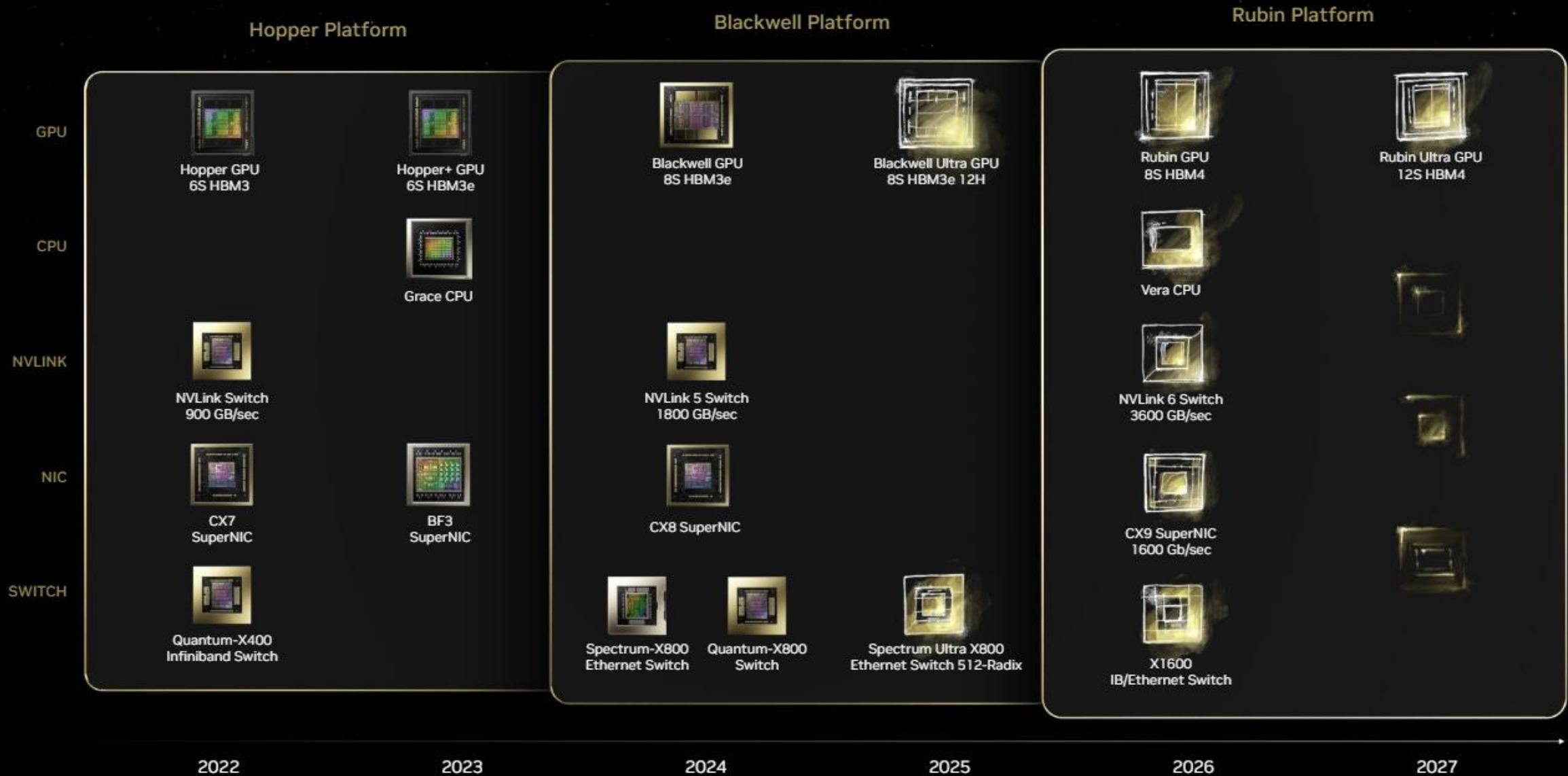
Hopper

>80 billion transistors
814 mm²
TSMC 4N



Blackwell

>208 billion transistors
>1600 mm²
TSMC 4NP



資料中心規模 · 一年節奏 · 技術限制 · 一個架構

DATACENTER SCALE · ONE-YEAR RHYTHM · TECHNOLOGY LIMITS · ONE ARCHITECTURE

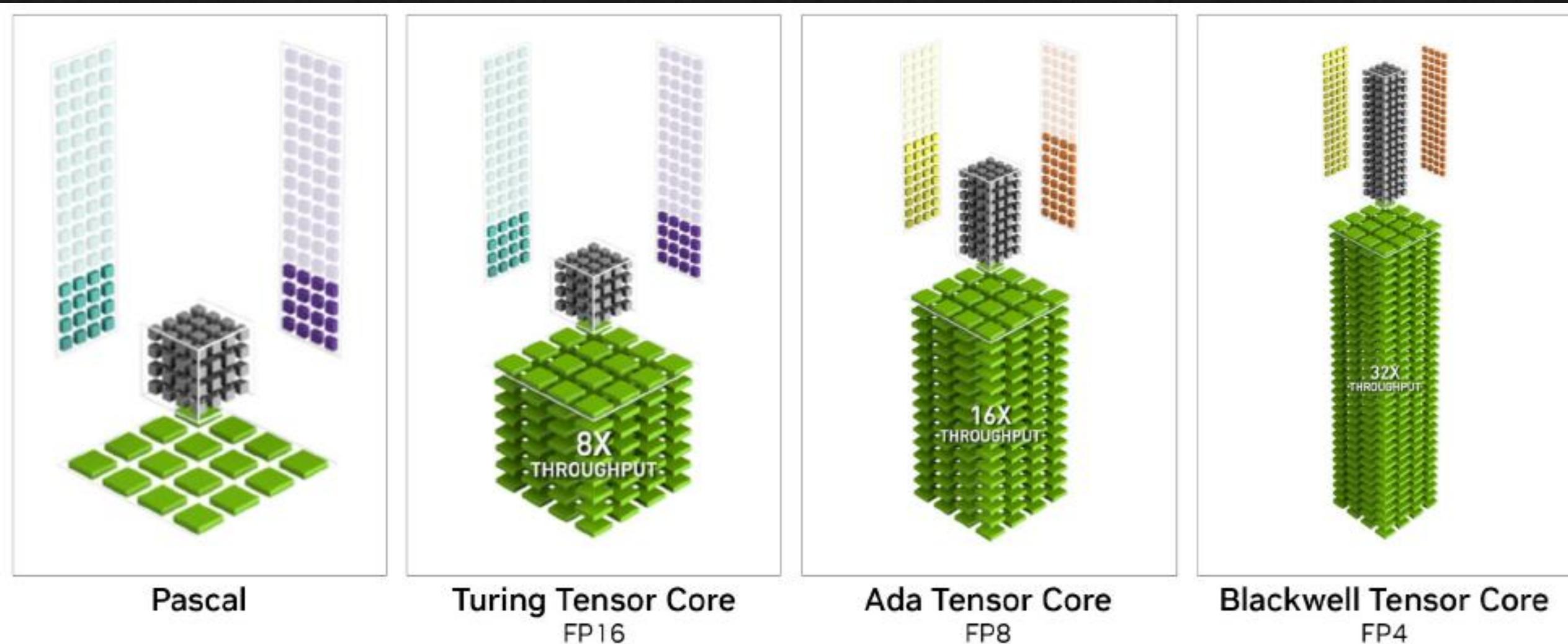


Figure 1 Blackwell 5th Generation Tensor Cores with FP4, double throughput of FP8

	Hopper	Ampere	Turing	Volta
--	--------	--------	--------	-------

Supported Tensor Core precisions

FP64, TF32, bfloat16, FP16, FP8, INT8

FP64, TF32, bfloat16, FP16, INT8, INT4, INT1

FP16, INT8, INT4, INT1

FP16

Supported CUDA* Core precisions

FP64, FP32, FP16, bfloat16, INT8

FP64, FP32, FP16, bfloat16, INT8

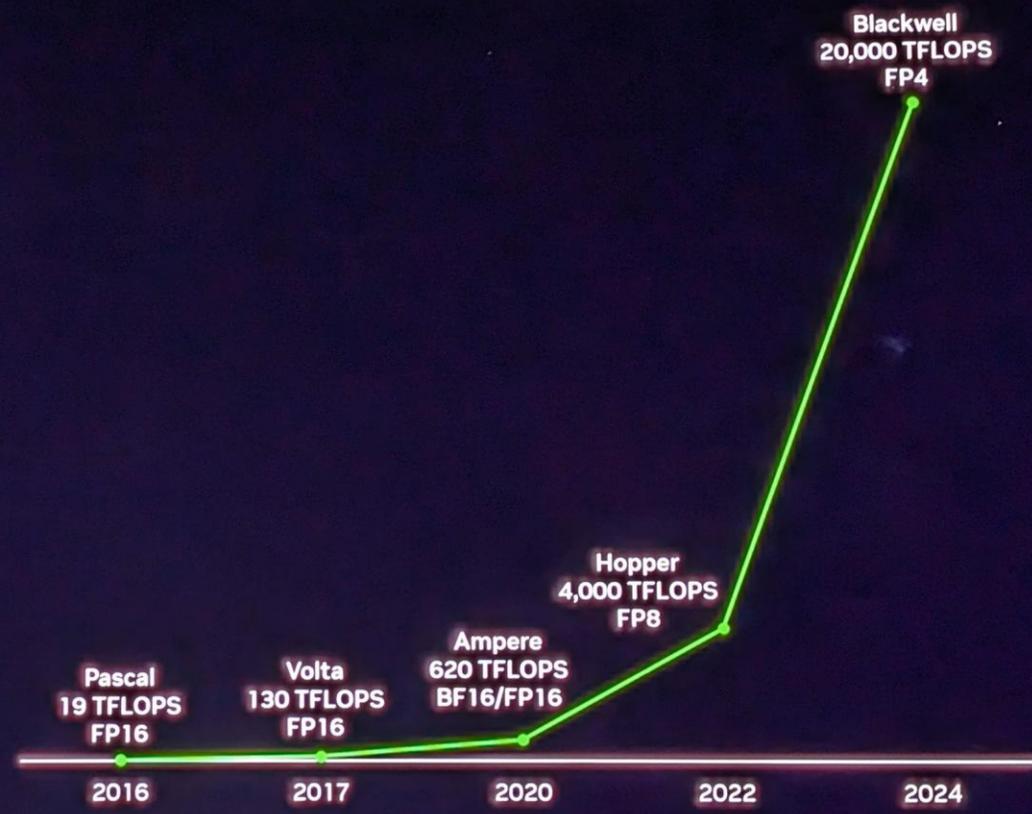
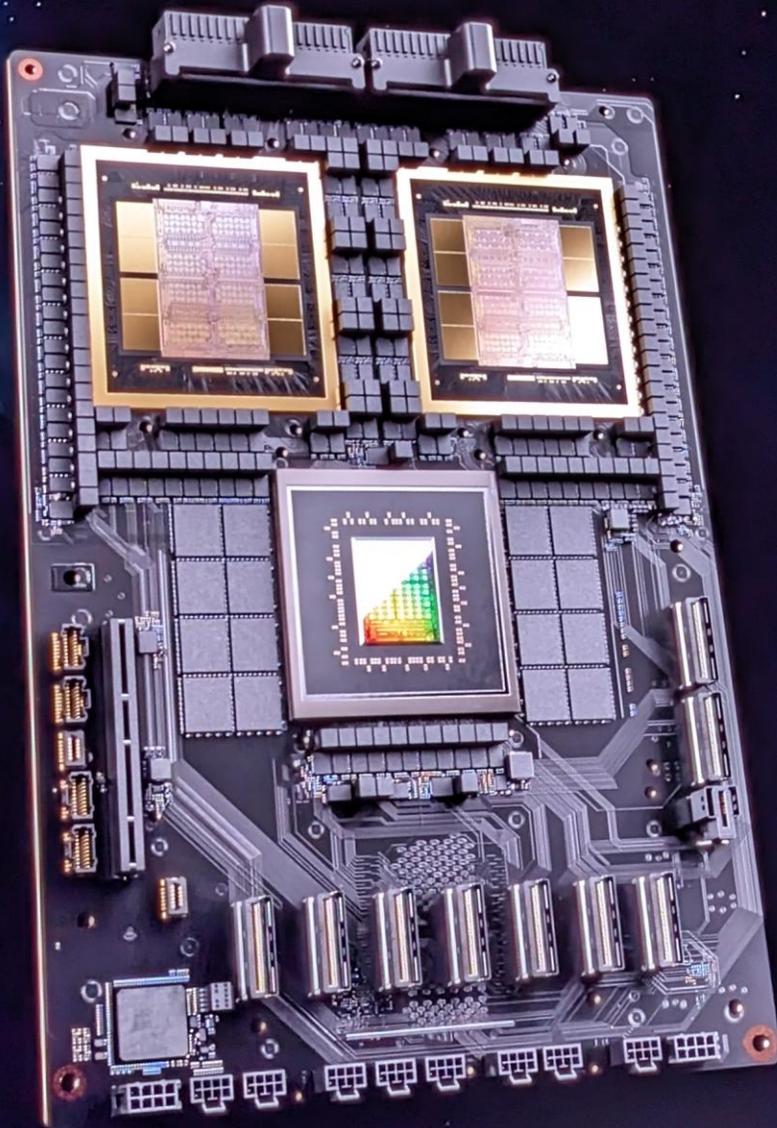
FP64, FP32, FP16, INT8

FP64, FP32, FP16, INT8

NVIDIA GPU Architecture Comparison: A100 vs H100 vs Blackwell			
Feature	NVIDIA A100 (Ampere)	NVIDIA H100 (Hopper)	NVIDIA Blackwell
Tensor Core	312 TFLOPS (FP16) 3rd Gen	989.4 TFLOPS (FP16) 4th Gen	Higher throughput 5th Gen
Key Precision	BF16, TF32	BF16, TF32, FP8	BF16, FP8, FP4
Transformer Acceleration	None	Transformer Engine (9x training)	Enhanced Transformer Engine
Max SM Count	108 SMs	132 SMs (5XMS)	Increased SMs
CUDA Cores (FP32)	6,912	16,896	Higher core count
Memory	HBM2 (40 GB)	HBM3 (80 GB, 3 TB/s)	HBM4 (Faster, improved cap)
Asynchronous Execution	Basic support	Tensor Memory Accelerator (TMA)	Enhanced async data transfer
MIG	1st Gen	2nd Gen, Secure MIG	Improved isolation & perf

		RTX 40 Series	RTX 30 Series	RTX 20 Series
NVIDIA Architecture	Architecture Name	Ada Lovelace	Ampere	Turing
	Streaming Multiprocessors	2x FP32	2x FP32	1x FP32
	Ray Tracing Cores	Gen 3	Gen 2	Gen 1
	Tensor Cores (AI)	Gen 4	Gen 3	Gen 2
	NVIDIA DLSS	3	2	2

