

[Korea University AI Seminar]

# The AI Architecture Decoded - Orchestrating Technology and Hardware Innovation

Sunghee Yun

Co-founder & CTO - AI Technology & Biz Dev @ [Erudio Bio, Inc.](#)

## About Speaker

- *Co-founder & CTO @ Erudio Bio, San Jose & Novato, CA, USA*
- Advisor & Evangelist @ CryptoLab, Inc., San Jose, CA, USA
- Chief Business Development Officer @ WeStory.ai, Cupertino, CA, USA
- Advisory Professor, Electrical Engineering and Computer Science @ DGIST, Korea
- Adjunct Professor, Electronic Engineering Department @ Sogang University, Korea
- Global Advisory Board Member @ Innovative Future Brain-Inspired Intelligence System Semiconductor of Sogang University, Korea
- *KFAS-Salzburg Global Leadership Initiative Fellow @ Salzburg Global Seminar*, Salzburg, Austria
- Technology Consultant @ Gerson Lehrman Group (GLG), NY, USA
- *Co-founder & CTO & Head of Global R&D & Chief Applied Scientist & Senior Fellow @ Gauss Labs, Inc., Palo Alto, CA, USA* *2020 – 2023*

- Senior Applied Scientist @ Mobile Shopping Team, Amazon.com, Inc., Vancouver, BC, Canada – 2020
- Principal Engineer @ Software R&D Center of DS Division, Samsung, Korea – 2017
- Principal Engineer @ Strategic Marketing & Sales Team, Samsung, Korea – 2016
- Principal Engineer @ DT Team of DRAM Development Lab, Samsung, Korea – 2015
- Senior Engineer @ CAE Team - Samsung, Korea – 2012
- MS & PhD - Electrical Engineering @ Stanford University, CA, USA – 2004
- Development Engineer @ Vyan, Santa Clara, CA, USA – 2001
- BS - Electrical Engineering @ Seoul National University, Seoul, Korea – 1998

## Highlight of Career Journey

- BS in EE @ SNU, MS & PhD in EE @ Stanford University
  - *Convex Optimization - Theory, Algorithms & Software*
  - advised by *Prof. Stephen P. Boyd*
- Principal Engineer @ Samsung Semiconductor, Inc.
  - AI & Convex Optimization
  - collaboration with *DRAM/NAND Design/Manufacturing/Test Teams*
- Senior Applied Scientist @ Amazon.com, Inc.
  - e-Commerce AIs - time-series anomaly detection, deep reinforcement learning & recommender system
  - Jeff Bezos's project - increase sales by *\$200M* via Amazon Mobile Shopping App
- Co-founder & CTO & Head of Global R&D & Chief Applied Scientist & Senior Fellow @ Gauss Labs, Inc.
- Co-founder & CTO - AI Technology & Business Development @ Erudio Bio, Inc.

# Today

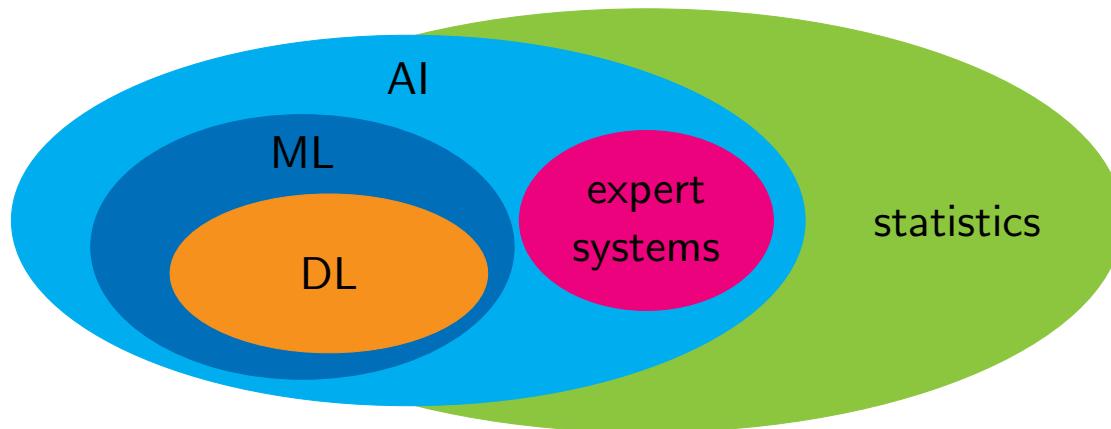
- Artificial Intelligence - 5
  - history & recent significant achievements
  - industry & market indices, is AI hype?
- AI Agents - 25
  - LLM & genAI, future of society powered by AI agents
- AI Hardware - 32
  - AI hardware industry, GPUs and AI accelerators
  - accelerator startups
- Global Semiconductor Industry - 48
  - US-China trade war, US re-shoring of manufacturing
- Appendices - 56
  - LLM
  - genAI
- Selected references - 107
- References - 109

# **Artificial Intelligence**

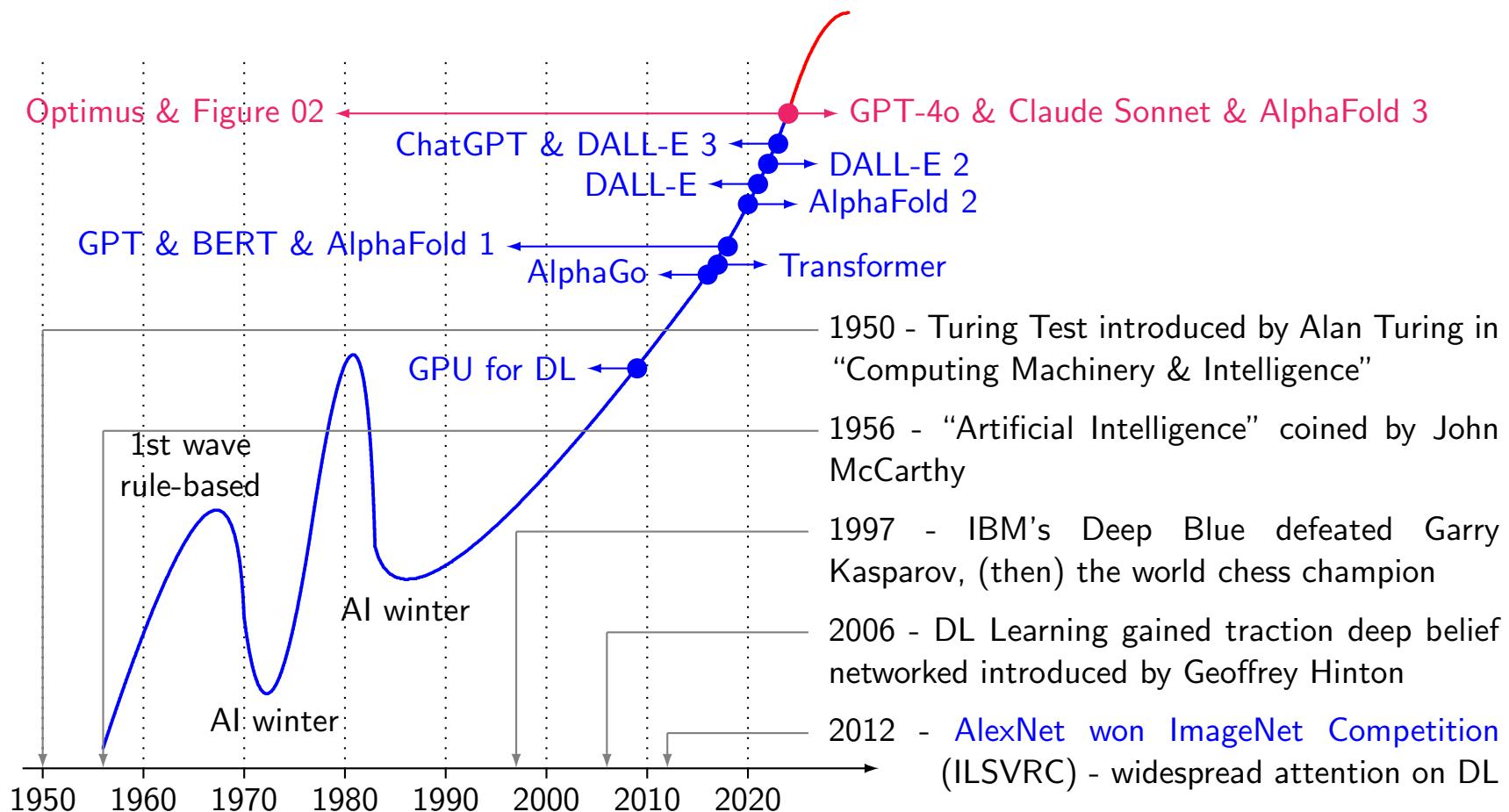
## **Definition and History**

## Definition of AI

- AI is
  - technology enabling machines to do tasks requiring human intelligence, such as learning, problem-solving, decision-making & language understanding
  - *not one thing* - encompass range of technologies, methodologies & applications
- relationship of AI, statistics, ML, DL, NN & expert system [HGH<sup>+</sup>22]



# History of AI



# **Significant AI Achievements - 2014 – 2024**

## Deep learning revolution

- 2012 – 2015 - DL revolution<sup>1</sup>
  - CNNs demonstrated exceptional performance in image recognition, e.g., *AlexNet's victory in ImageNet competition*
  - widespread adoption of DL learning in CV transforming industries
- 2016 - AlphaGo defeats human Go champion
  - DeepMind's AlphaGo defeated world champion in Go, extremely complex game *believed to be beyond AI's reach*
  - significant milestone in RL - AI's potential in solving complex & strategic problems



<sup>1</sup>DL: deep learning, CNN: convolutional neural network, CV: computer vision, RL: reinforcement learning

## Transformer changes everything

- 2017 – 2018 - Transformers & NLP breakthroughs<sup>2</sup>
  - *Transformer (e.g., BERT & GPT) revolutionized NLP*
  - major advancements in, *e.g.*, machine translation & chatbots
- 2020 - AI in healthcare – AlphaFold & beyond
  - DeepMind's *AlphaFold solves 50-year-old protein folding problem* predicting 3D protein structures with remarkable accuracy
  - accelerates drug discovery and personalized medicine - offering new insights into diseases and potential treatments



<sup>2</sup>NLP: natural language processing, GPT: generative pre-trained transformer