1000X AI Compute in 8 Years



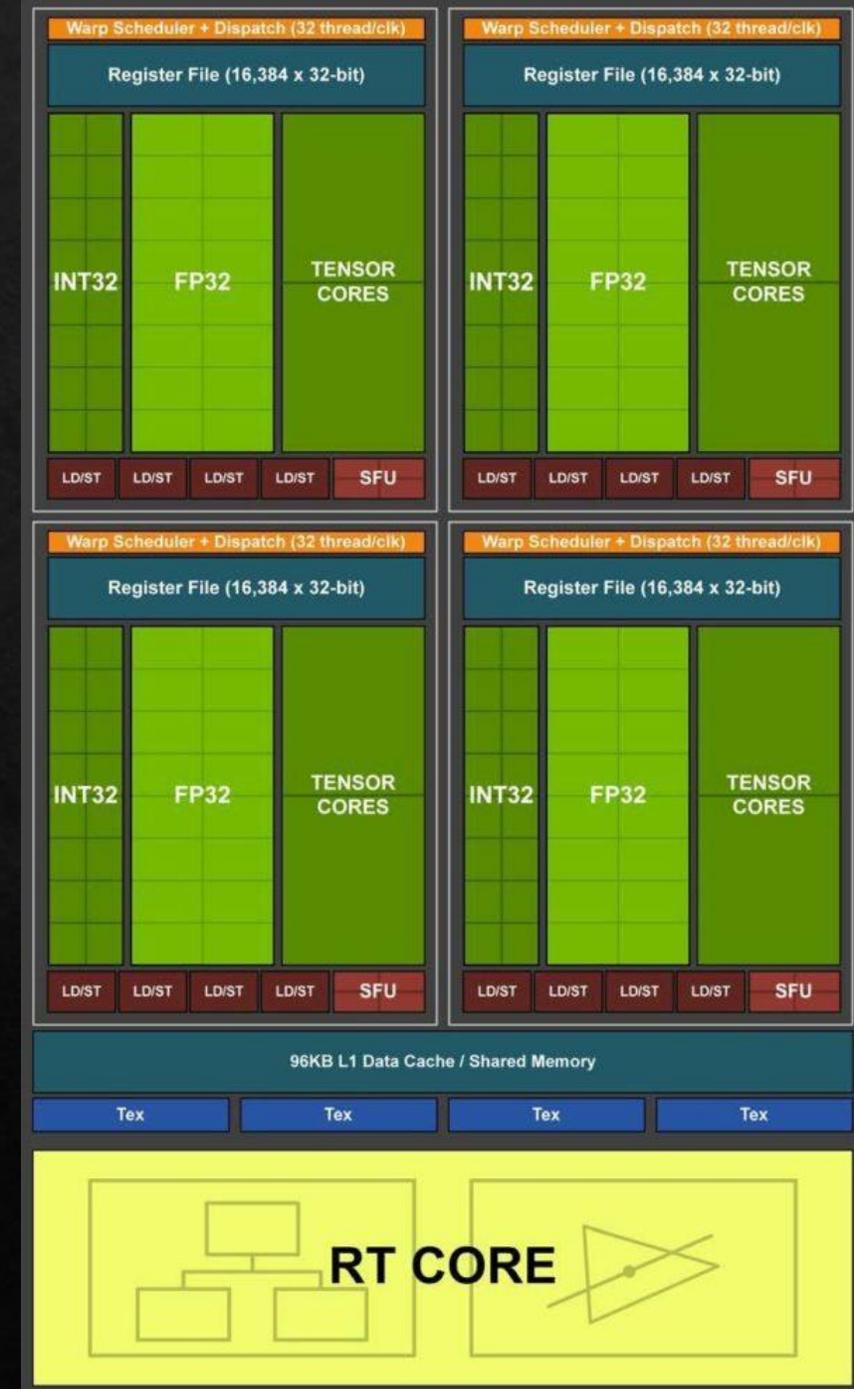
نسل های هسته های

پردازنده های گرافیکی (GPU)

Streaming Multiprocessor (SM)



SM



NVIDIA RTX FULL-STACK INVENTIONS

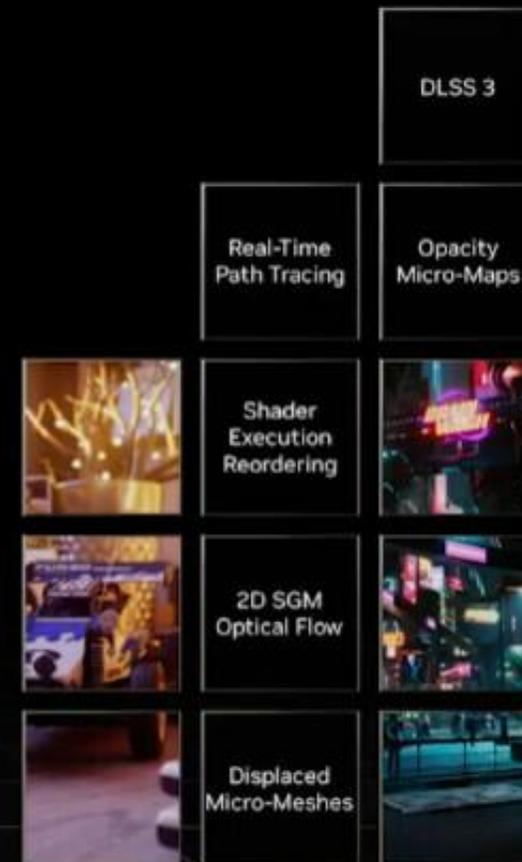
1ST GEN RTX



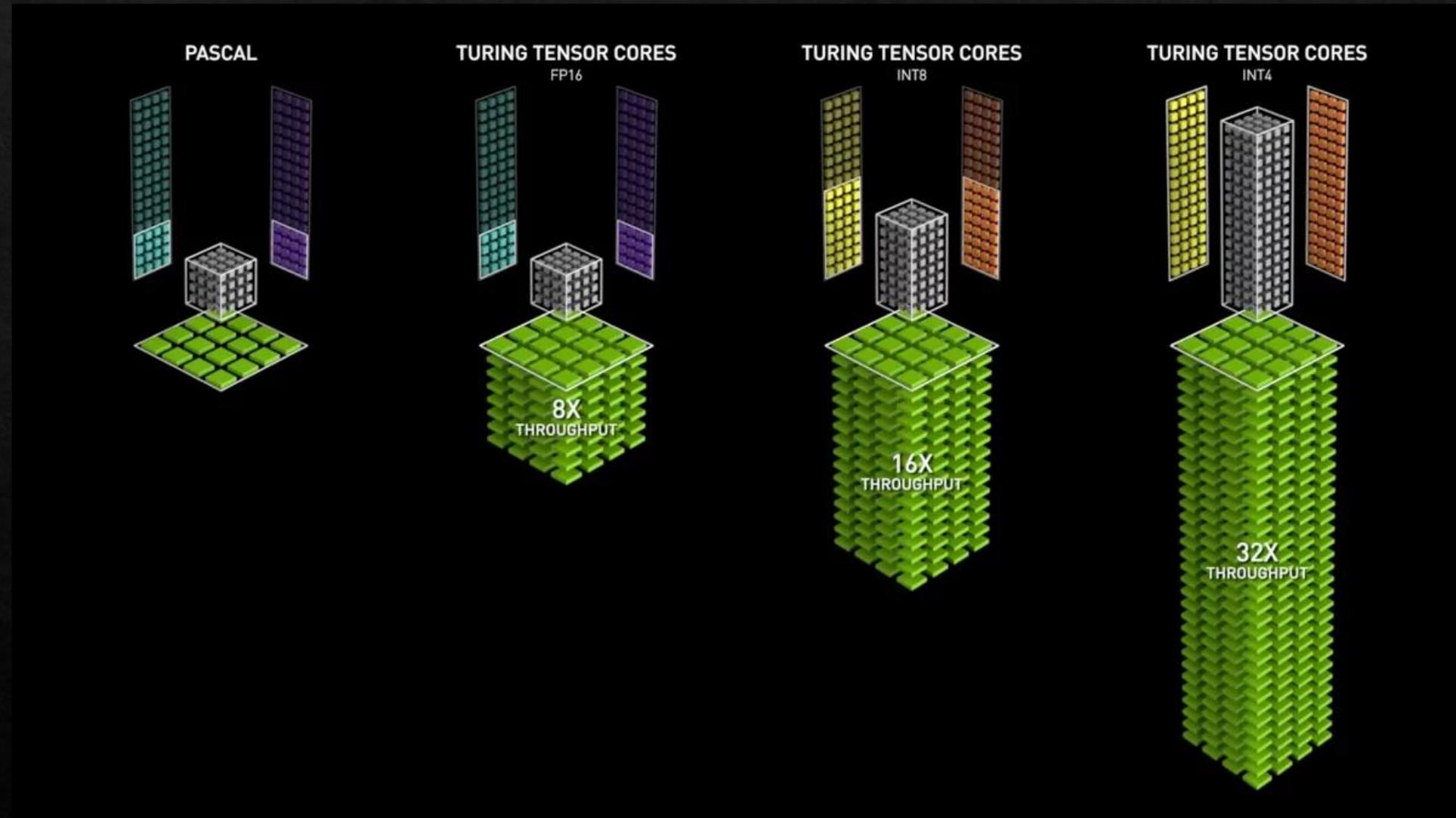
2ND GEN RTX



3RD GEN RTX

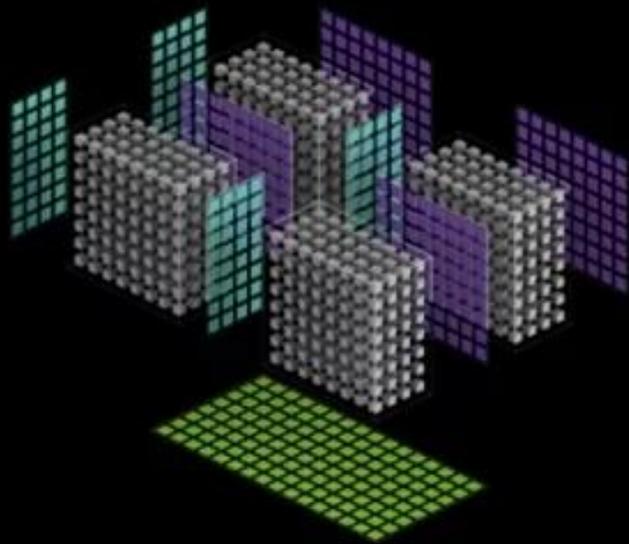


The Evolution of Tensor Cores in NVIDIA GPUs

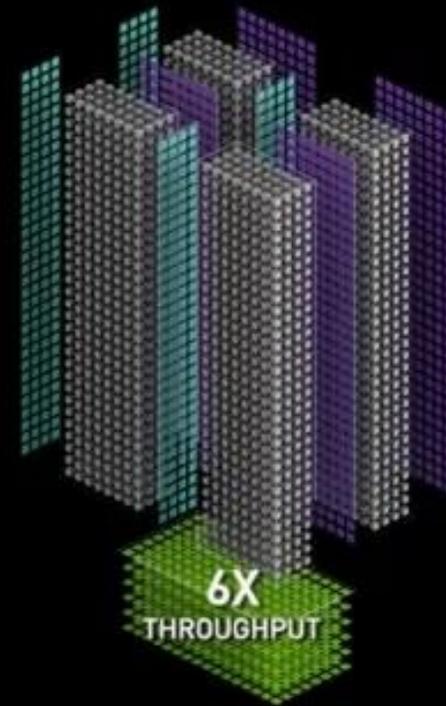


Fourth Generation: Hopper Architecture

A100 FP16



H100 FP8



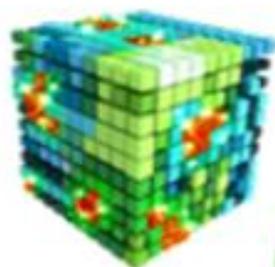
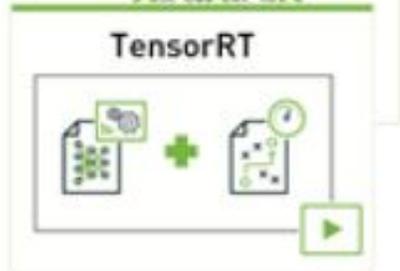
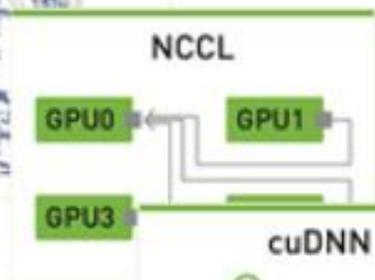
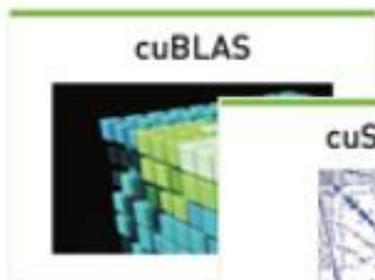
CUDA



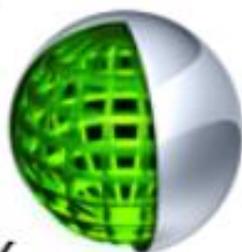
Fortran



OpenACC
Directives for Accelerators



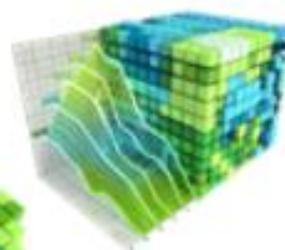
CUDA
MEMCHECK



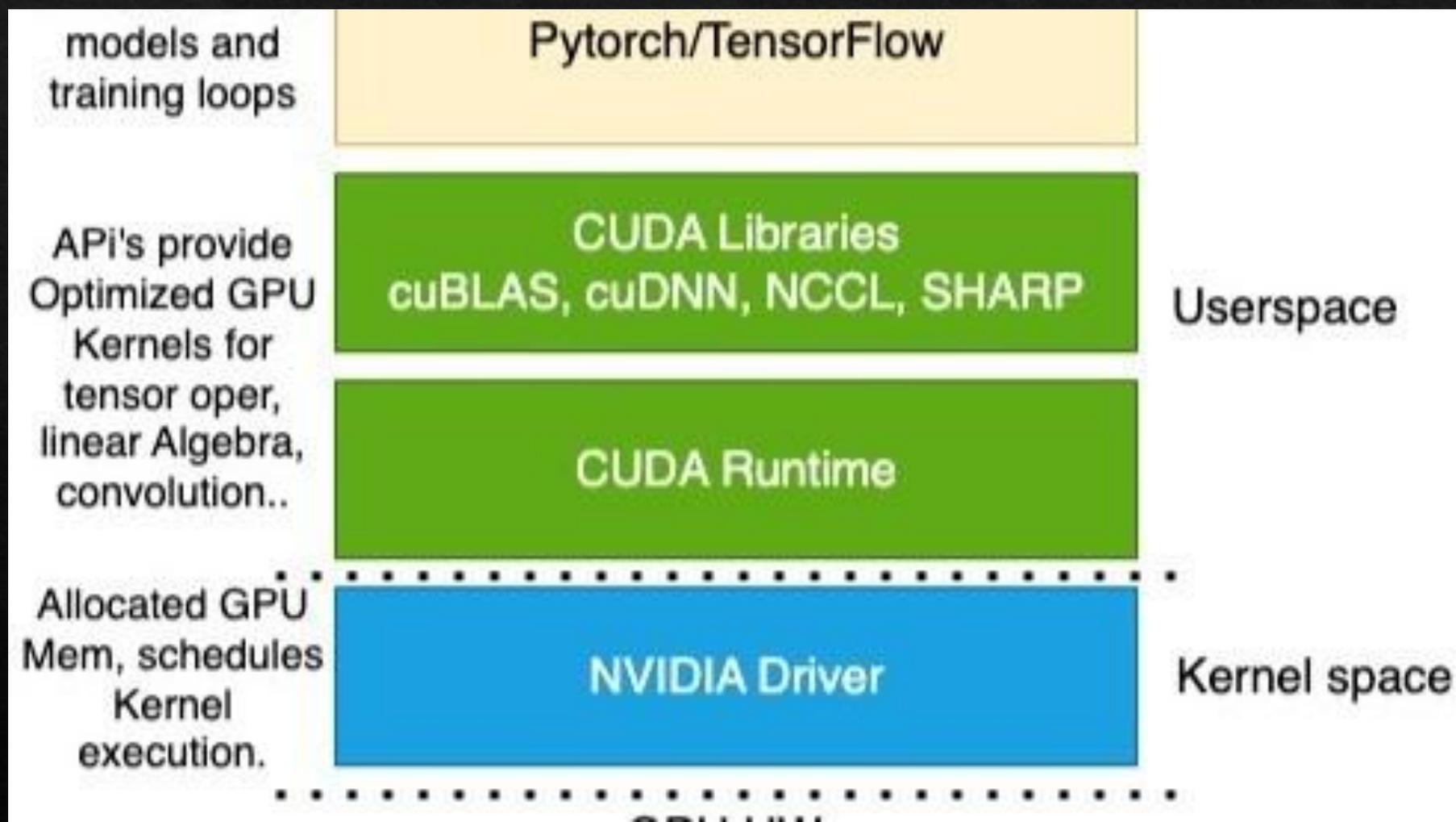
Nsight IDE

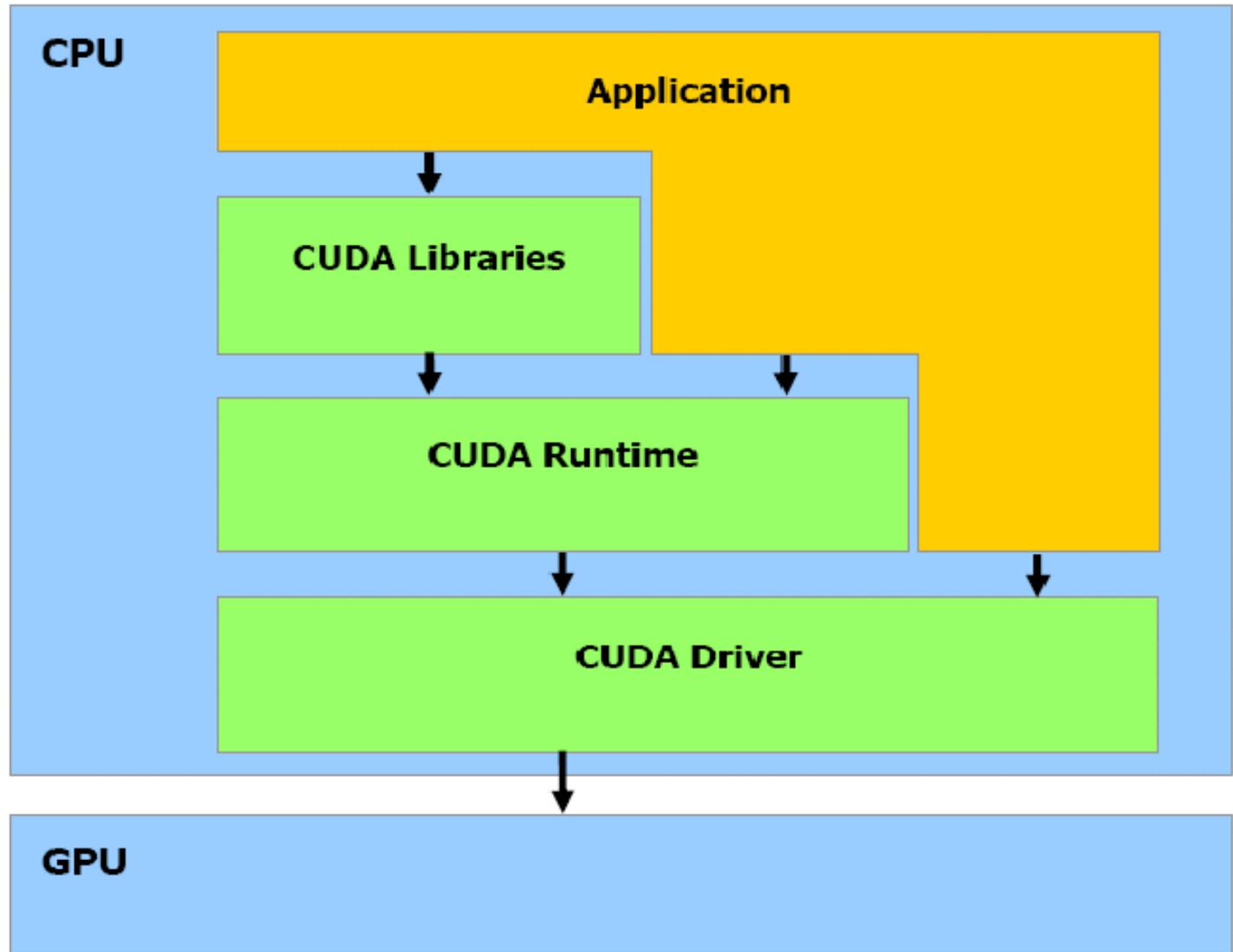


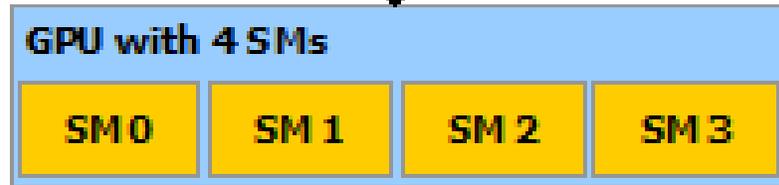
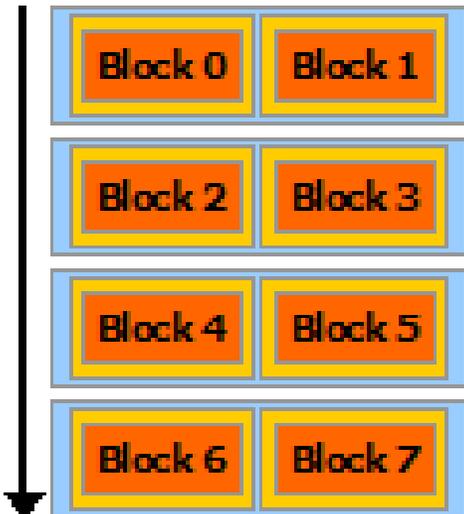
CUDA-GDB
Debugger