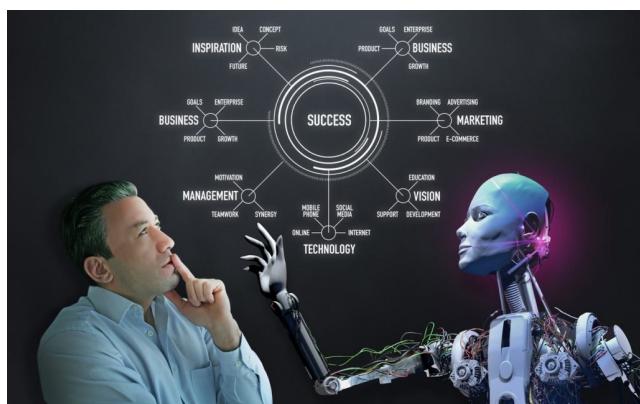
## Lots of breakthroughs in AI technology and applications in 2024

- proliferation of advanced AI models
  - GPT-4o, Claude Sonnet, Llama 3, Sora
  - *transforming industries* such as content creation, customer service, education, etc.
- breakthroughs in specialized AI applications
  - Figure 02, Optimus, AlphaFold 3
  - driving unprecedented advancements in automation, drug discovery, scientific understanding - *profoundly affecting healthcare, manufacturing, scientific research*



## **Transformative impact of AI - reshaping industries, work & society**

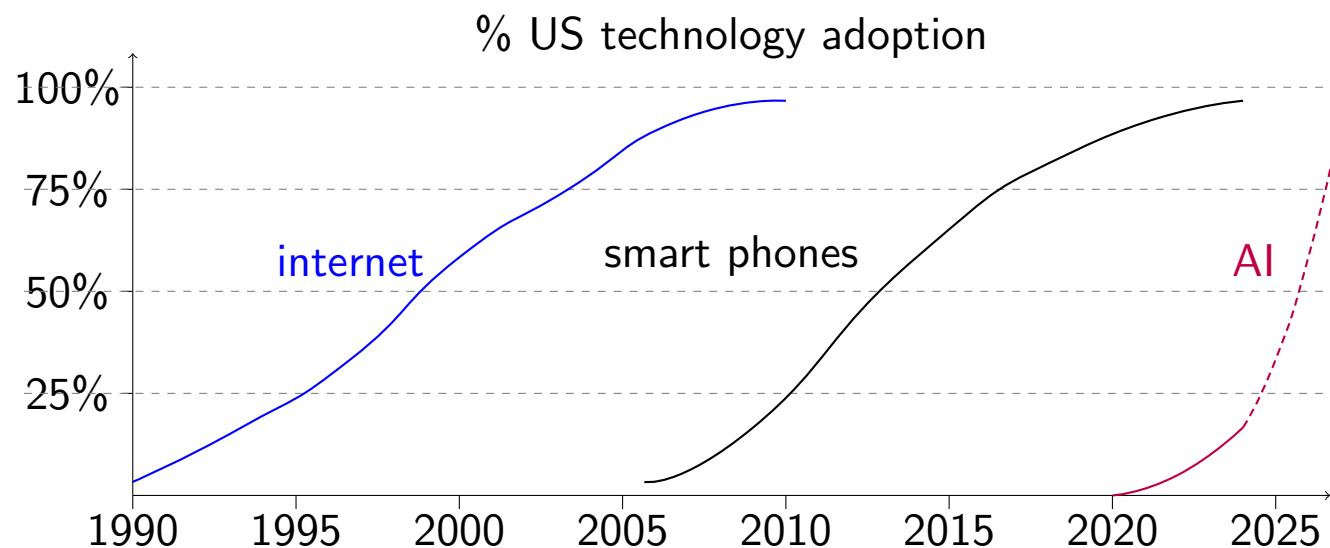
- accelerating human-AI collaboration
    - not only reshaping industries but *altering how humans interact with technology*
    - AI's role as collaborator and augmentor redefines productivity, creativity, the way we address global challenges, e.g., *sustainability & healthcare*
  - AI-driven automation *transforms workforce dynamics* - creating new opportunities while challenging traditional job roles
  - *ethical AI considerations* becoming central not only to business strategy, but to society as a whole - *influencing regulations, corporate responsibility & public trust*



# **Recent Advances in AI**

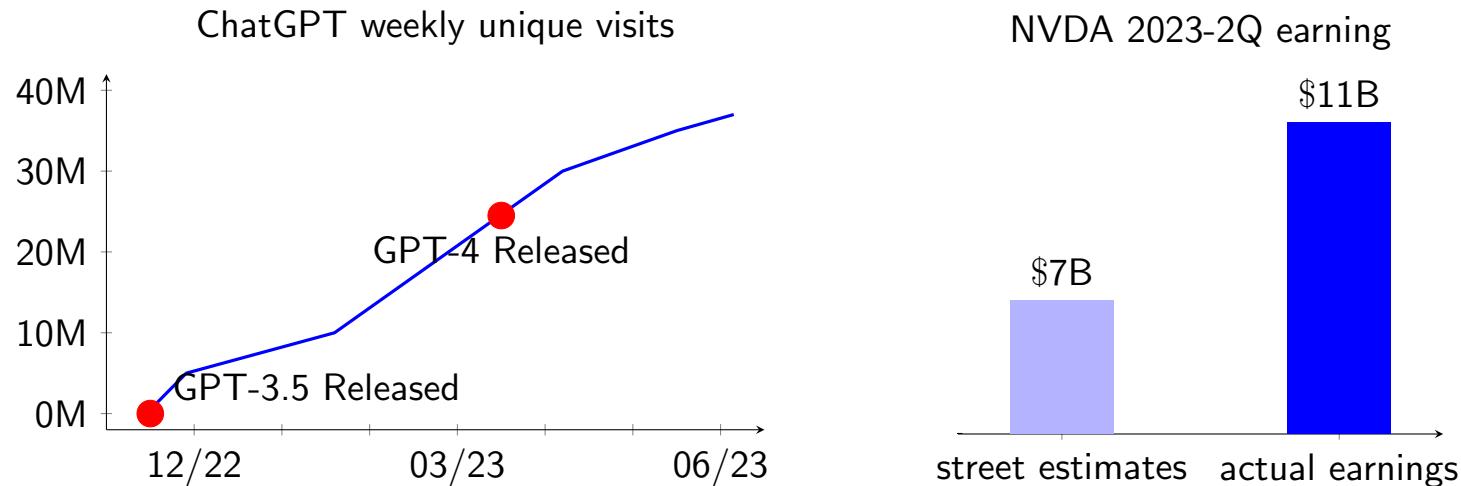
## Where are we in AI today?

- sunrise phase - currently experiencing dawn of AI era with significant advancements and increasing adoption across various industries
- early adoption - in early stages of AI lifecycle with widespread adoption and innovation across sectors marking significant shift in technology's role in society



## Explosion of AI ecosystems - ChatGPT & NVIDIA

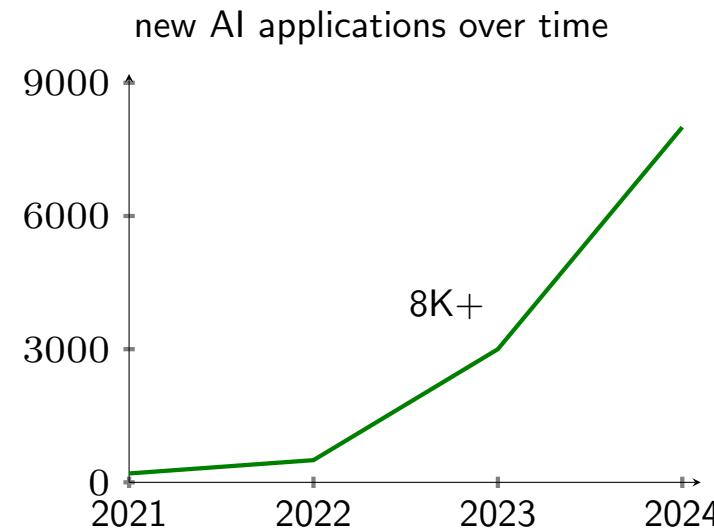
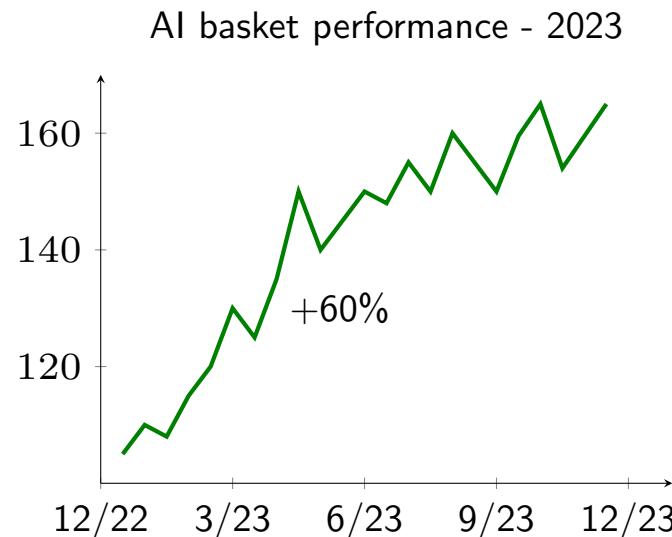
- took only *5 months for ChatGPT users to reach 35M*
- NVIDIA 2023 Q2 earning exceeds market expectation by big margin - \$7B vs \$13.5B
  - surprisingly, *101% year-to-year growth*
  - even more surprisingly *gross margin was 71.2%* - up from 43.5% in previous year<sup>3</sup>



<sup>3</sup>source - Bloomberg

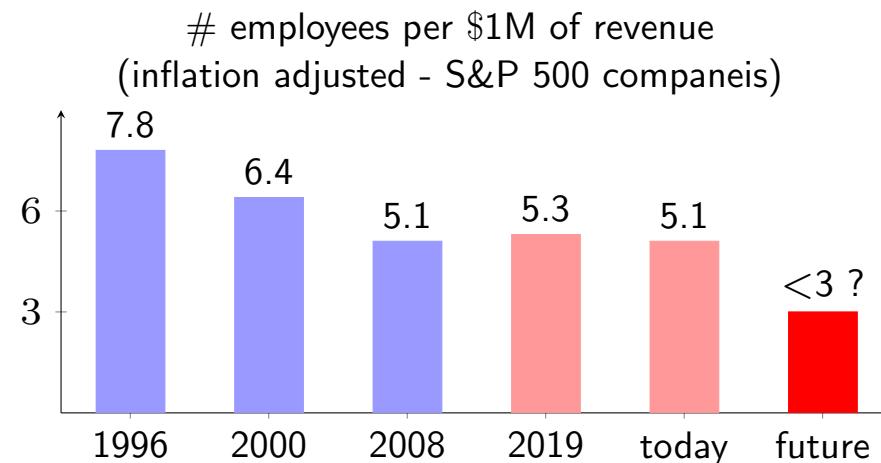
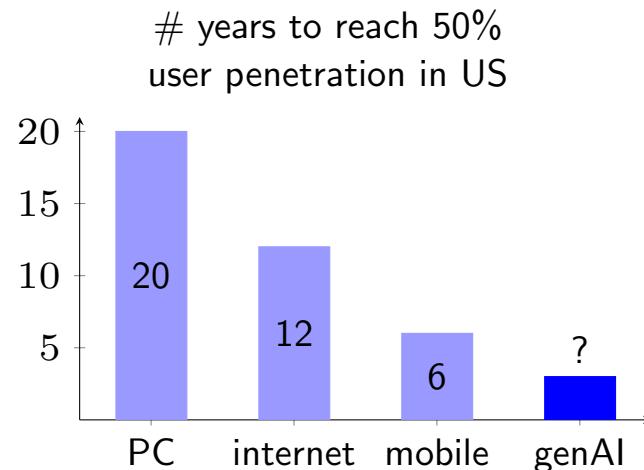
## Explosion of AI ecosystems - AI stock market