NVIDIA
Visual Profiler



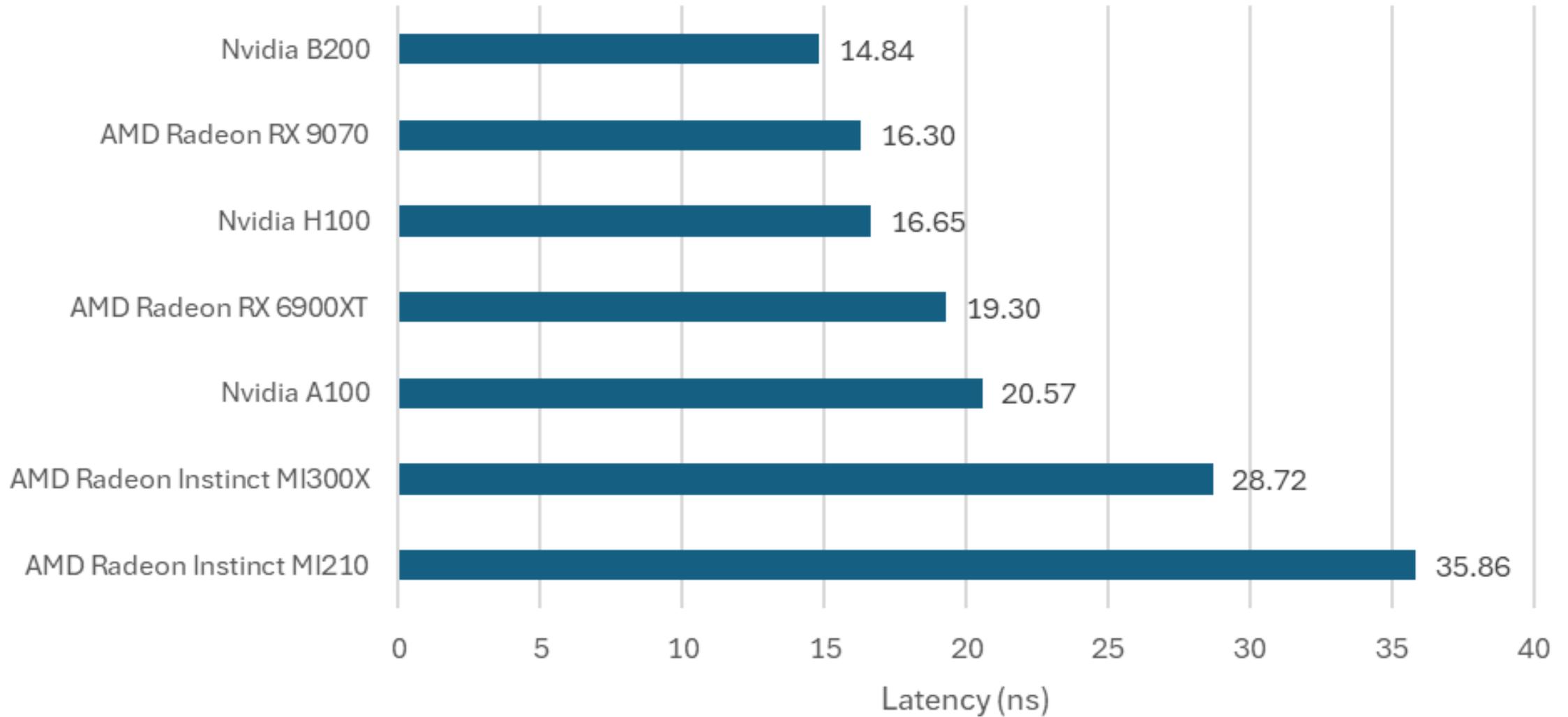




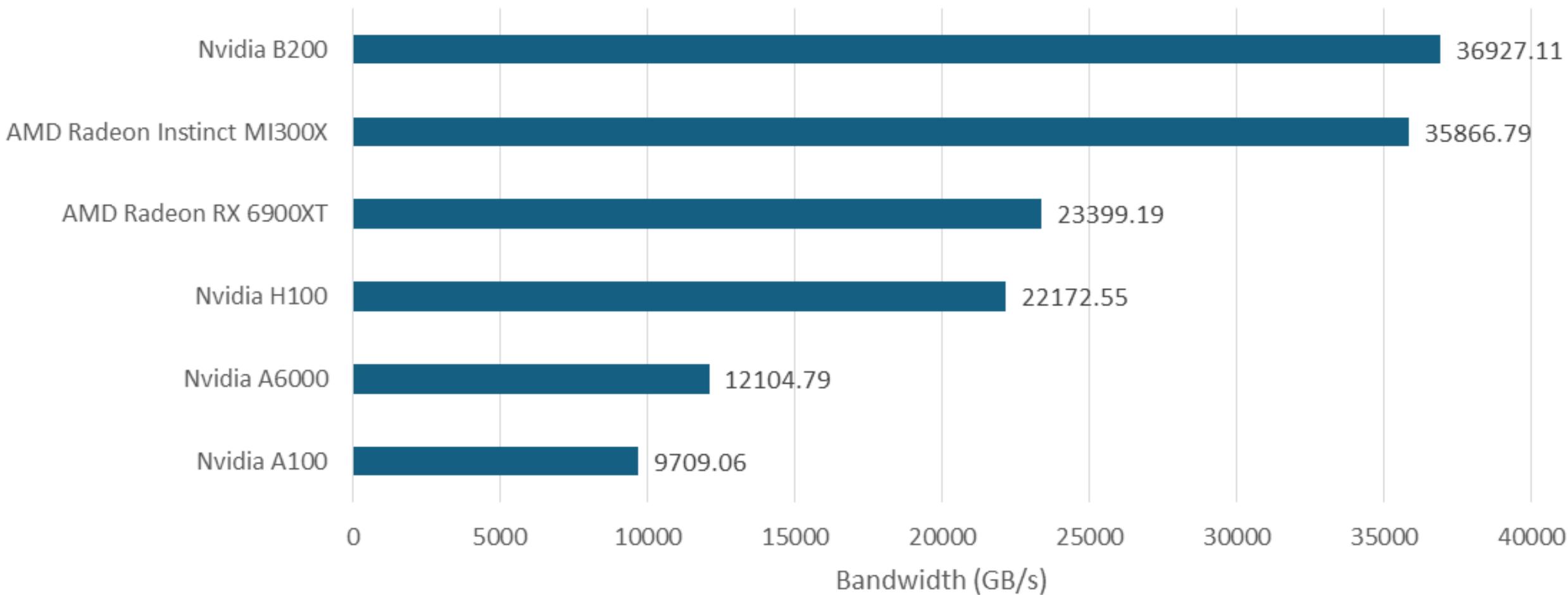
روندهای تکنولوژی

GPU

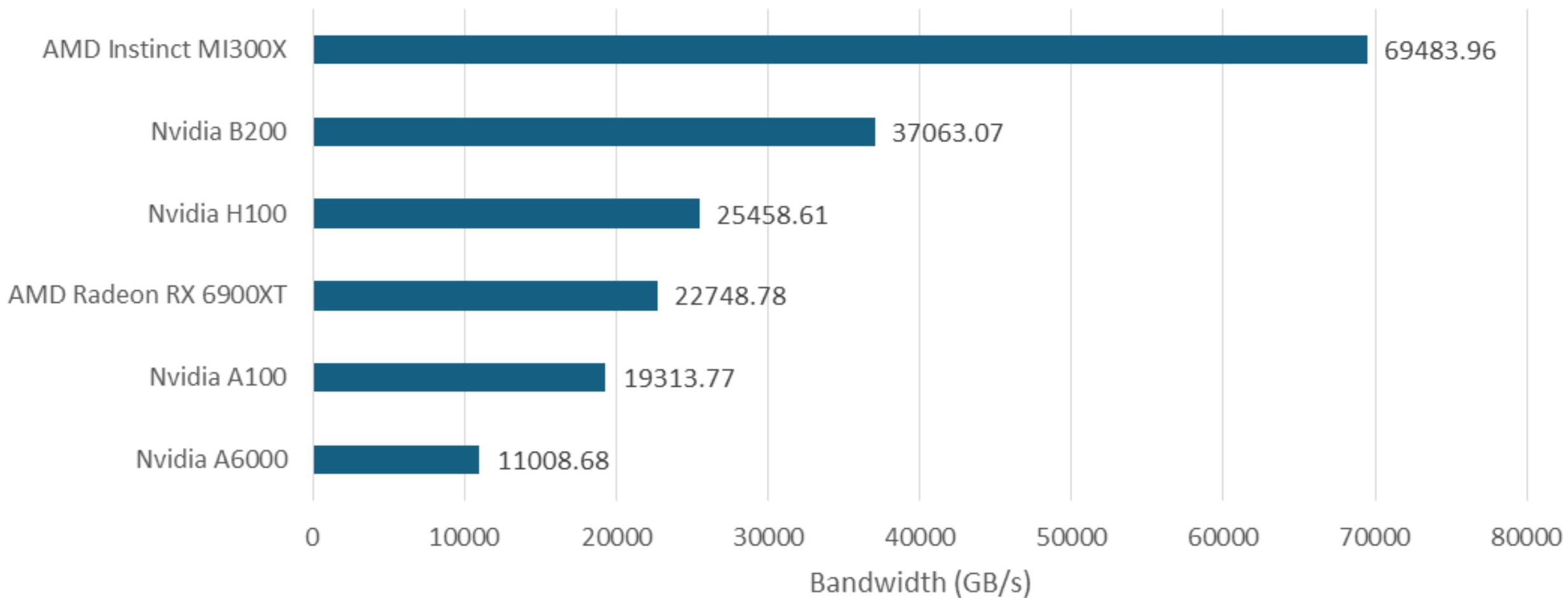
Local Memory Latency



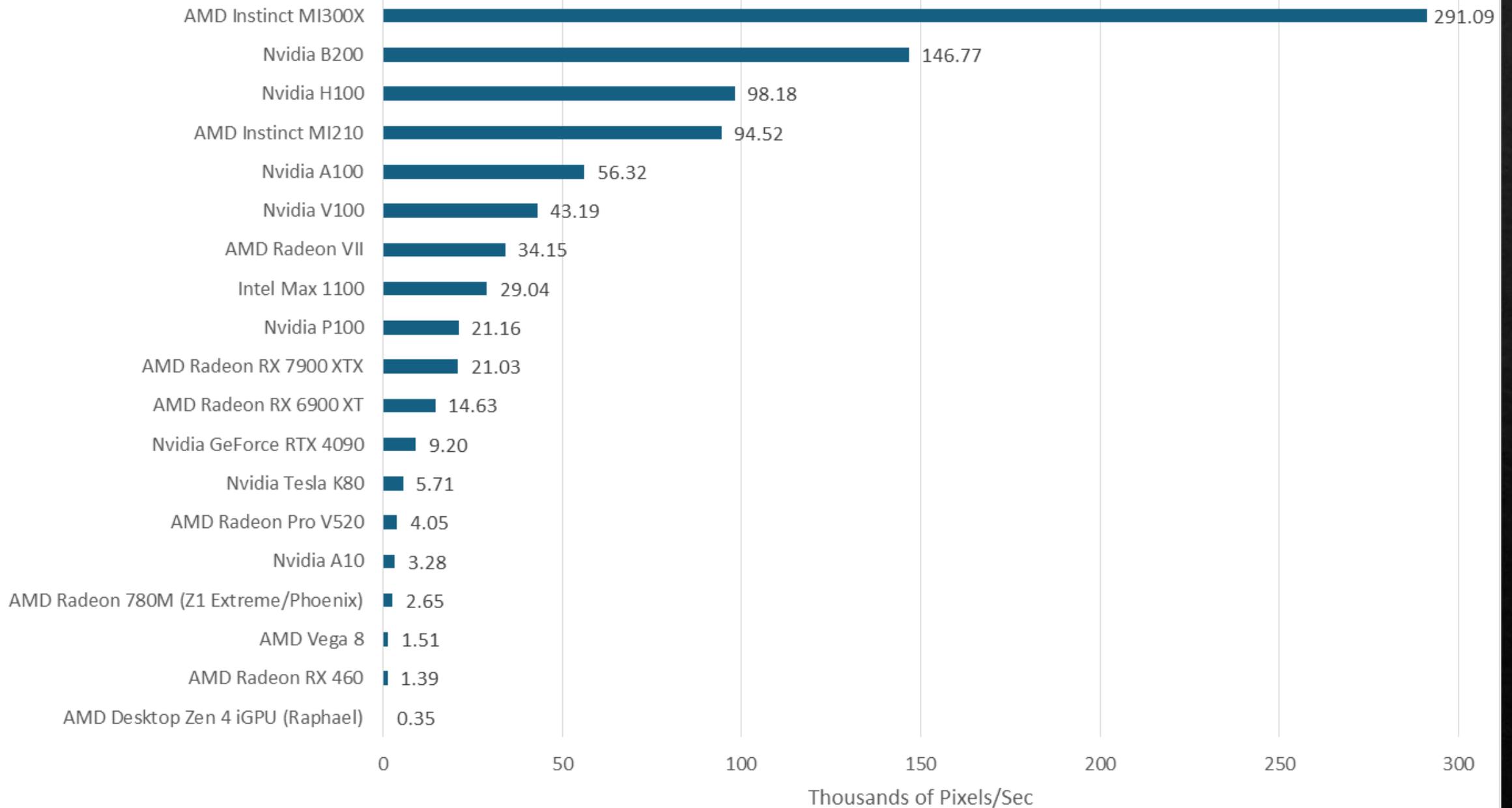
First Level Cache Bandwidth



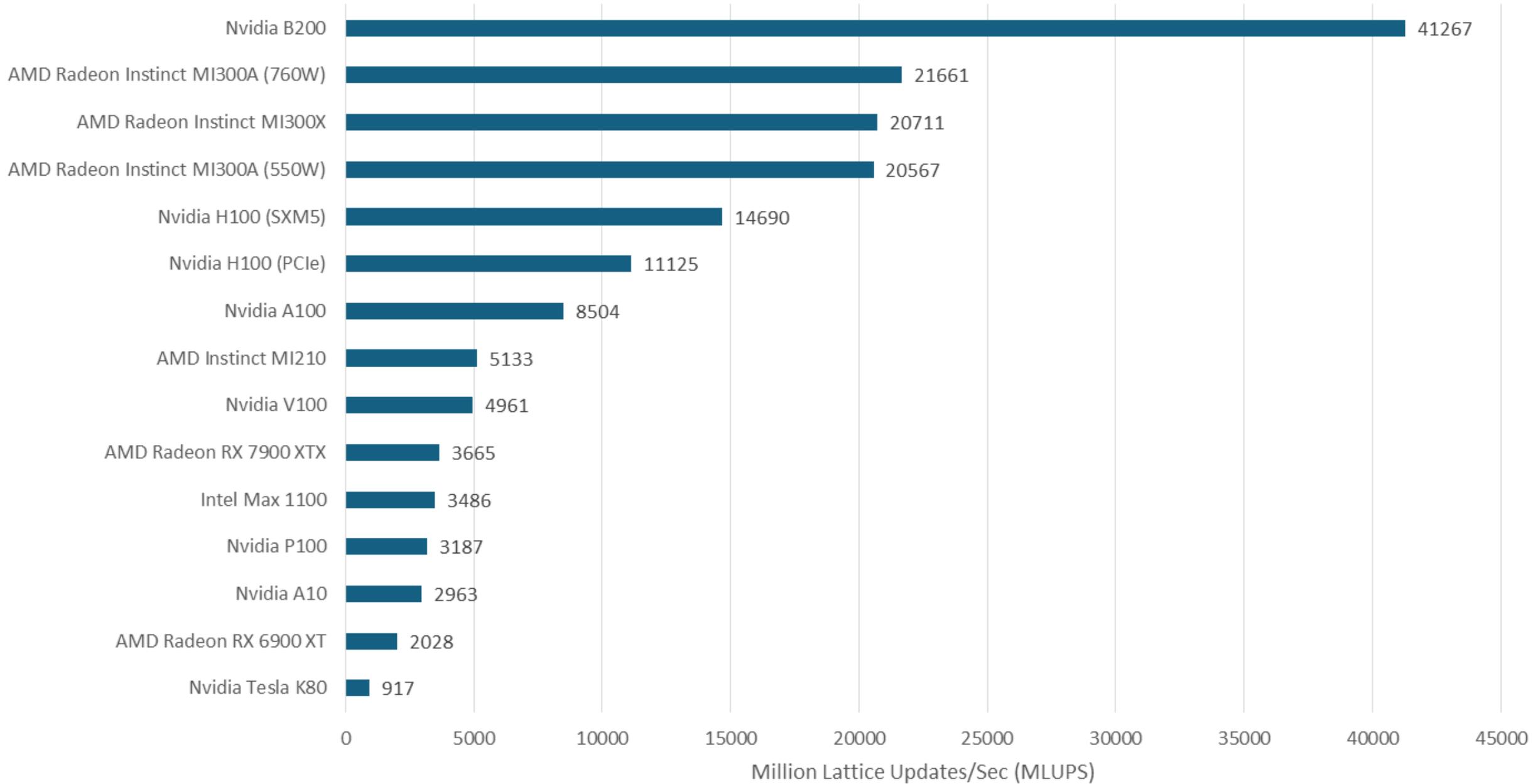
Local Memory Bandwidth



Compute Gravitational Potential (FP64)



FluidX3D, FP32



برنامه آینده

NVIDIA



*Oberon Rack
Liquid Cooled*

Blackwell Ultra NVL72

Second Half 2025

1.1 EF Dense FP4 Inference
0.36 EF FP8 Training
1.5X GB200 NVL72

New Attention Instructions
2X

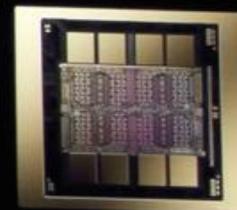
20 TB HBM | 40 TB Fast Memory
1.5X

14.4 TB/s CX8
2X

Grace



Blackwell Ultra



2 Reticle-Sized GPUs
1 Dense FP4 | 288GB HBM3e





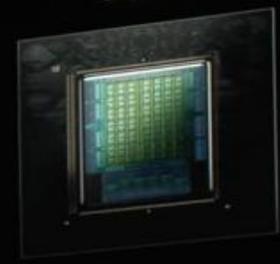
Oberon Rack
Liquid Cooled

Vera Rubin NVL144

Second Half 2026

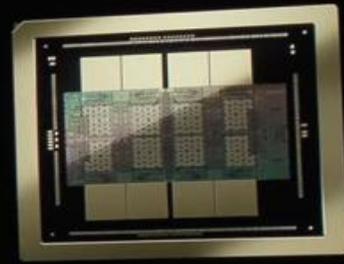
- 3.6 EF FP4 Inference
- 1.2 EF FP8 Training
- 3.3X GB300 NVL72
- 13 TB/s HBM4
- 75 TB Fast Memory
- 1.6X
- 260 TB/s NVLink6
- 2X
- 28.8 TB/s CX9
- 2X

Vera



- 88 Custom Arm Cores
- 176 Threads
- 1.8 TB/s NVLink-C2C

Rubin



- 2 Reticle-Sized GPUs
- 50PF FP4 | 288GB HBM4



*Kyber Rack
Liquid Cooled*

Rubin Ultra NVL576

Second Half 2027

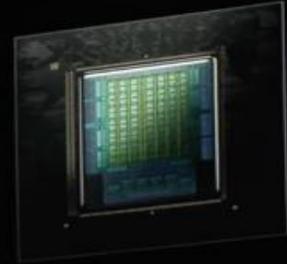
15 EF FP4 Inference
5 EF FP8 Training
14X GB300 NVL72

4.6 PB/s HBM4e
365 TB Fast Memory
8X

1.5 PBs NVLink7
12X

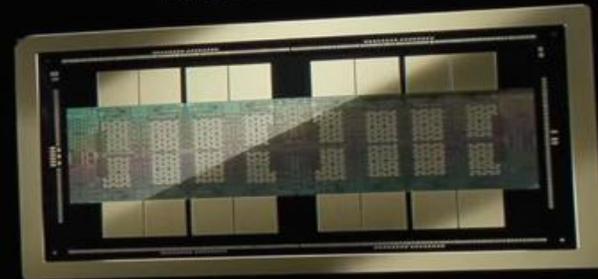
115.2 TB/s CX9
8X

Vera



88 Custom Arm Cores
176 Threads
1.8 TB/s NVLink-C2C

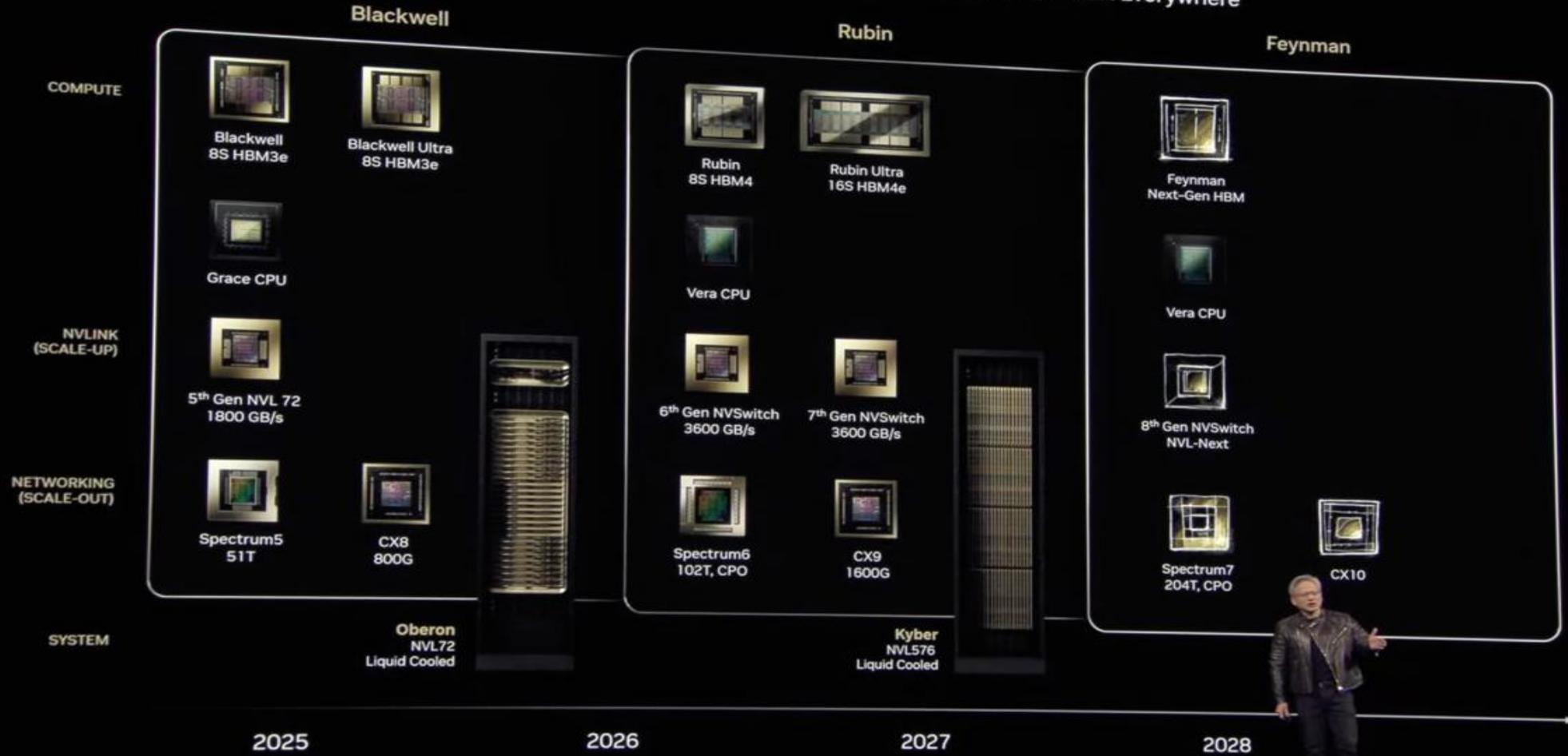
Rubin Ultra



4 Reticle-Sized GPUs
100PF FP4 | 1TB HBM4e

NVIDIA Paves Road to Gigawatt AI Factories

One-Year Rhythm | Full-Stack | One Architecture | CUDA Everywhere



NVIDIA AI Infrastructure for Enterprise Computing

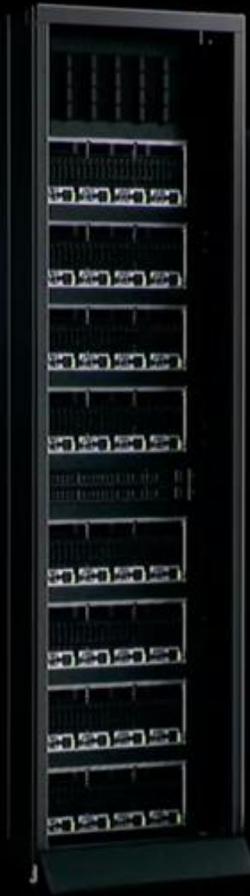
DGX Spark
1 PFLOPS



RTX PRO
Workstation



DGX
Station
20 PFLOPS



RTX PRO
Enterprise Server



DGX B200

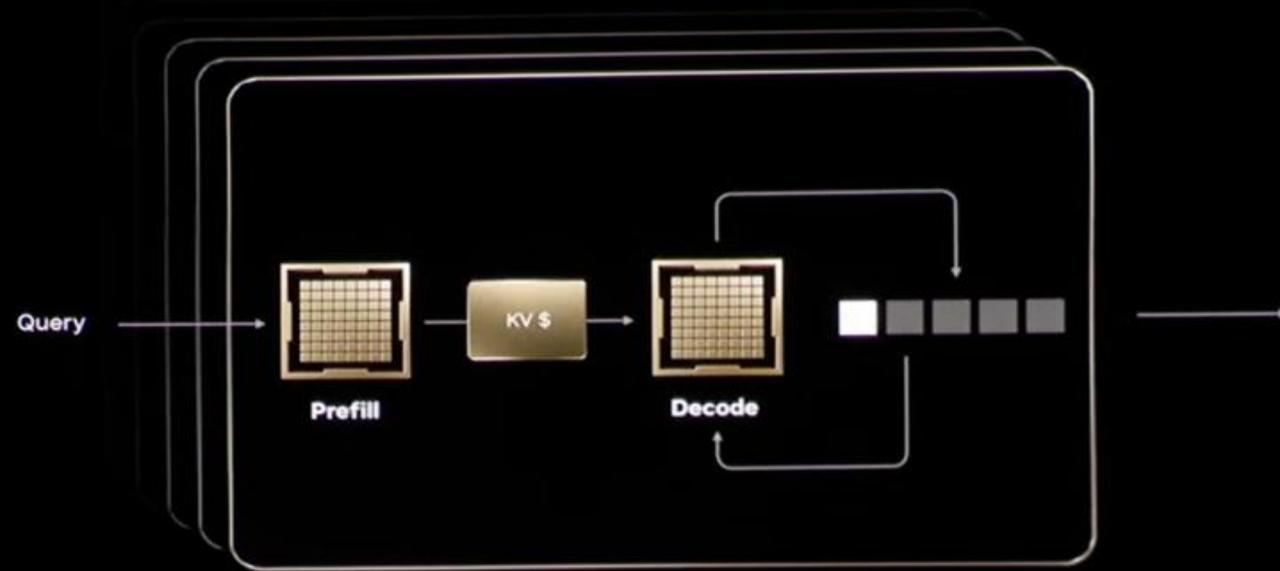


DGX GB300
15 EFLOPS



Announcing NVIDIA Dynamo

Distributed Inference Serving Library



Disaggregated Inference

GPU Resource Allocation

KV Cache Routing

Communication Library (NIXL)

cohere

Fireworks AI

Meta

Microsoft Azure

perplexity

PyTorch

SGL

together.ai

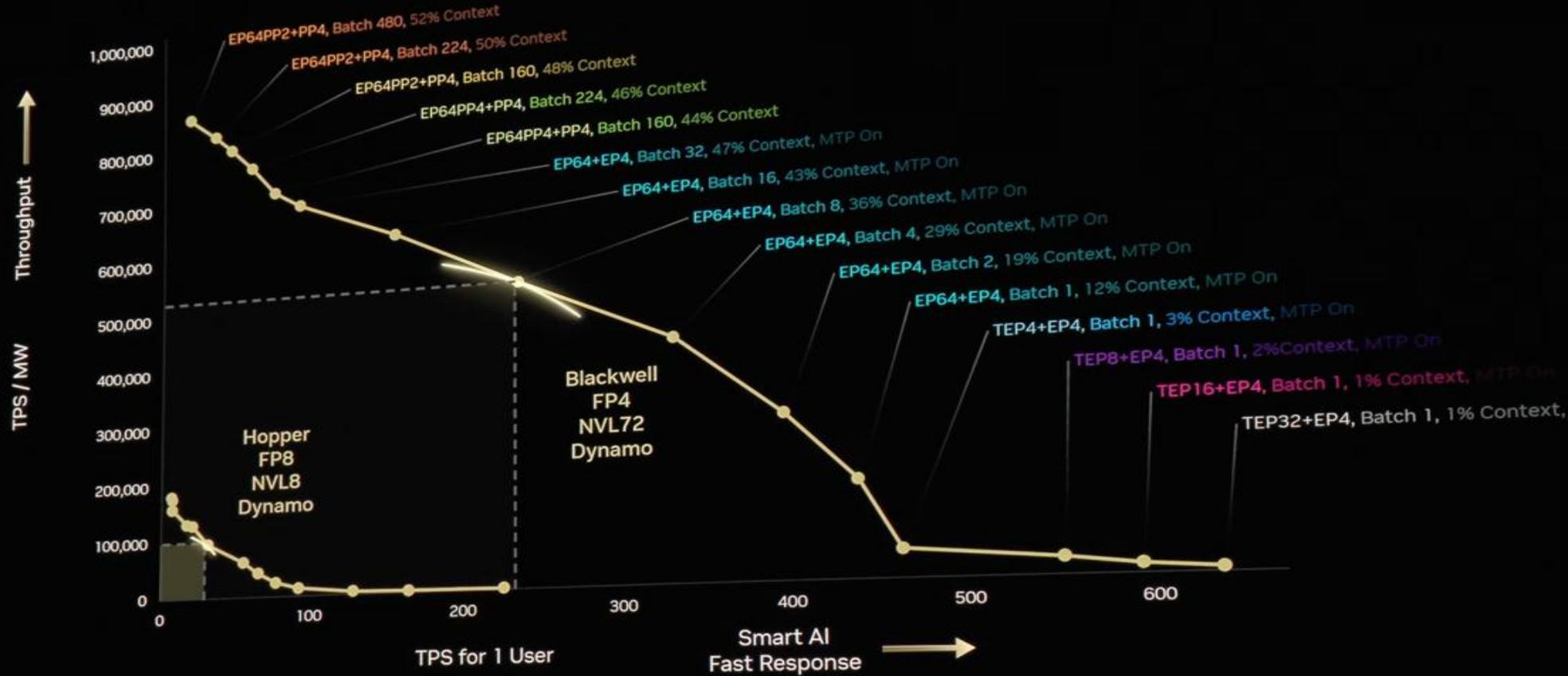
vAST

vLLM

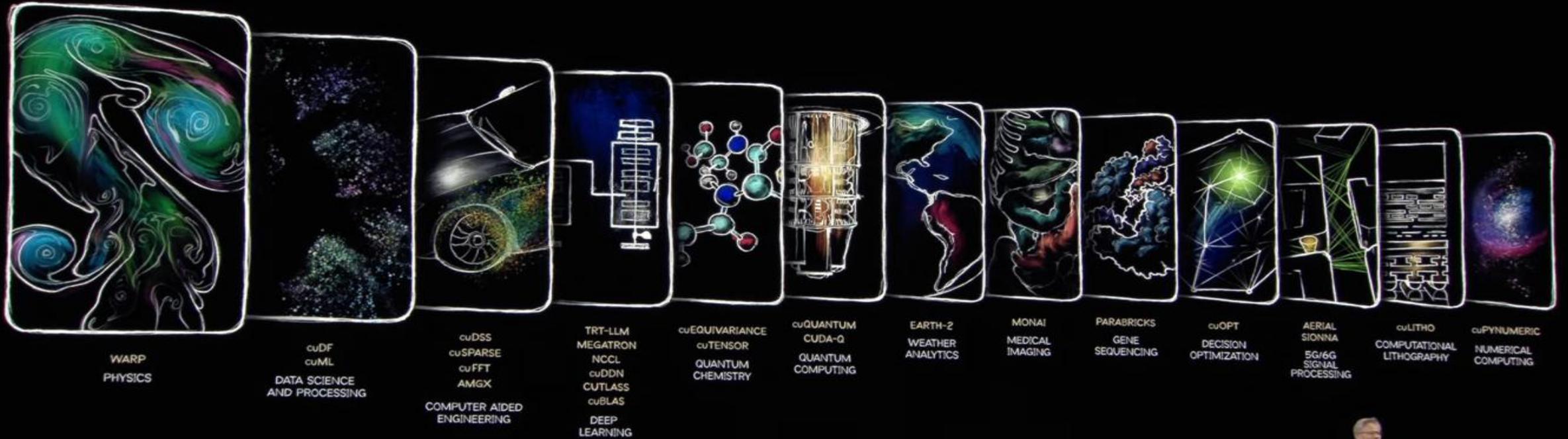


Blackwell 40X Hopper

FP4, NVL72, Dynamo, and TRT-LLM Continuous Optimization
32K ISL / 8K OSL



CUDA-X FOR EVERY INDUSTRY



LLM

به زبان ساده

Artificial intelligence (AI)

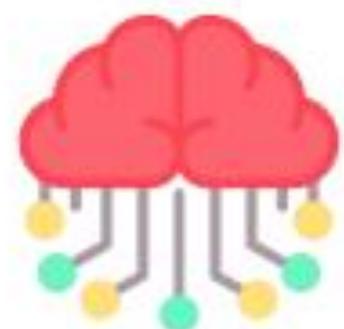
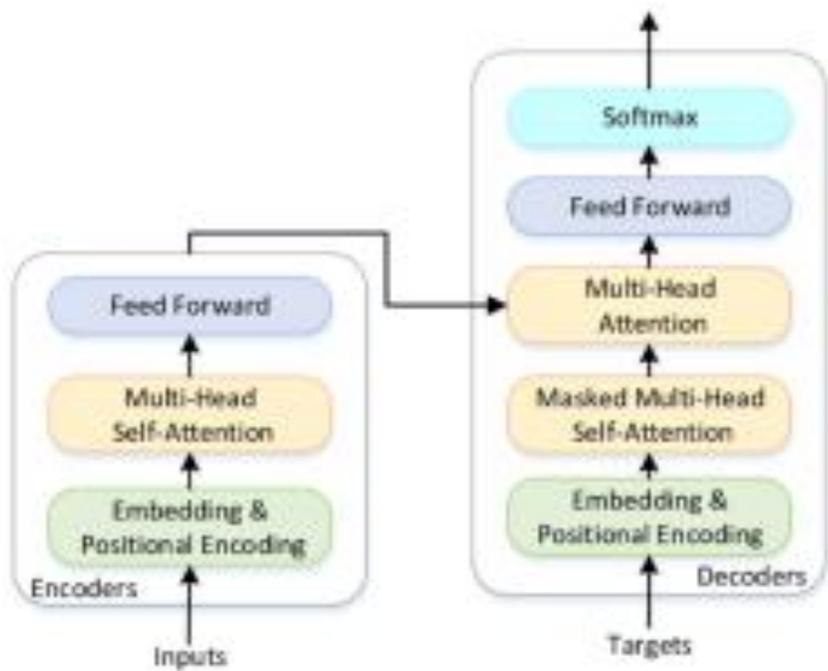
Machine learning (ML)

Natural language processing (NLP)

Large language models (LLMs)

Deep learning (DL)





Data



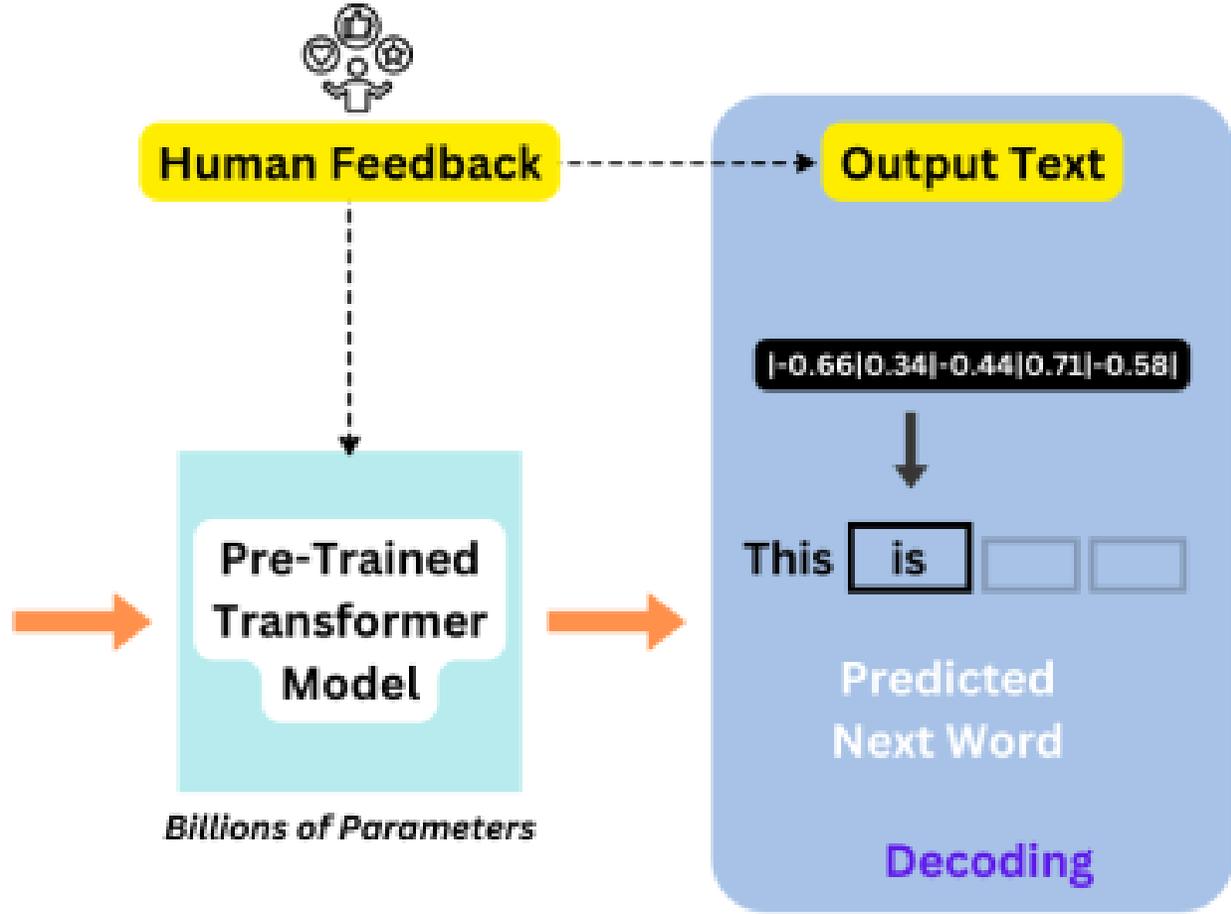
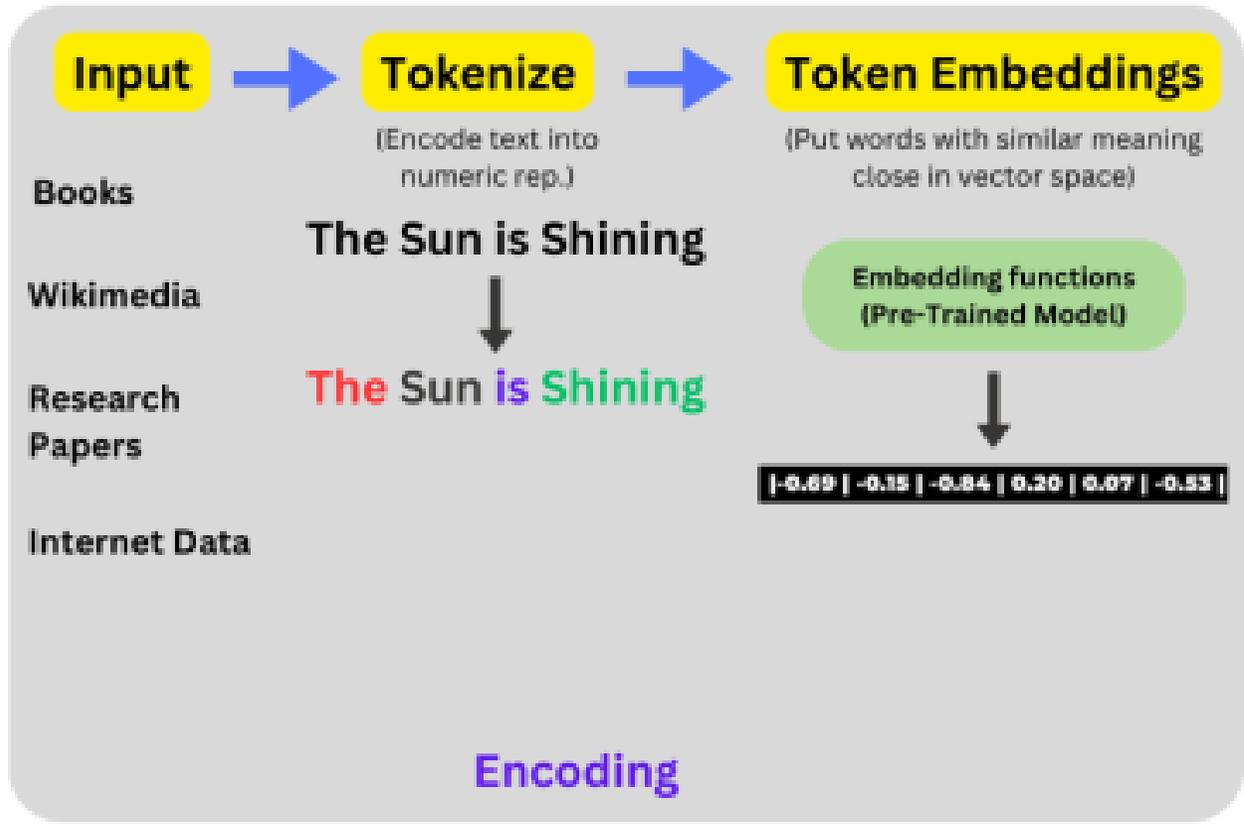
Architecture