- *AI investment surge in 2023 - portfolio performance soars by 60%*
  - AI-focused stocks significantly outpaced traditional market indices
- *over 8,000 new AI applications* developed in last 3 years
  - applications span from healthcare and finance to manufacturing and entertainment



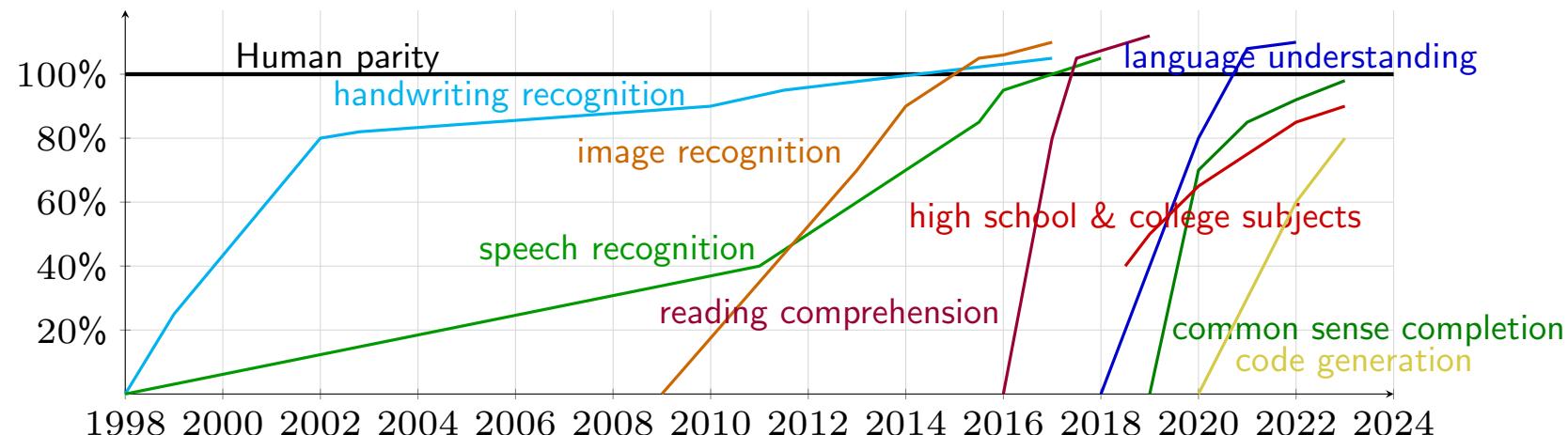
## AI's transformative impact - adoption speed & economic potential

- adoption - has been twice as fast with platform shifts suggesting
  - increasing demand and readiness for new technology improved user experience & accessibility
- AI's potential to drive economy for years to come
  - 35% improvement in productivity driven by introduction of PCs and internet
  - greater gains expected with AI proliferation



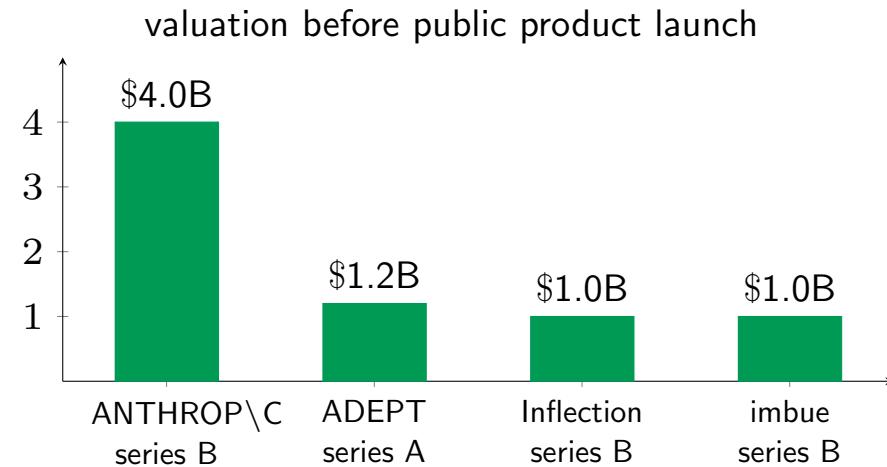
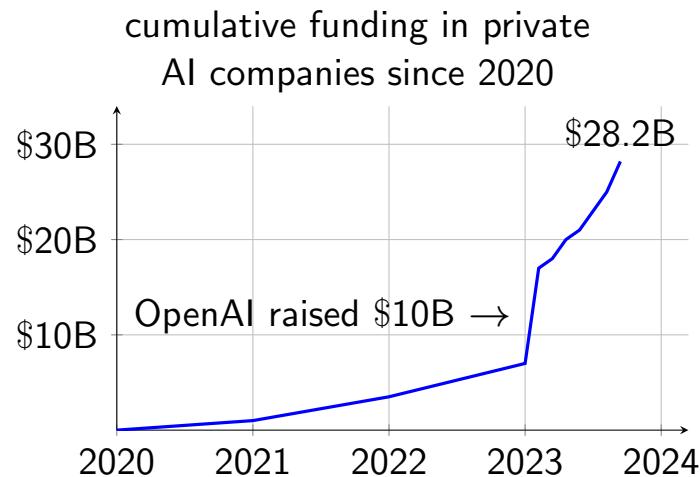
## AI getting more & more faster

- steep upward slopes of AI capabilities highlight accelerating pace of AI development
  - period of exponential growth with AI potentially mastering new skills and surpassing human capabilities at ever-increasing rate
- closing gap to human parity - some capabilities approaching or arguably reached human parity, while others having still way to go
  - achieving truly human-like capabilities in broad range remains a challenge



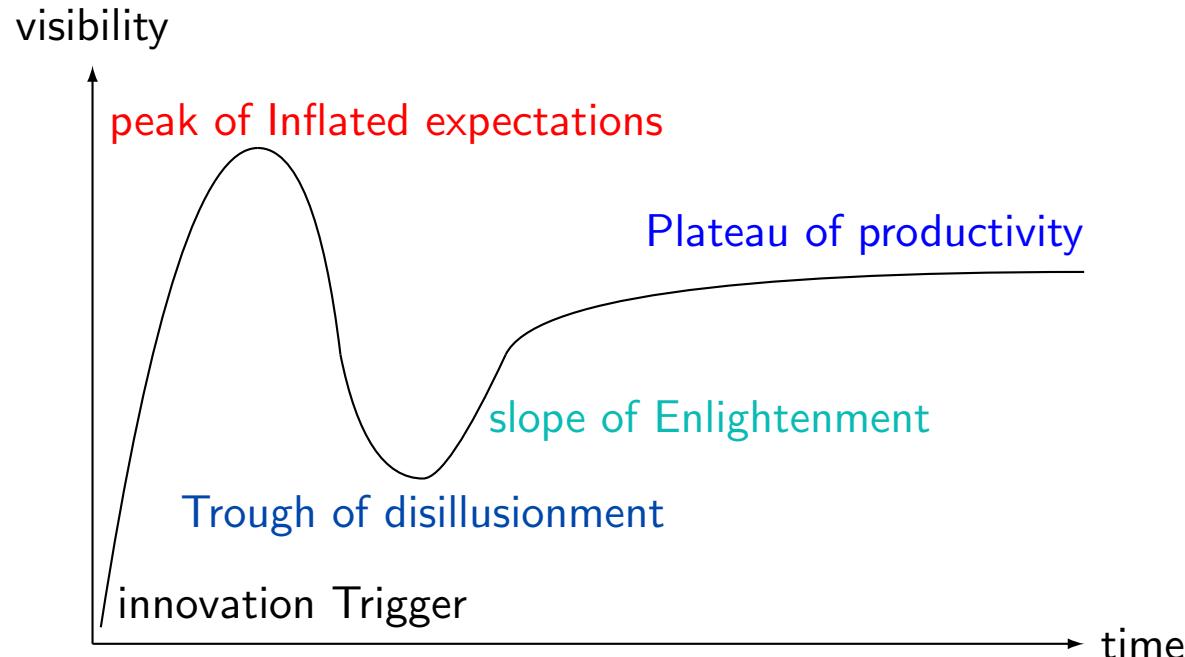
## Massive investment in AI

- *explosive growth* - cumulative funding skyrocketed reaching staggering \$28.2B
- OpenAI - significant fundraising (= \$10B) fueled rapid growth
- *valuation surge* - substantial valuations even before public products for stellar companies
- *fierce competition for capital* among AI startups driving innovation & accelerating development
- massive investment indicates *strong belief in & optimistic outlook for potential of AI* to revolutionize industries & drive economic growth



**Is AI hype?**

## Technology hype cycle



- innovation trigger - technology breakthrough kicks things off
- peak of inflated expectations - early publicity induces many successes followed by even more
- trough of disillusionment - expectations wane as technology producers shake out or fail
- slope of enlightenment - benefit enterprise, technology better understood, more enterprises fund pilots

## Fiber vs cloud infrastructure

- fiber infrastructure - 1990s
  - Telco Co's raised \$1.6T of equity & \$600B of debt
  - bandwidth costs decreased 90% within 4 years
  - companies - Covage, NothStart, Telligent, Electric Lightwave, 360 networks, Nextlink, Broadwind, UUNET, NFS Communications, Global Crossing, Level 3 Communications
  - became *public good*
- cloud infrastructure - 2010s
  - entirely new computing paradigm
  - mostly public companies with data centers
  - *big 4 hyperscalers generate \$150B + annual revenue*



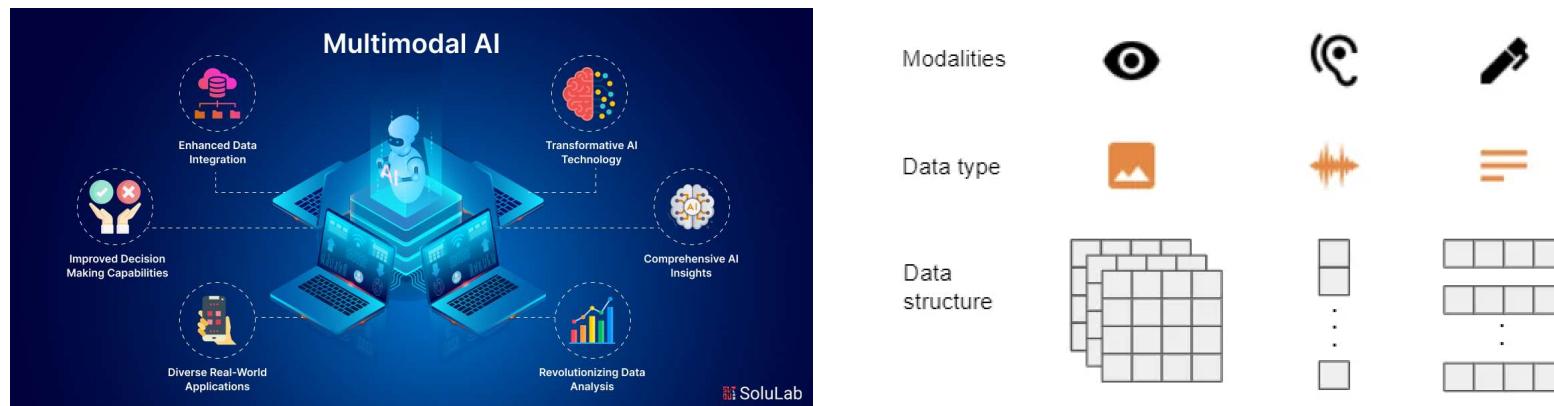
## Yes & No

characteristics of hype cycles	speaker's views
value accrual misaligned with investment	<ul style="list-style-type: none"><li>• OpenAI still operating at a loss; business model <i>still</i> not clear</li><li>• gradual value creation across broad range of industries and technologies (<i>e.g.</i>, CV, LLMs, RL) unlike fiber optic bubble in 1990s</li></ul>
overestimating timeline & capabilities of technology	<ul style="list-style-type: none"><li>• self-driving cars delayed for over 15 years, with limited hope for achieving level 5 autonomy</li><li>• AI, however, has proven useful within a shorter 5-year span, with enterprises eagerly adopting</li></ul>
lack of widespread utility due to technology maturity	<ul style="list-style-type: none"><li>• AI already providing significant utility across various domains</li><li>• vs quantum computing remains promising in theory but lacks widespread practical utility</li></ul>

# AI Agents

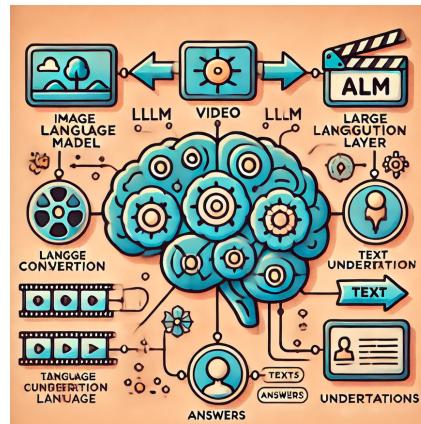
# Multimodal learning

- understand information from multiple modalities, *e.g.*, text, images, audio, video
- representation learning methods
  - combine multiple representations or learn multimodal representations simultaneously
- applications
  - images from text prompt, videos with narration, musics with lyrics
- collaboration among different modalities
  - understand image world (open system) using language (closed system)



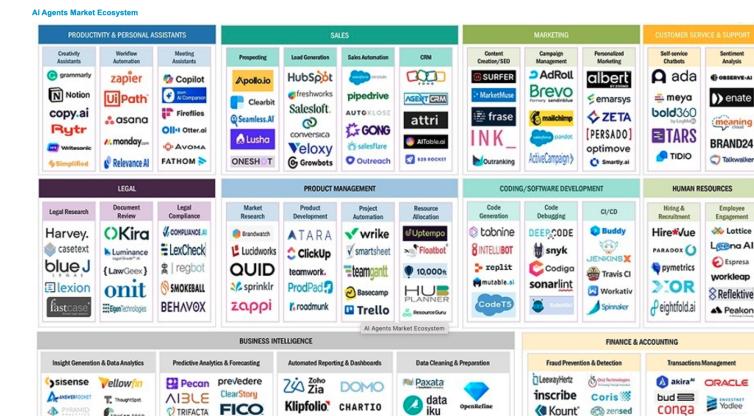
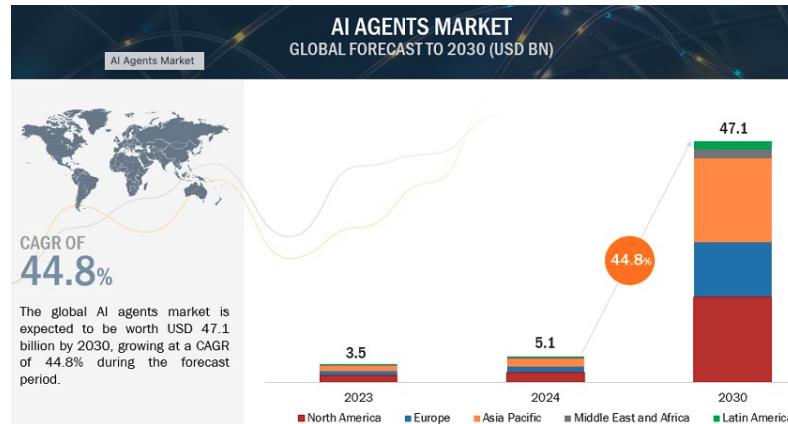
## Implications of success of LLMs

- many researchers change gears towards LLM
  - from computer vision (CV), speech, music, video, even reinforcement learning
- *LLM is not only about NLP . . . humans have . . .*
  - evolved to optimize natural language structures for eons
  - handed down knowledge using *this natural languages* for thousands of years
  - internal structure (or equivalently, representation) of natural languages optimized via *thousands of generation by evolution*
- *LLM connects non-linguistic world (open system) via natural languages (closed system)*



## Multimodal AI (mmAI) - definition & history

- mmAI - systems processing & integrating data from multiple sources & modalities, to generate unified response / decision
- 1990s – 2000s - early systems - initial research combining basic text & image data
- 2010s - CNNs & RNNs enabling more sophisticated handling of multimodality
- 2020s - modern multimodal models - Transformer-based architectures handling complex multi-source data at highly advanced level
- mmAI *mimics human cognitive ability* to interpret and integrate information from various sources, leading to holistic decision-making

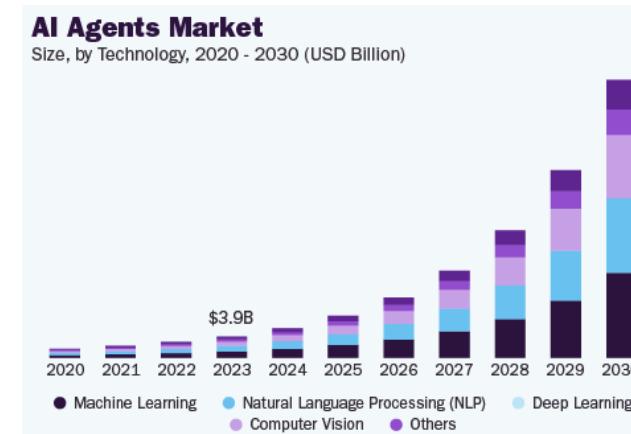
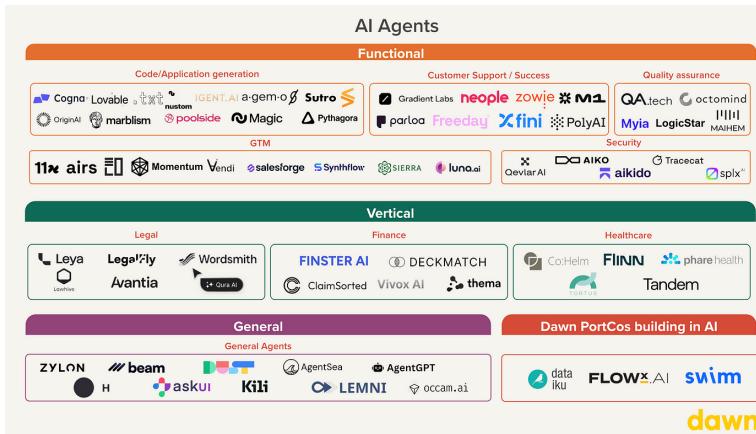


## mmAI Technology

- core components
  - data preprocessing - images, text, audio & video
  - architectures - unified Transformer-based (*e.g.*, ViT) & cross-attention mechanisms / hybrid architectures (*e.g.*, CNNs + LLMs)
  - integration layers - fusion methods for combining data representations from different modalities
- technical challenges
  - data alignment - accurate alignment of multimodal data
  - computational demand - high-resource requirements for training and inferencing
  - diverse data quality - manage variations in data quality across modalities
- advancements
  - multimodal embeddings - shared feature spaces interaction between modalities
  - self-supervised learning - leverage unlabeled data to learn representations across modalities

# AI agents powered by multimodal LLMs

- foundation
  - integrate multimodal AI capabilities for enhanced interaction & decision-making
- components
  - perceive environment through multiple modalities (visual, audio, text), process using LLM technology, generate contextual responses & take actions
- capabilities
  - understand complex environments, reason across modalities, engage in natural interactions, adapt behavior based on context & feedback



## AI agents - Present & Future

- emerging applications
  - scientific research - agents analyzing & running experiments & generating hypotheses
  - creative collaboration - AI partners in design & art combining multiple mediums
  - environmental monitoring - processing satellite sensor data for climate analysis
  - healthcare - enhanced diagnostic combining imaging, *e.g.*, MRI, with patient history
  - customer experience - virtual assistants understanding spoken language & visual cues
  - autonomous vehicles - integration of visual, radar & audio data
- future
  - ubiquitous AI agents - seamless integration into everyday devices
  - highly tailored personalized experience - in education, entertainment & healthcare