Training



LLM





کاربرد های LLMs

Exceptional Large Language Model Use Cases



Application of LLMs

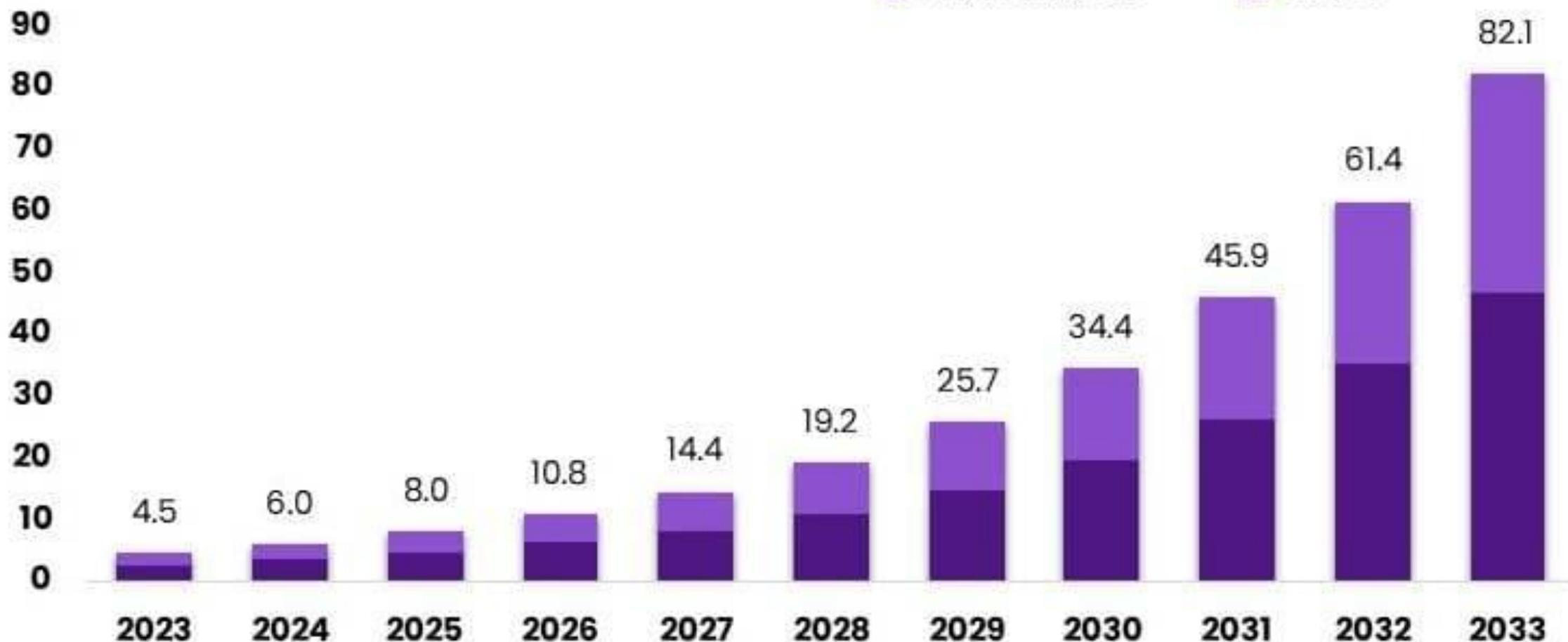


اهمیت LLMs

Global Large Language Model (LLM) Market

Size, by Deployment, 2023-2033 (USD Billion)

■ On-Premises ■ Cloud



The Market will Grow
At the CAGR of:

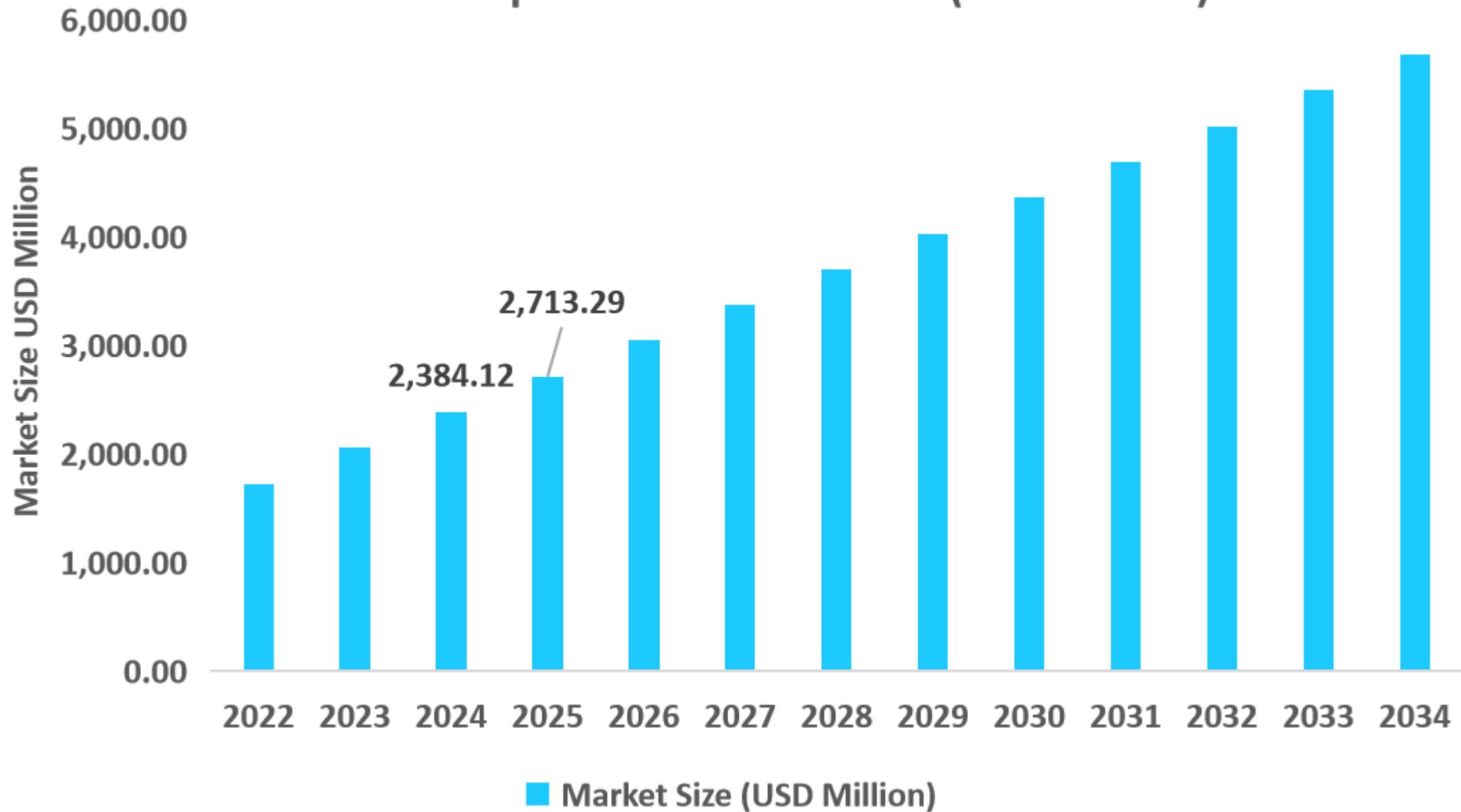
33.7%

The Forecasted Market
Size for 2033 in USD:

\$82.1 B

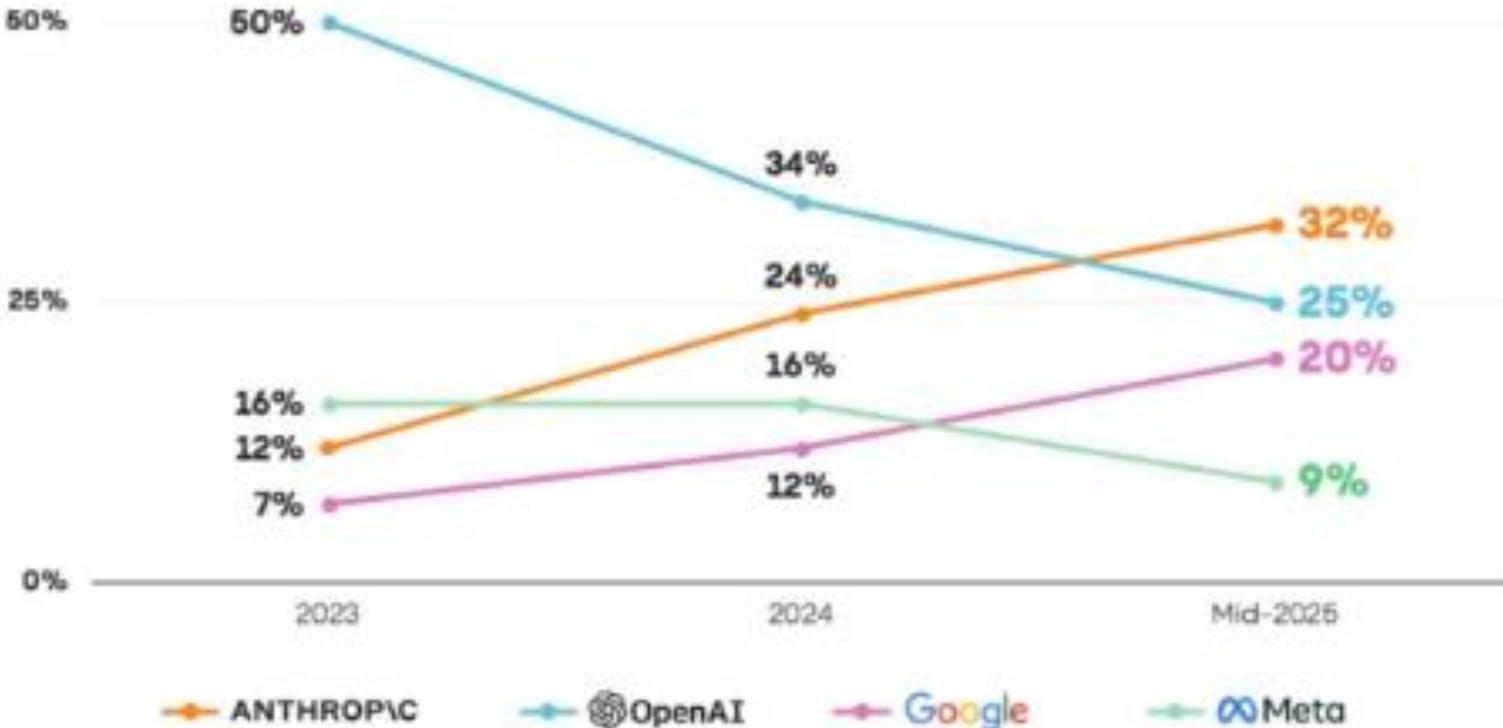
 **market.us**
ONE STOP SHOP FOR THE REPORTS

U.S Enterprise LLM Market Size (USD Million)