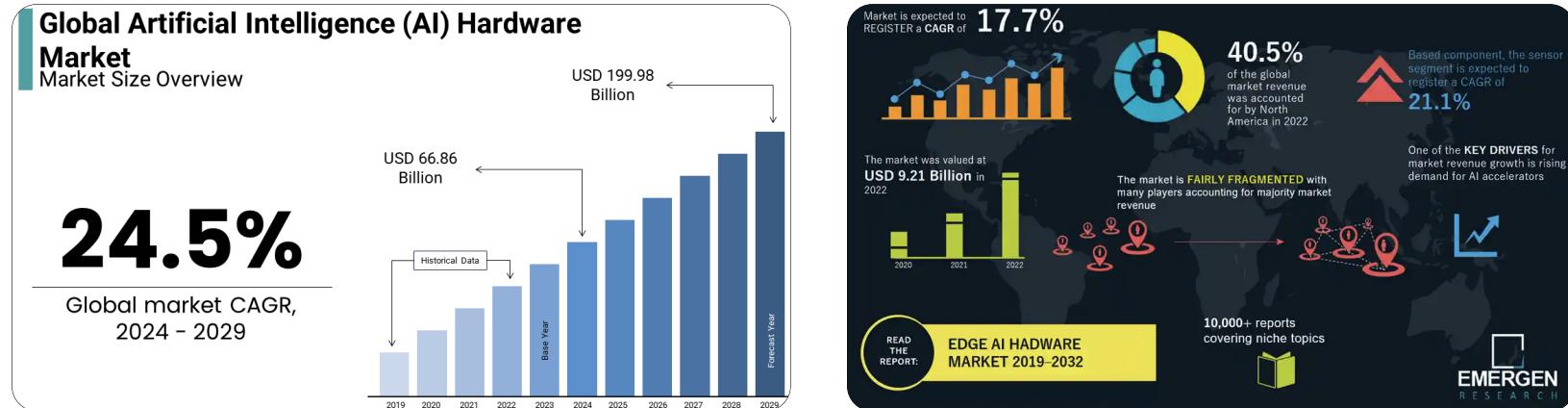


# **AI Hardware**

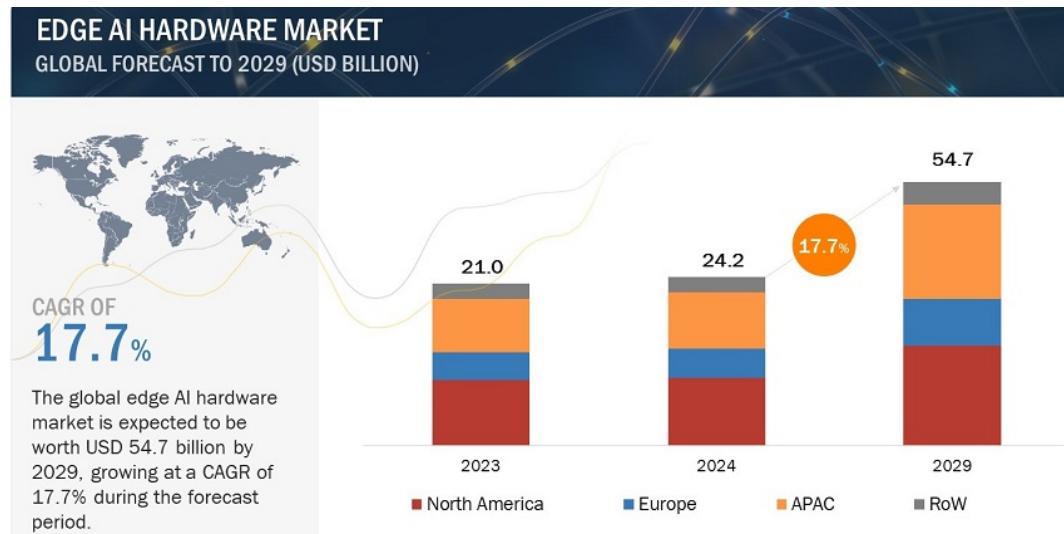
# **AI Hardware Industry**

## Landscape of AI hardware industry

- global AI hardware market valued at \$66.96B in 2024, projected to grow significantly
- major companies - Nvidia, Intel, AMD, Qualcomm, and IBM w/ Nvidia holding substantial market share



- North America leading market - high R&D investments & key industry players
- Asia Pacific rapidly expanding - strong semiconductor industries in South Korea, China & Japan
- demand for advanced processors such as GPUs, TPUs & AI accelerators rising due to complexity of AI algorithms & high computational power



## Predictions for future of AI hardware market

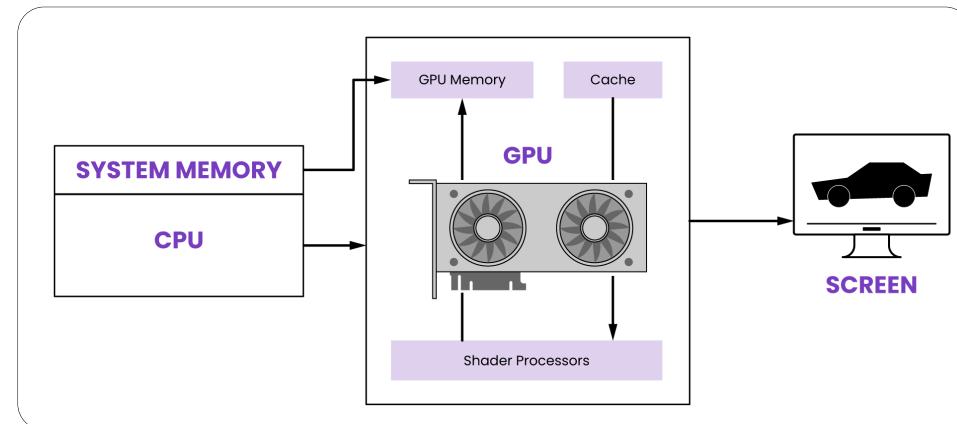
- AI hardware market expected to reach \$382B by 2032 - significant growth in data center AI chips
- integration of AI w/ 5G & increased use of AI in edge computing anticipated to drive future demand
- AI hardware becoming crucial in sectors such as autonomous vehicles, robotics & medical devices
- need to address challenges such as heat and power management along with technical complexities



# **GPUs and AI Accelerators**

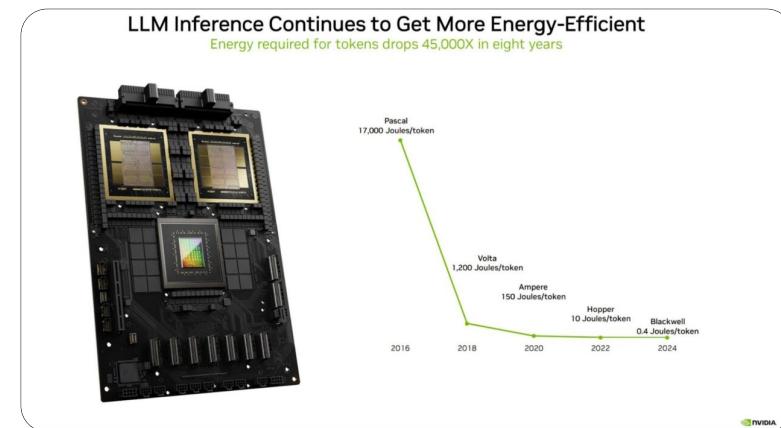
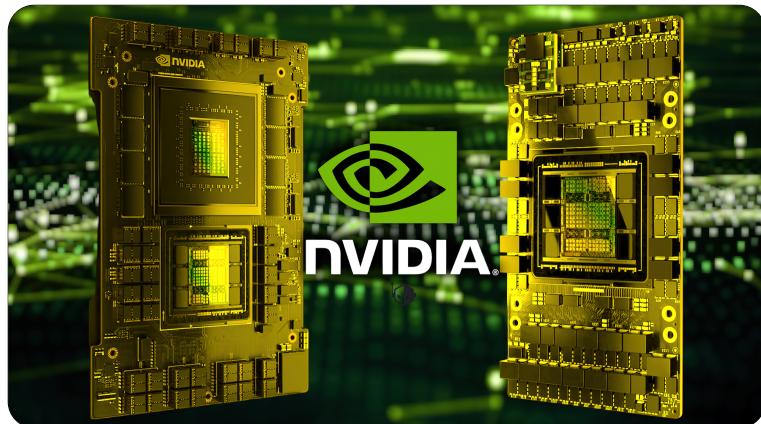
## Technical challenges of GPUs & AI accelerators

- facing challenges in scaling to handle increasingly large AI models and datasets - traditional architectures struggling w/ massive parallel processing demands of modern AI applications
- AI applications require extensive memory bandwidth often leading to bottlenecks - efficient memory management is crucial
- AI accelerators consume significant power - high operational costs and environmental concerns for both cloud-based & edge AI applications



## Potential solutions for overcoming challenges

- development of AI-specific architectures such as tensor cores and custom ASICs to improve efficiency and performance - novel architectures like FPGAs for specific AI tasks, *e.g.*, for RAG & vectorDB
- implementing software optimizations to enhance hardware usability and performance - use of compilers and frameworks that maximize efficiency of existing hardware
- encouraging market competition to drive innovation and reduce monopolistic control - exploring alternative hardware solutions and improving energy efficiency standards



## Big tech's in-house chip development

- shift towards in-house AI hardware - major tech companies increasingly developing their own AI chips - move to enhance AI capabilities and reduce dependence
- collaboration with specialized partners - partnering with specialized firms for manufacturing and technology blending in-house expertise with external innovation

	Microsoft	Google	Amazon	Meta
Chip	Maia 100	TPU v5e	Inferentia2	MTIA v1
Launch Date	November, 2023	August, 2023	Early 2023	2025
IP	ARM	ARM	ARM	RISC-V
Process Technology	TSMC 5nm	TSMC 5nm	TSMC 7nm	TSMC 7nm
Transistor Count	105 billion	-	-	-
INT8	-	393 TOPS	-	102.4 TOPS
FP16	-	-	-	51.2 TFLOPS
BF16	-	197 TFLOPS	-	-
Memory	-	-	-	LPDDR5
TDP	-	-	-	25W
Packaging Technology	CoWoS	CoWoS	CoWoS-S	2D
Collaborating Partners	Global Unichip Corp.	Broadcom	Alchip Technologies	Andes Technology
Application	Training/Inference	Inference	Inference	Training/Inference
LLM	GPT-3.5, GPT-4	BERT, PaLM, LaMDA	Titan FM	Llama, Llama2