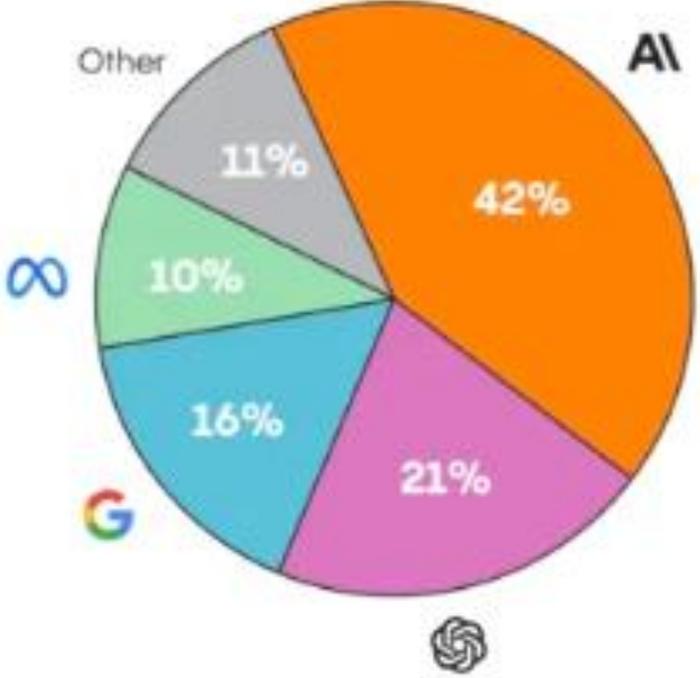


Enterprise LLM API Market Share by Usage

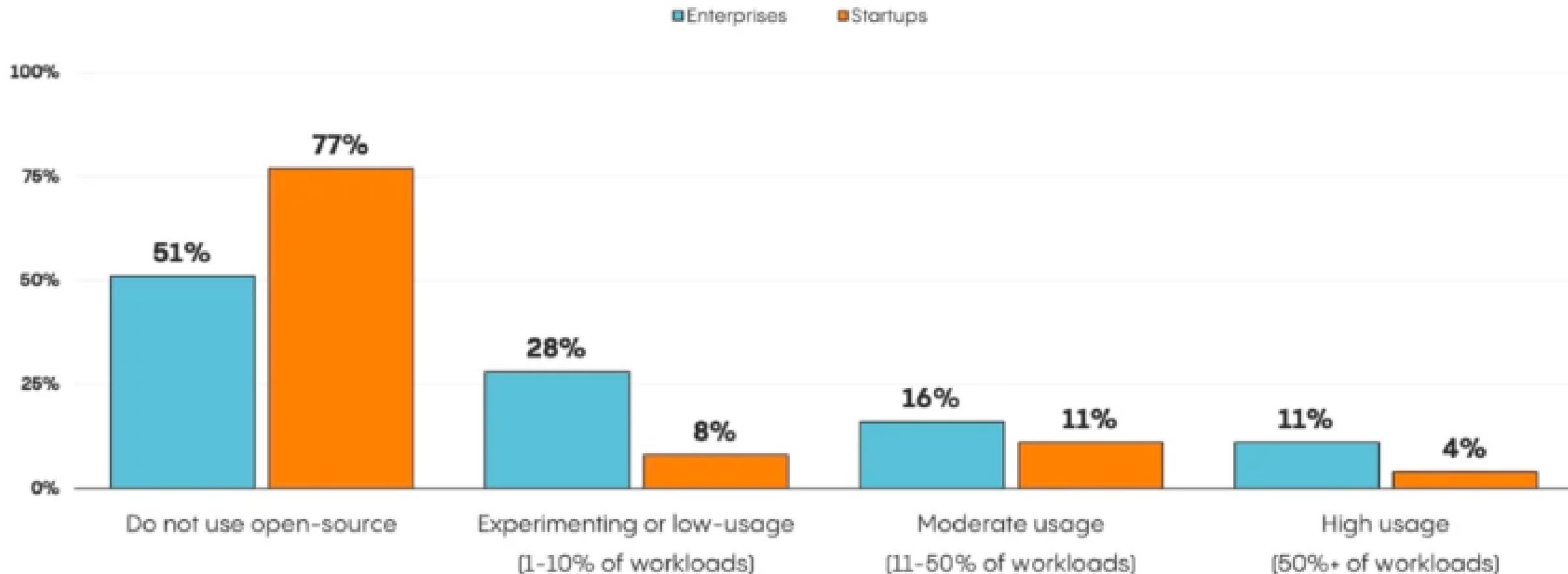
Change in Enterprise LLM API Market Share



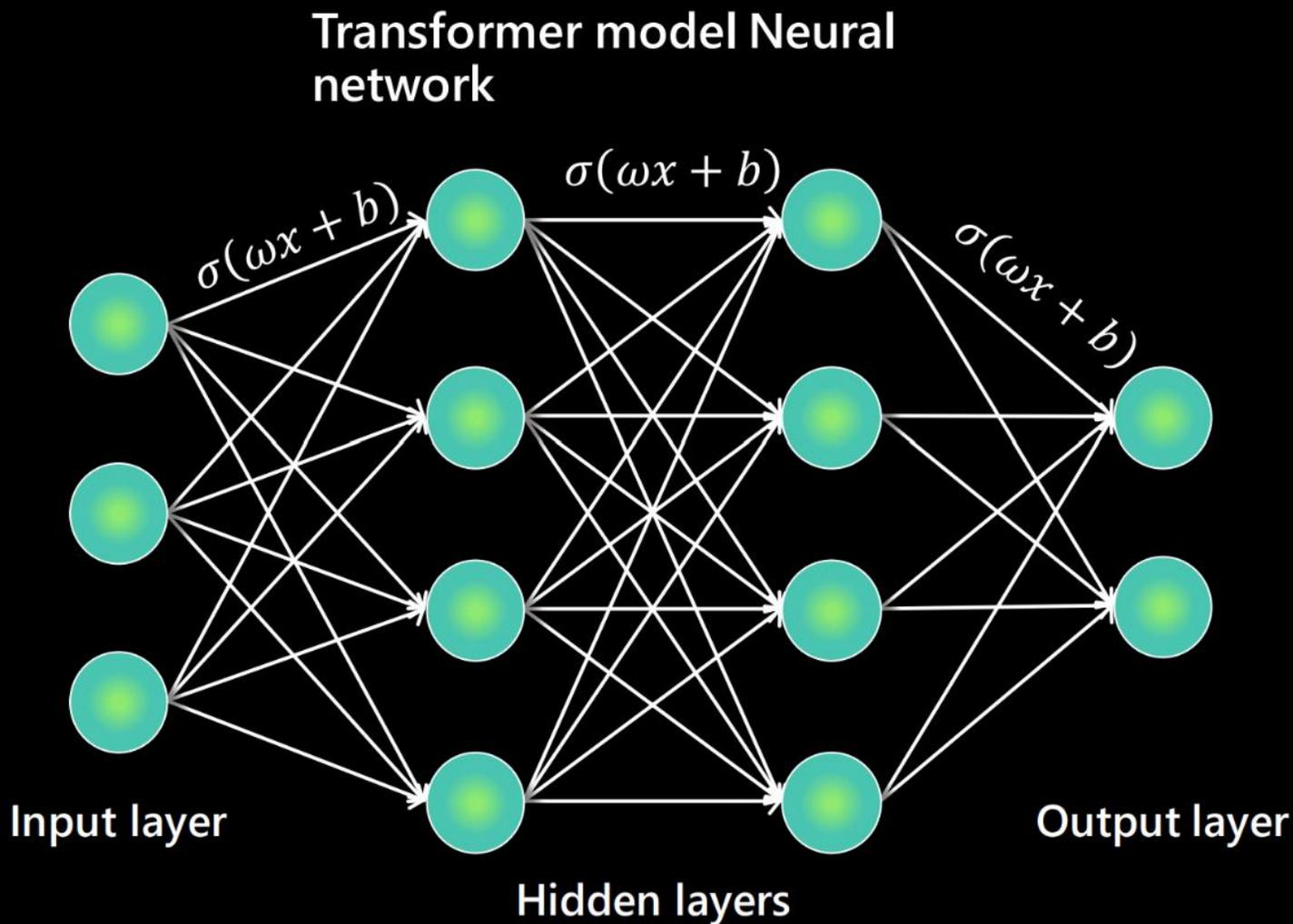
Coding Market Share



Companies Choose **Closed-Source**



How large are they?



BERT Large - 2018

345M

GPT2 - 2019

1.5B

GPT3 - 2020

175B

Turing Megatron NLG
2021

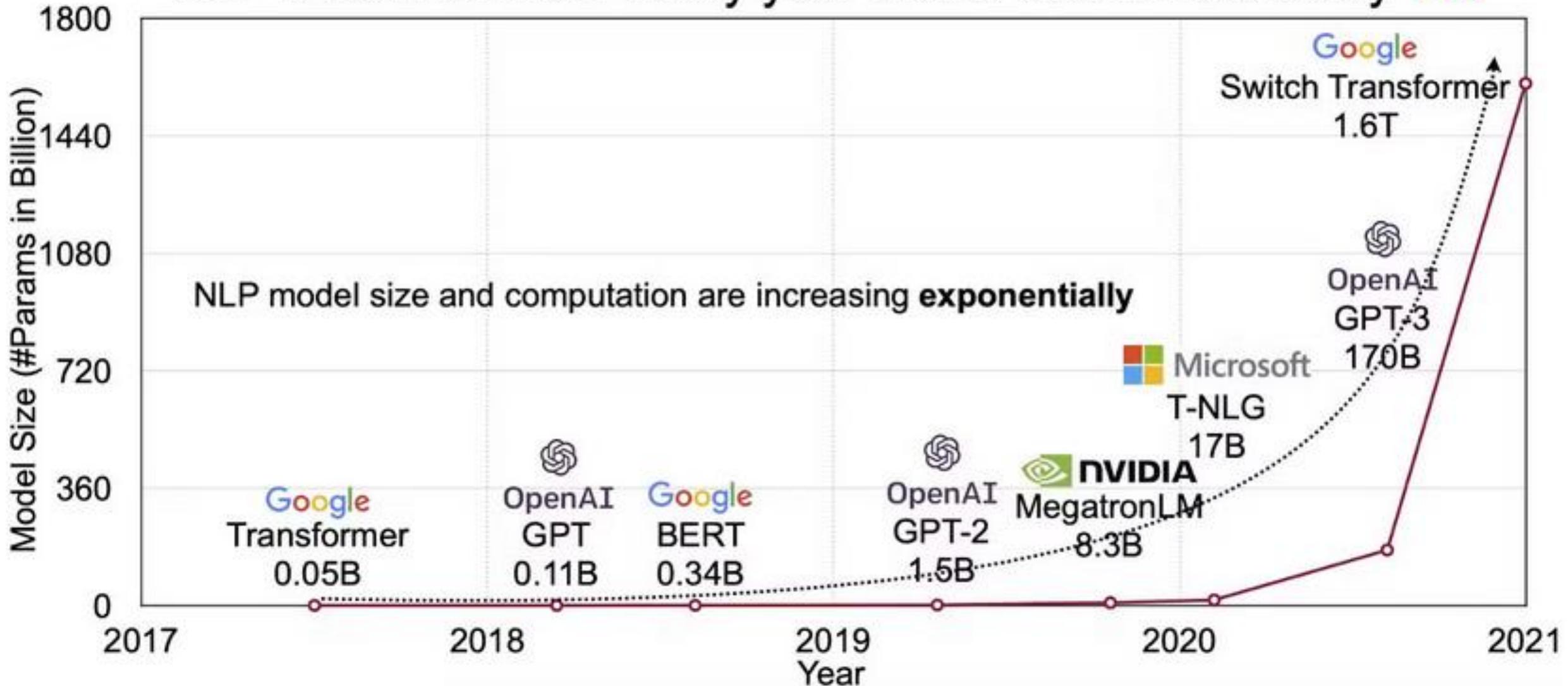
530B

GPT4 - 2023

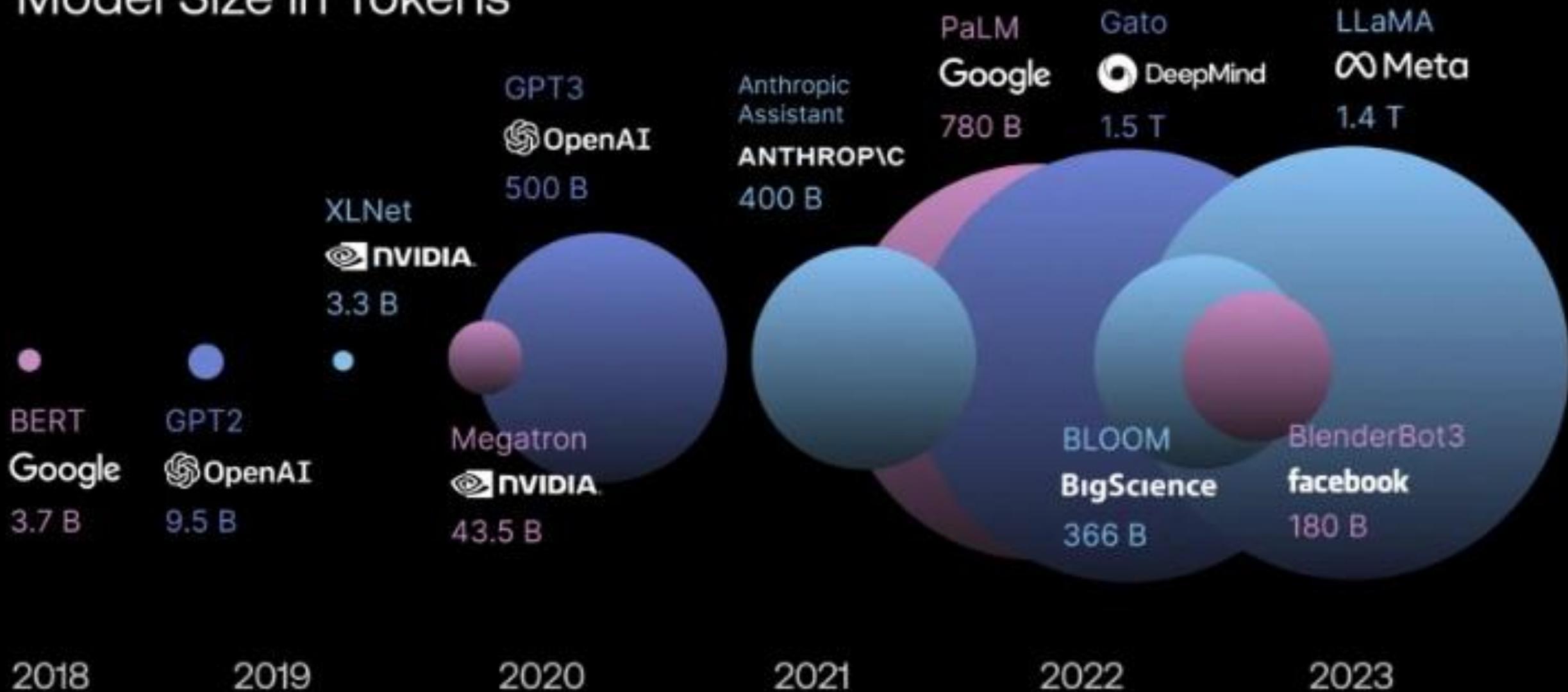
1.4T (estimated)

Rapid change in model size

NLP's Moore's Law: Every year model size increases by **10x**



Model Size in Tokens



CodeLlama

-  WizardCoder
-  Phind-CodeLlama
-  OpenAssistant-Codellama
-  NexusRaven

Mistral

-  Mistralic
-  SQLCoder-7B
-  Mistral-OpenOrca

StarCoder

-  SQLCoder-15B
-  WizardCoder-15B

Llama2

-  OpenChat
-  CodeLlama
-  Nous-Hermes
-  Vicuna
-  NSQL
-  Redmond-Puffin
-  Llama2-32K
-  OpenAssistant-Llama2
-  WizardLM

GPT

-  GPT-3.5-Turbo
-  GPT-4

Stable

-  StableLM
-  StableCode
-  OpenAssistant-StableLM

Pythia

-  Pythia-ChatBase
-  OpenAssistant-Pythia

Claude

-  Claude-2

PaLM

-  PaLM-2

Falcon

-  Falcon

Qwen

-  Qwen

Map of LLMs October 2023



Top 10 Large Language Models in 2025



Top 15 LLM to boost your performance, capabilities and innovation



GPT-4

GPT-4: OpenAI



PaLM 2

PaLM 2: Google



Orca: Microsoft



GPT-3

GPT-3: OpenAI



Bard

Bard: Google



Guanaco

Guanaco:
HuggingFace/TUM



GPT-3.5

GPT-3.5: OpenAI



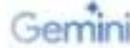
Claude

Claude v1: Anthropic



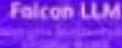
VICUNA

Vicuna: LMSYS



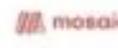
Gemini

Gemini:
Google DeepMind



Falcon LLM

Falcon: TII (UAE)



mosaicML

MPT-30B: MosaicML



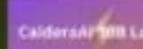
LLaMA
by Meta

LLaMA: Meta AI



cohere

Cohere: Cohere.ai



CalderaAI

30B Lazarus:
CalderaAI

LLM Market Share

LLM	Global	US	India	Germany	France	Taiwan	Hong Kong	Russia
ChatGPT	78%	78%	80%	81%	87%	79%	49%	57%
DeepSeek	8%	4%	5%	6%	2%	2%	24%	25%
Gemini	6%	7%	6%	5%	5%	9%	4%	2%
Grok	2%	2%	4%	1%	0%	3%	6%	3%
Perplexity	2%	3%	2%	3%	2%	2%	11%	5%
Claude	2%	3%	2%	1%	1%	4%	1%	0%
Copilot	1%	2%	1%	2%	2%	1%	4%	0%
Qwen	0%	0%	0%	0%	0%	0%	0%	5%
Mistral	0%	0%	0%	0%	1%	0%	0%	2%
Hey & Go	0%	0%	0%	1%	0%	0%	0%	0%