## AMD - Nvidia's new competitor

- key points
  - AMD launched new AI accelerator chip, *Instinct MI300X*, on Dec 6, 2023
  - CDNA 3 architecture, mix of 5nm and 6nm IPs, delivering 153B transistors
  - *outperforms Nvidia's H100 TensorRT-LLM* by 1.6X higher memory bandwidth and 1.3X FP16 TFLOPS
  - up to 40% faster vs Nvidia's Llama-2 70B model in 8x8 server configurations
- market impact
  - significant challenge to Nvidia's dominance in AI accelerator market
  - performance gains over Nvidia's offerings could drive *customer adoption and market share for AMD*
- future prediction
  - *AMD stocks soared* since launch indicating investor confidence in their competitiveness
  - Lisa Su, AMD's CEO, categorized Instinct MI300X as “next big thing” in tech industry
  - potential risks include need to *manage ROCm vs CUDA software ecosystem* & ensure rapid customer adoption and production coverage

# **AI Accelerator Startups**

## AI accelerator startups

- innovative architectures - startups like Groq, SambaNova & Graphcore leading with *novel architectures designed to accelerate AI workloads*
  - *Groq* - tensor streaming processor (TSP) offering ultra-low latency & high throughput, high-performance AI inference chips enhancing speed & efficiency
  - *SambaNova* - reconfigurable dataflow architecture optimizing for various AI workloads
  - *Graphcore* - intelligence processing unit (IPU) tailored for graph-based computation excelling in sparse data processing
  - *Cerebras Systems* - develop wafer scale engine (WSE), largest chip built for AI workloads, unmatched computational power revolutionizing AI hardware capabilities
  - *Hailo* - specialize for edge devices optimizing AI processes for real-time applications, raised \$120M emphasizing potential to disrupt traditional AI chip markets

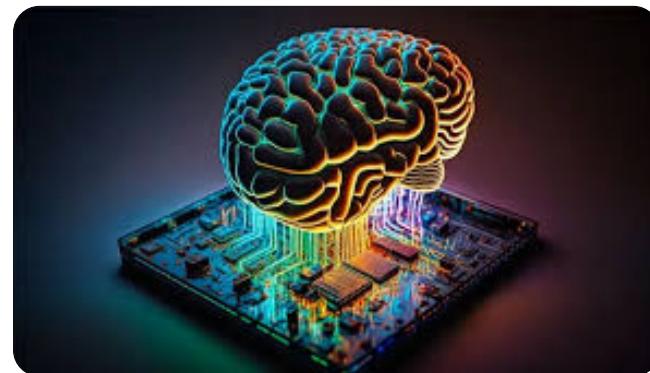


## Technological competitiveness

- energy efficiency
  - energy-efficient designs crucial for scalability in data centers and edge devices
  - startups developing solutions significantly reducing power consumption without compromising performance
- customization & flexibility
  - AI accelerators from startups often offer greater customization options for specific AI tasks compared to traditional GPUs
  - flexibility in hardware allows for tailored solutions that can outperform general-purpose accelerators in certain applications
- software integration
  - robust software ecosystems critical - startups investing in developing software stacks that optimize performance for their hardware
  - compatibility with existing AI frameworks is competitive advantage, *e.g.*, TensorFlow & PyTorch

## Industry and market influence

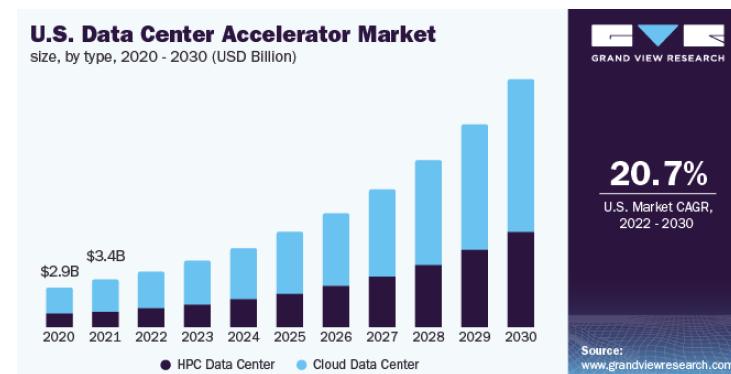
- disruption of traditional players
  - challenging dominance of established players like NVIDIA & Intel
  - unique architectures providing specialized solutions traditional GPUs and CPUs cannot efficiently handle
- driving down costs
  - offering competitive alternatives pushing down cost of AI computation
  - could lead to democratization of AI w/ more companies affording high-performance AI capabilities



- accelerating AI innovation
  - contributing to rapid innovation providing hardware that can handle emerging AI models & workloads
  - adaptability and specialization enable advancements in AI research & faster development cycles
- strategic partnerships & acquisitions
  - big techs increasingly forming strategic partnerships or acquiring startups to stay competitive
  - collaborations can speed up integration of advanced AI hardware into mainstream products



- market growth & opportunities
  - AI accelerator market expected to grow significantly driven by demand in data centers, edge computing & autonomous systems
  - startups well-positioned to capture significant share of growing market particularly in niche applications
- future outlook
  - dependency on Asia for fabrication might lead to strategic shifts in global tech policies and investments in local manufacturing
  - increasing demand for efficient AI processing on edge devices and in data center.



# **Global Semiconductor Industry**

## Hard-to-predict AI hardware markets

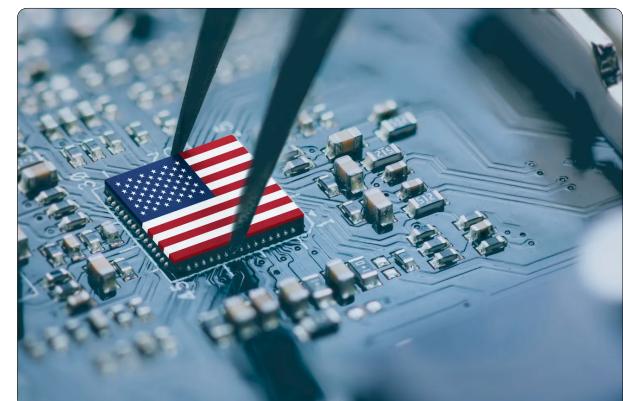
- US
  - birthplace for modern semiconductor chips driving PC market, internet, multi-media, mobile phones, and AI . . .
    - Intel, Texas Instrument (TI), Global Foundry
  - traditionally strong with design houses - NVIDIA, AMD, Broadcom, Apple, . . .
  - threatened experiencing global chip shortage & vulnerable supply chain via COVID
  - national security concerns & economic competitiveness
- China
  - strong fast followers - SMIC<sup>4</sup>, Huawei, Hua Hong Semiconductor (foundry)
- South Korea
  - best memory chip makers - Samsung, SK hynix
  - struggling with LSI and foundry business

---

<sup>4</sup>SMIC - Semiconductor Manufacturing International Corporation

## Reshoring semiconductor manufacturing industry

- trade & semiconductor WAR between US & China
  - export controls on advanced chips and equipment
- CHIPS & Science Act (Aug, 2022)
  - \$52B in subsidies for domestic production, 25% investment tax credit for chip plants
  - (coerce) world-best semiconductor manufacturers build factories in US with support
    - GlobalFoundries - \$1.5B @ Feb-2024
    - Intel - \$8.5B @ Apr-2024 - Ohio - two fabs expandable to \$100B
    - Samsung - \$6.4B @ Apr-2024 - Talor, Texas
    - TSMC - \$6.6B @ Apr-2024 - Phoenix, Arizona
      - two foundry fabs (3nm & 4nm)



## Turmoils in global semiconductor business

- global context
  - EU Chips Act - €43B to boost European chip production
  - Japan & South Korea - significant investments in domestic capacity
- industry dynamics
  - Intel's foundry ambitions - targeting 50% global market share by 2030
  - TSMC expanding global footprint (US, Japan, possibly Germany)
- future outlook
  - projected shift in global semiconductor manufacturing landscape
  - increased geographical diversification of chip production

## Export controls on US chip technology to China



- goal - limit China's access to advanced semiconductor tech to maintain US strategic advantage
- impacts on
  - China - advanced chips and equipment not allowed, domestic innovation increased
  - US - short-term - US lose market share and revenue in China
  - US - long-term - potential decline in US global competitiveness
- Chinese response - circumvent controls and adapt supply chains
- conclusion
  - US-China chip rivalry transforms global supply chains with deep implications for *security & industry*
  - US success hinges on better coordination and policy analysis
- reference - [Balancing the Ledger - Center for Strategic & International Studies \(CSIS\)](#)

## China strikes back on US sanction

- **Huawei's launch of Mate 60 Pro smartphone**
  - these domestically produced chips represent major breakthrough against US sanctions
  - its success with *advanced 7nm Kirin 9000S chip* demonstrates significant progress in China's self-reliance in high-tech manufacturing - narrowing the technological gap with global leaders
- **Huawei case highlights potential failure of US sanctions potentially leading to more aggressive US measures**
  - US export controls on China's semiconductor industry are effective in the short term but insufficient to halt China's progress especially in legacy chip manufacturing
  - to maintain technological edge, US must balance further restrictions with supporting its semiconductor industry to avoid overreliance on export controls



## Chinese semiconductor companies

- Chinese major semiconductor companies
  - SMIC - China's largest chip foundry, advancing 7nm technology
  - HiSilicon - Huawei's chip design arm, crucial for the Kirin processors
  - YMTC - leader in 3D NAND memory chip production
  - Huahong Group, CXMT, SMEE, GigaDevice, UnilC Semiconductors, ASMC, etc.
- *SMIC shows significant progress in producing 7nm chips* & YMTC leads memory chip manufacturer - both face challenges from US export controls
- industry faces internal challenges, e.g., corruption & misallocation of resources
- but remains crucial to China's goal of technological self-reliance



# Appendices

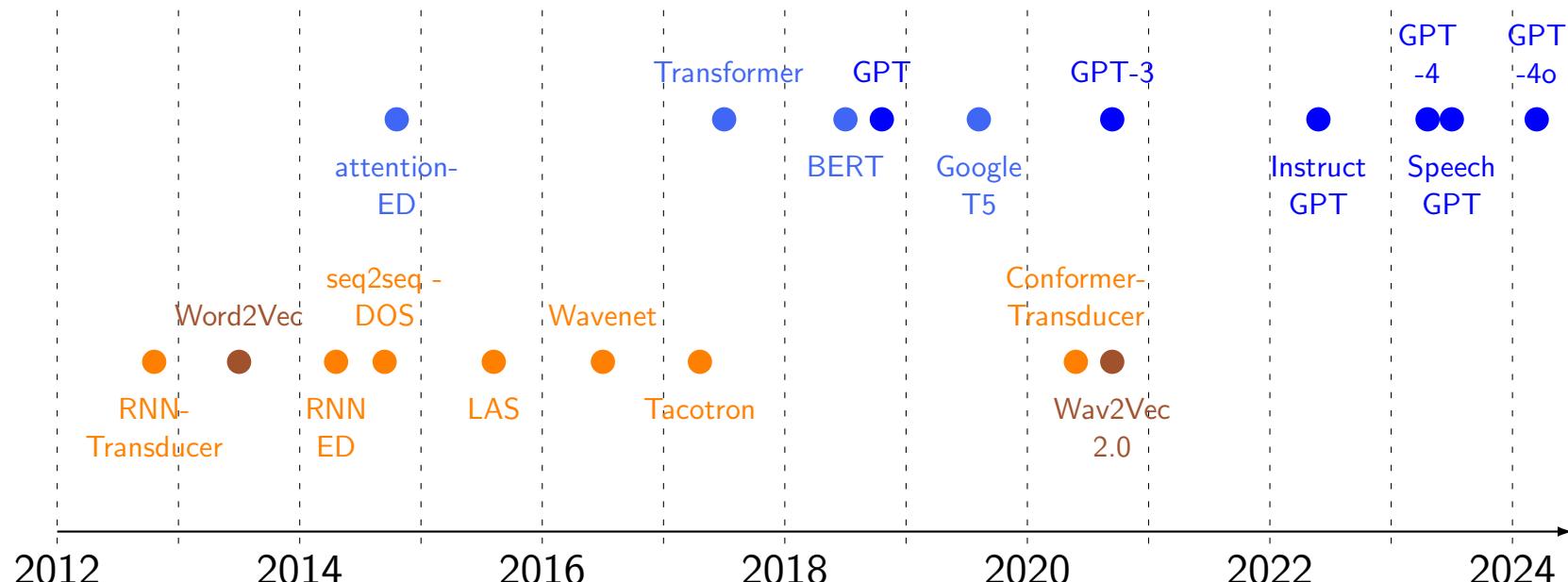
**LLM**

# **Language Models**

## History of language models

- bag of words - first introduced – 1954
- word embedding – 1980
- RNN based models - conceptualized by David Rumelhart – 1986
- LSTM (based on RNN) – 1997
- 380M-sized seq2seq model using LSTMs proposed – 2014
- 130M-sized seq2seq model using gated recurrent units (GRUs) – 2014
- Transformer - Attention is All You Need - A. Vaswani et al. @ Google – 2017
  - 100M-sized encoder-decoder multi-head attention model for machine translation
  - non-recurrent architecture, handle arbitrarily long dependencies
  - parallelizable, *simple* (linear-mapping-based) attention model

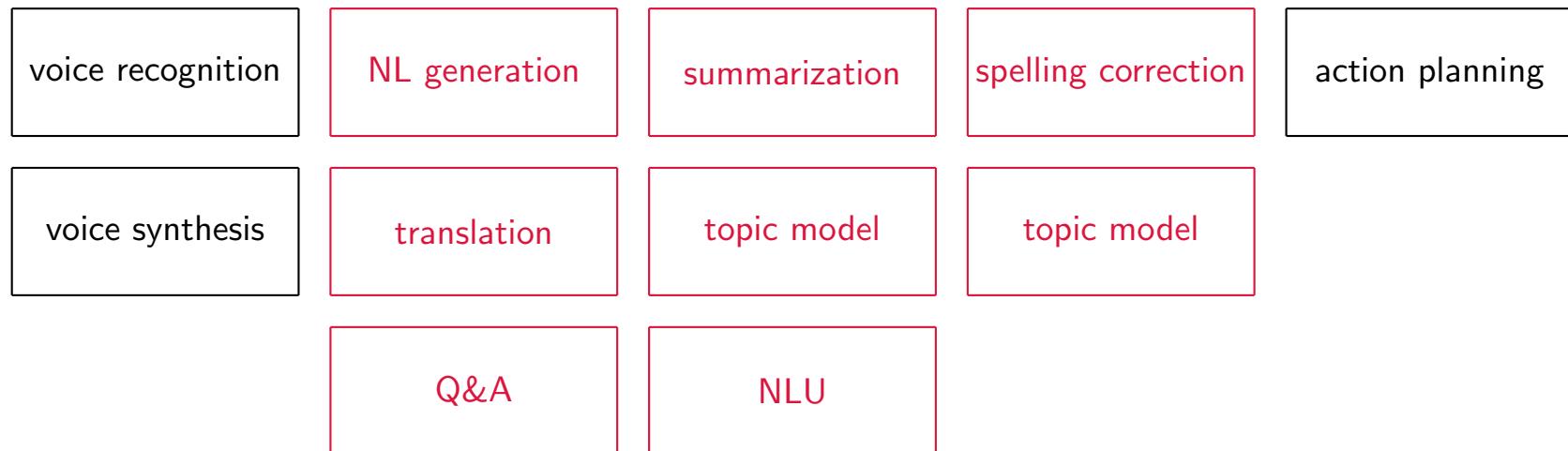
## Recent advances in speech & language processing



- LAS: listen, attend, and spell, ED: encoder-decoder, DOS: decoder-only structure

## Types of language models

- many of language models have **common requirements** - language representation learning
- can be learned via pre-training *high performing model* and fine-tuning/transfer learning/domain adaptation
- this *high performing model* learning essential language representation *is* (language) foundation model
  - actually, same for other types of learning, e.g., CV



**NLP Market**

## NLP market size

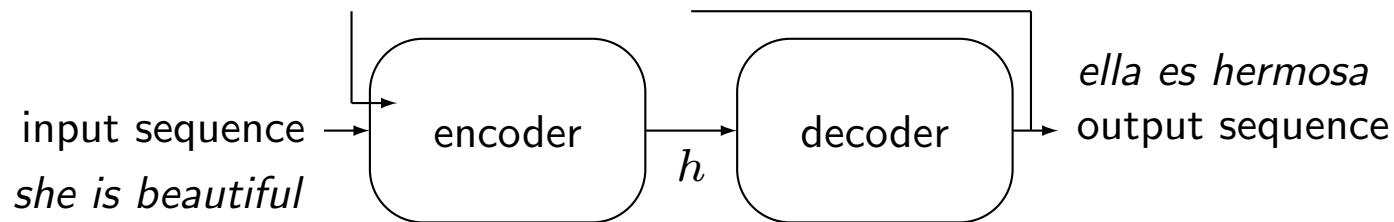
- global NLP market size estimated at USD 16.08B in 2022, is expected to hit USD 413.11B by 2032 - *CAGR of 38.4%*
- in 2022
  - north america NLP market size valued at USD 8.2B
  - high tech and telecom segment accounted revenue share of over 23.1%
  - healthcare segment held a 10% market share
  - (by component) solution segment hit 76% revenue share
  - (deployment mode) on-premise segment generated 56% revenue share
  - (organizational size) large-scale segment contributed highest market share
- source - [Precedence Research](#)



# **Sequence-to-Sequence Models**

## Sequence-to-sequence (seq2seq) model

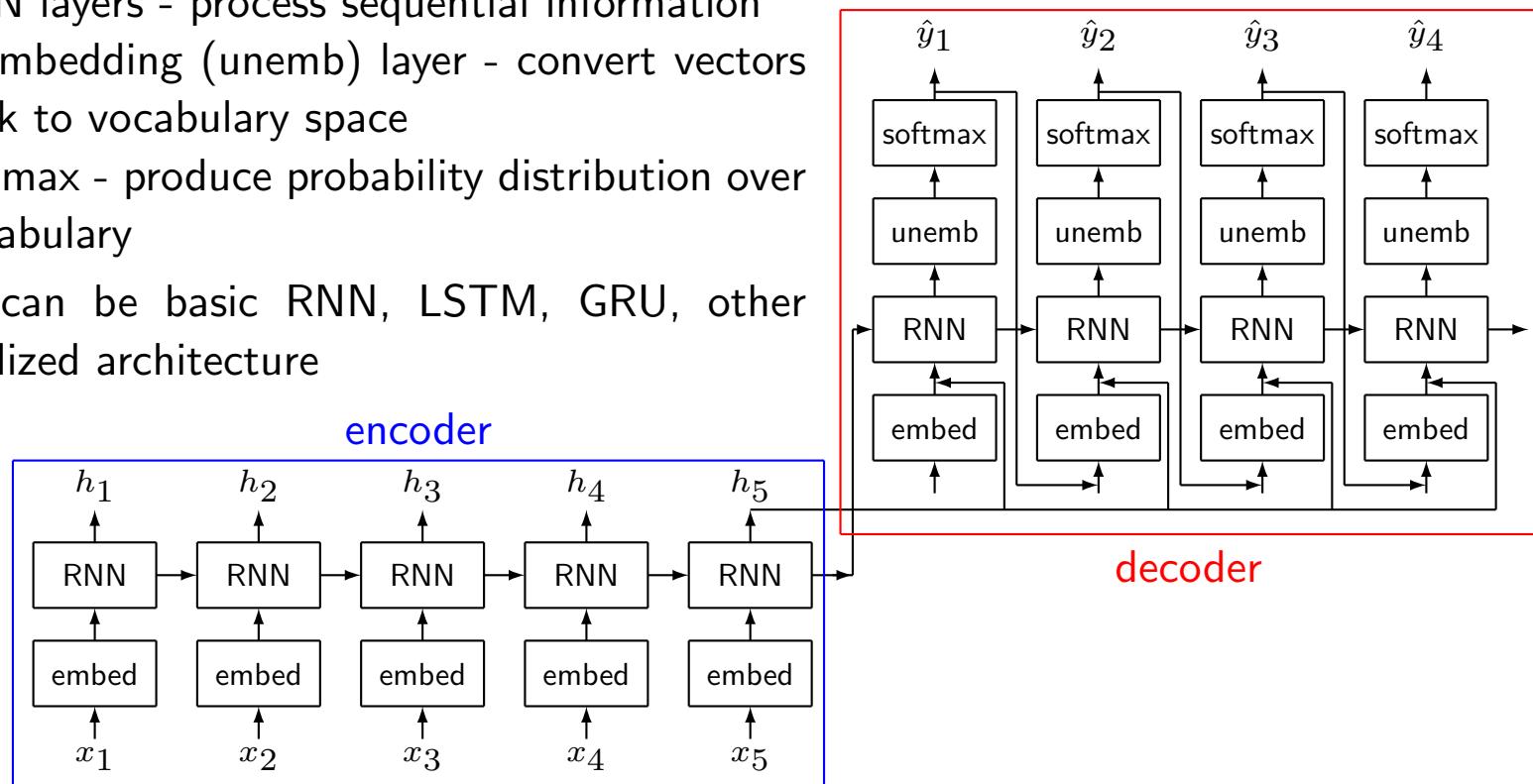
- seq2seq - take sequences as inputs and spit out sequences
- encoder-decoder architecture