Sessions (web+mobile), SimilarWeb, without apps, without APIs, Jan-Apr 2025, Malte Landwehr

LLMs: NEXT PUBLIC RELEASES IN 2026

ESTIMATES AS OF
DECEMBER 2025

Jan-Mar

Meta AI
Avocado

xAI
Grok-5 (6T)

Google DeepMind
Gemma 4

Apr-Jun

Anthropic
Claude 5

Google DeepMind
Gemini 4

Jul-Sep

Meta AI
Next

xAI
Grok-6

Anthropic
Claude 5.5

OpenAI
GPT-6

Microsoft
MAI-2

Oct-Dec

OpenAI
Next

Google DeepMind
Gemma 5

Baidu
ERNIE 6

Estimates only, selected highlights only. Full models table at: <https://lifearchitct.ai/models-table/> Alan D. Thompson, December 2025. <https://lifearchitct.ai/>



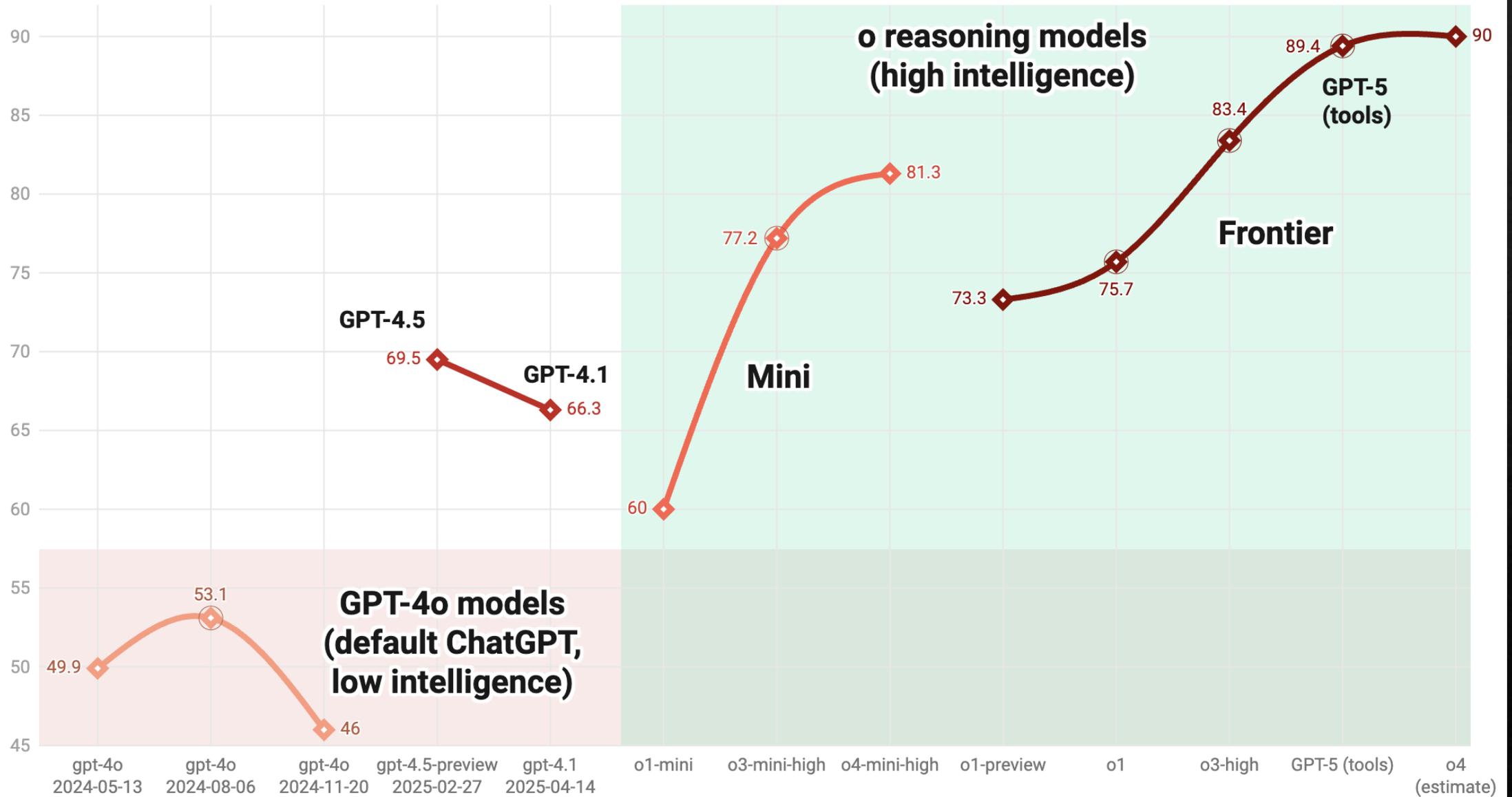
Model	Training end	Chip type	TFLOP/s (max)	Chip count	Wall clock (days)	Total time (years)	Retail (US\$)	MMLU
GPT-3 175B	Apr/2020	V100	130	10,000	15 days	405y	\$9M	43.9
Llama 1 65B	Jan/2023	A100	312	2,048	21 days	118y	\$4M	63.4
Llama 2 70B	Jun/2023	A100	312	2,048	35 days	196y	\$7M	68.0
Titan 200B	Apr/2023	A100	312	13,760	48 days	1,319y	\$45M	70.4
GPT-4 1.7T	Aug/2022	A100	312	25,000	95 days	6,507y	\$224M	86.4
Gemini	Nov/2023	TPUv4	275	57,000	100 days	15,616y	\$440M	90.0
Llama 3 405B	Apr/2024	H100	989	24,576	50 days	3,366y	\$125M	85+
GPT-5	Apr/2024	H100	989	50,000	120 days	16,438y	\$612M	
Grok 2	Jun/2024	H100	989	20,000	50 days	6,571y	\$245M	
Olympus	Aug/2024	H100	989					
Gemini 2	Nov/2024	TPUv6	1,847					
Grok 3	Dec/2024	H100	989	100,000	50 days	32,855y	\$1.2B	

Alan D. Thompson. May/2024. LifeArchitect.ai

Table. Model training compute (see working, with sources⁸).

The GPT-4x and o Model Family: Varied Intelligence Scores 2024–2025

OpenAI's GPT-4x and o models on GPQA scores • LifeArchitect.ai

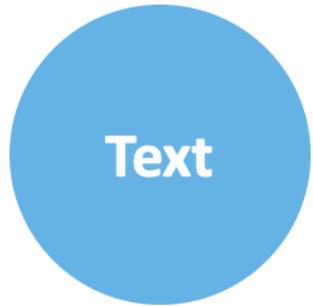


Source: LifeArchitect.ai/GPT-5 based on official eval data from GitHub.com/openai/simple-evals

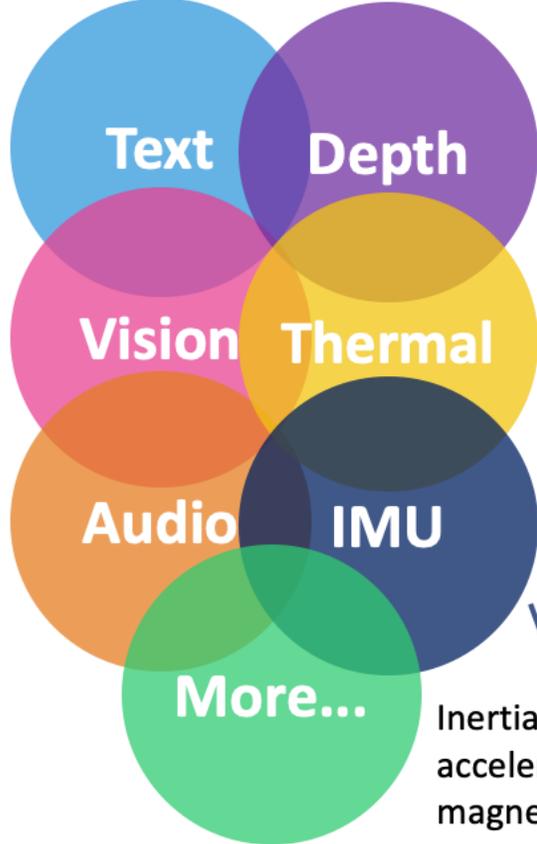
Chart: Alan D. Thompson • Aug/2025 • Source: LifeArchitect.ai • Get the data • Created with Datawrapper

GPT MODALITIES (2024)

GPT-3

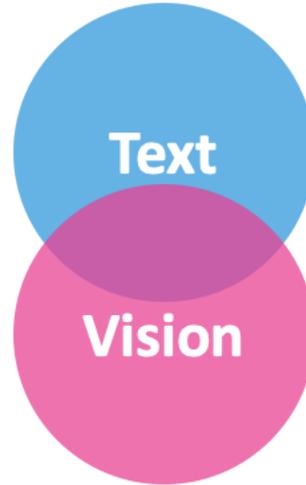


PandaGPT,
Meta-Transformer,
NEXT-GPT, 4M-21...

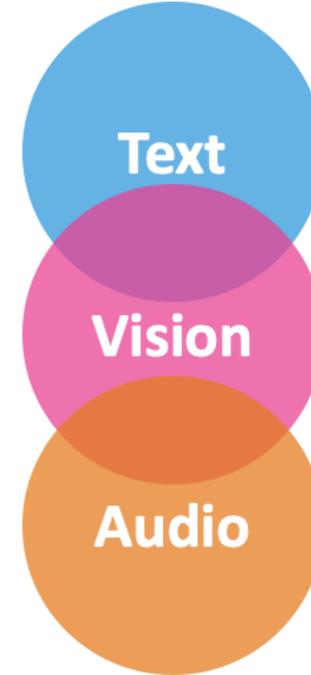


Inertial measurement unit:
accelerometer, gyroscope,
magnetometer or compass...

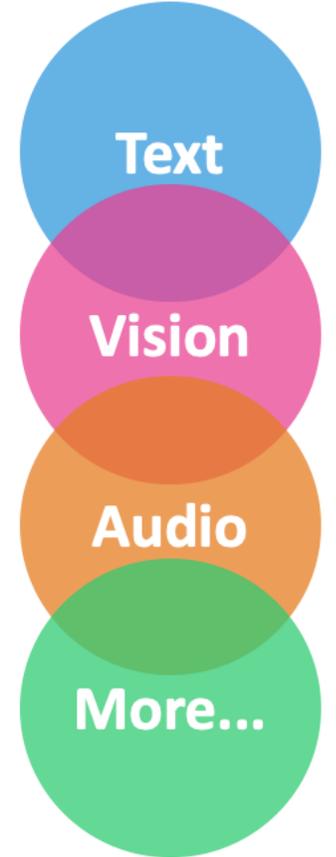
GPT-4



GPT-4o



GPT-5



FRONTIER AI MODELS + HIGHLIGHTS

2026
H1



poe.com



Google Gemma
Alibaba Qwen/QwQ
Mistral
Microsoft phi

+ hundreds more...



Meta AI Llama
OpenAI GPT-OSS
Kimi
DeepSeek R

+ hundreds more...

Selected highlights only. Model highlights: <https://life architect.ai/models-table/> Total models: https://huggingface.co/models?pipeline_tag=text-generation Some images by Flaticon.com. Alan D. Thompson. 2021-2026.



LifeArchitect.ai/models-table

(700 model highlights of 300,000 total models)

LARGE LANGUAGE MODEL MARKET MAP

LARGE LANGUAGE MODEL SOFTWARE PROVIDERS

LLM APIs	VECTOR DATABASES	LLM FRAMEWORKS	TEXT-TO-SPEECH	LLM MONITORING TOOLS
 OpenAI ANTHROPIC cohere	  Pinecone  Chroma	 LlamaIndex  LangChain  FIXIE	 RESEMBLE.AI ElevenLabs WELLSAID	 DISTYL  Guardrails AI  Helicone

LARGE LANGUAGE MODEL SERVICE PROVIDERS

COMPUTE PLATFORM PROVIDERS	MODEL HUBS	FINE-TUNING/CUSTOM MODEL TRAINING FRAMEWORKS	MONITORING/OBSERVABILITY PLATFORM PROVIDERS	HOSTING SERVICE PROVIDERS
 Lambda  mosaicML  Azure	 Hugging Face  Replicate	 PyTorch  TensorFlow LAMINI	 Arthur  arize WHYLABS	 Replicate  Hugging Face

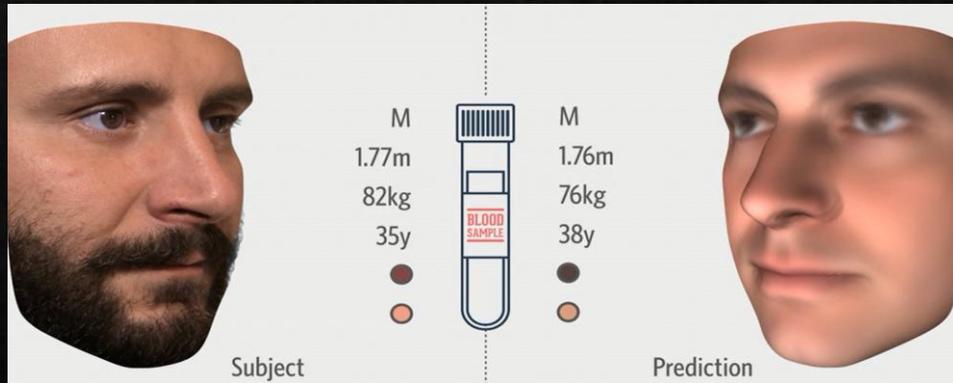
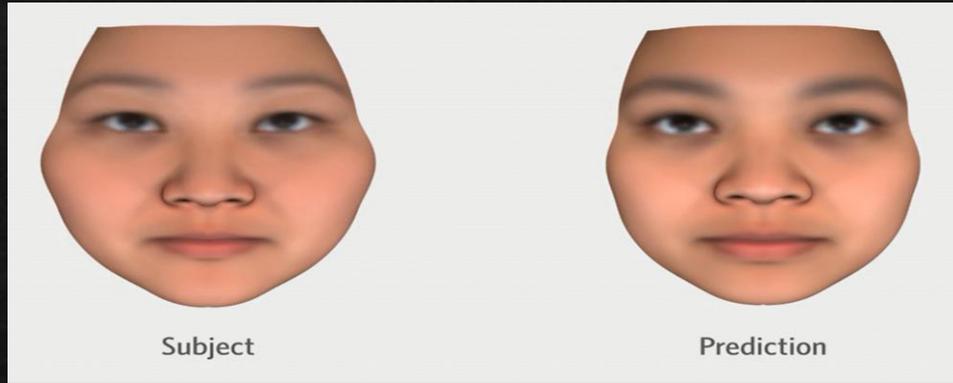
END-USERS

GOVERNMENT & REGULATORY BODIES

 edger.finance  summer health	 BEN™ BRANDED ENTERTAINMENT NETWORK  OXIDE.AI®	 NIST  ico. Information Commissioner's Office	 European Commission  NiTDA
--	--	---	--

چالش های هوش مصنوعی

شرح مساله - پیش بینی چهره در سال ۲۰۱۶



دستکاری ژنوم انسانی و ساخت یک موجود جدید!