- encoder & decoder can be RNN-type models
- $h \in \mathbf{R}^n$  - hidden state - *fixed length* vector
- (try to) condense and store information of input sequence (losslessly) in (fixed-length) hidden states
  - finite hidden state - not flexible enough, *i.e.*, cannot handle arbitrarily large information
  - memory loss for long sequences
  - LSTM was promising fix, but with (inevitable) limits

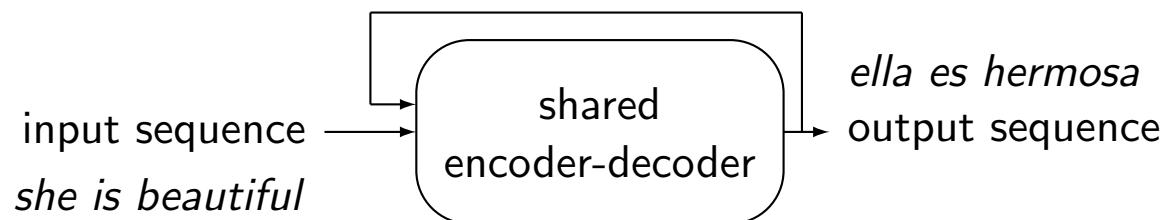
## RNN-type encoder-decoder architecture

- components
  - embedding layer - convert input tokens to vector representations
  - RNN layers - process sequential information
  - unembedding (unemb) layer - convert vectors back to vocabulary space
  - softmax - produce probability distribution over vocabulary
- RNN can be basic RNN, LSTM, GRU, other specialized architecture



## Shared encoder-decoder model

- single neural network structure can handle both encoding & decoding tasks
  - efficient architecture reducing model complexity
  - allow for better parameter sharing across tasks
- widely used in modern LLMs to process & generate text sequences
  - applications - machine translation, text summarization, question answering
- advantages
  - efficient use of parameters, versatile for multiple NLP tasks



# **Large Language Models**

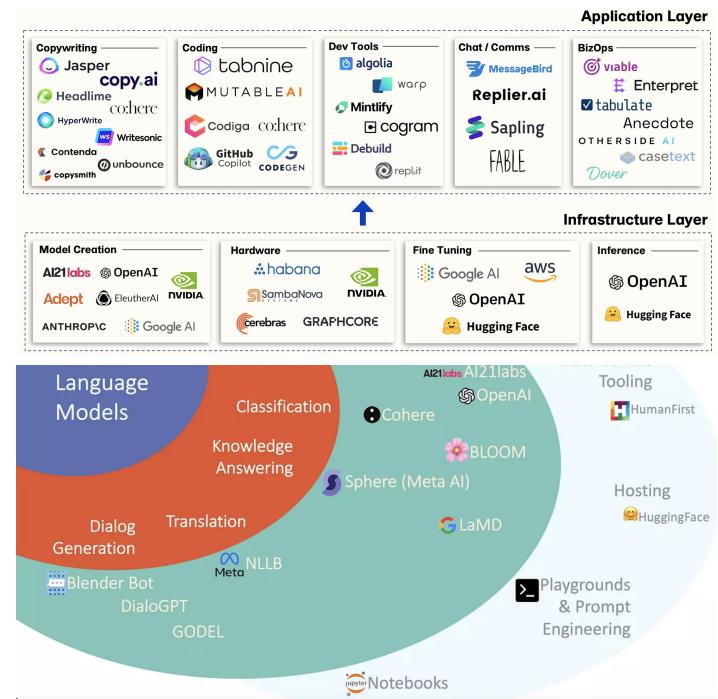
# LLM

- LLM
  - type of AI aimed for NLP trained on massive corpus of texts & programming code
  - allow learn statistical relationships between words & phrases, *i.e.*, conditional probabilities
  - *amazing performance shocked everyone - unreasonable effectiveness of data (Halevy et al., 2009)*
- applications
  - conversational AI agent / virtual assistant
  - machine translation / text summarization / content creation / sentiment analysis / question answering
  - code generation
  - market research / legal service / insurance policy / triage hiring candidates
  - + virtually infinite # of applications



# LLMs

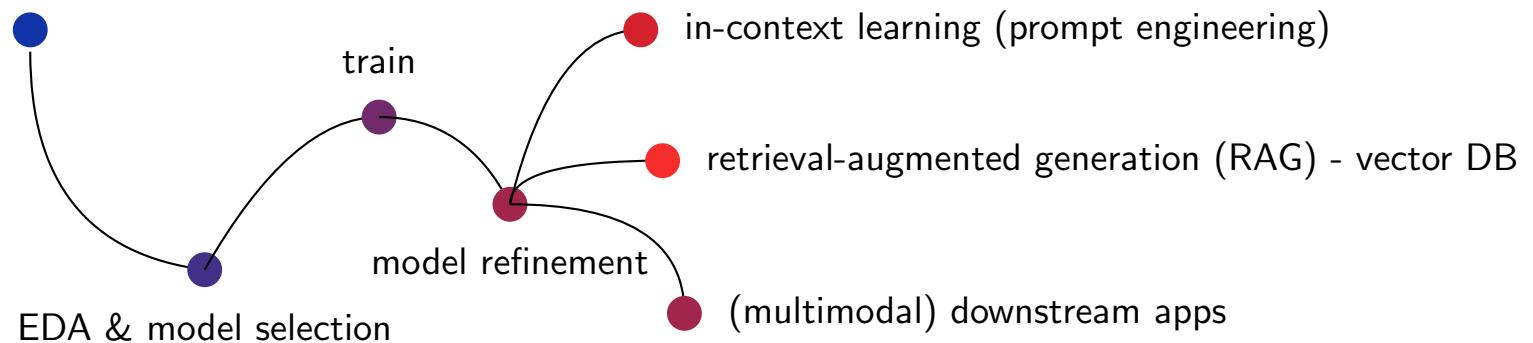
- Foundation Models
  - GPT-x/Chat-GPT - OpenAI, Llama-x - Meta, PaLM-x (Bard) - Google
- # parameters
  - generative pre-trained transformer (GPT) - GPT-1: 117M, GPT-2: 1.5B, GPT-3: 175B, GPT-4: 100T, GPT-4o: 200B
  - large language model Meta AI (Llama) - Llama1: 65B, Llama2: 70B, Llama3: 70B
  - scaling language modeling with pathways (PaLM) - 540B
- burns lots of cash on GPUs!
- applicable to many NLP & genAI applications



## LLM building blocks

- data - trained on massive datasets of text & code
  - quality & size critical on performance
- architecture - GPT/Llama/Mistral
  - can make huge difference
- training - self-supervised/supervised learning
- inference - generates outputs
  - in-context learning, prompt engineering

goal and scope of LLM project



# **Transformer**

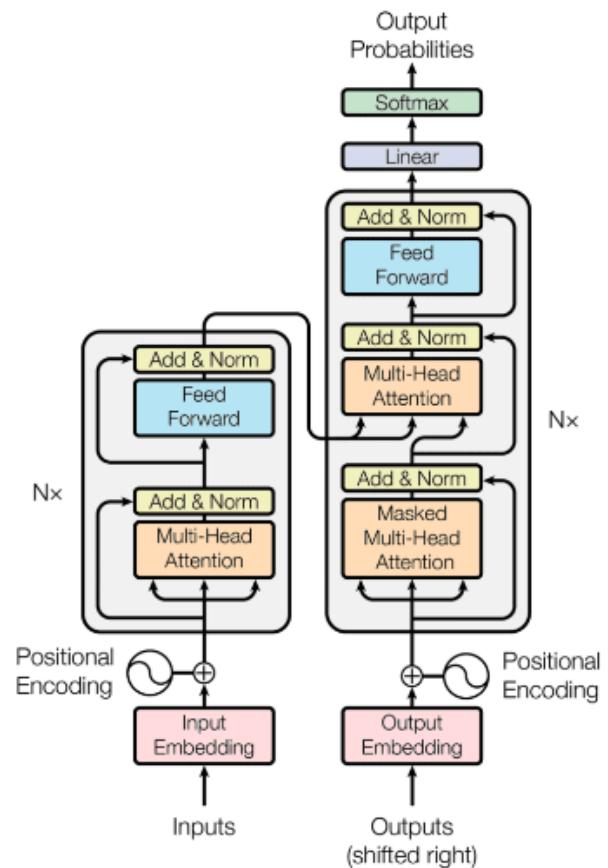
## **LLM architectural secret (or known) sauce**

### **Transformer - simple parallelizable attention mechanism**

A. Vaswani, et al. Attention is All You Need, 2017

# Transformer architecture

- encoding-decoding architecture
  - input embedding space → multi-head & mult-layer representation space → output embedding space
- additive positional encoding - information regarding order of words @ input embedding
- multi-layer and multi-head attention followed by addition / normalization & feed forward (FF) layers
- *(relatively simple) attentions*
  - single-head (scaled dot-product) / multi-head attention
  - self attention / encoder-decoder attention
  - masked attention
- benefits
  - *evaluate dependencies between arbitrarily distant words*
  - has recurrent nature w/o recurrent architecture → parallelizable → fast w/ additional cost in computation



## Single-head scaled dot-product attention

- values/keys/queries denote value/key/query *vectors*,  $d_k$  &  $d_v$  are lengths of keys/queries & vectors
- we use *standard* notions for matrices and vectors - not transposed version that (almost) all ML scientists (wrongly) use
- output: weighted-average of values where weights are attentions among tokens
- assume  $n$  queries and  $m$  key-value pairs

$$Q \in \mathbf{R}^{d_k \times n}, K \in \mathbf{R}^{d_k \times m}, V \in \mathbf{R}^{d_v \times m}$$

- attention! outputs  $n$  values (since we have  $n$  queries)

$$\text{Attention}(Q, K, V) = V \text{softmax} \left( K^T Q / \sqrt{d_k} \right) \in \mathbf{R}^{d_v \times n}$$

- *much simpler attention mechanism than previous work*
  - attention weights were output of complicated non-linear NN

## Single-head - close look at equations

- focus on  $i$ th query,  $q_i \in \mathbf{R}^{d_k}$ ,  $Q = [ \quad - \quad q_i \quad - \quad ] \in \mathbf{R}^{d_k \times n}$
- assume  $m$  keys and  $m$  values,  $k_1, \dots, k_m \in \mathbf{R}^{d_k}$  &  $v_1, \dots, v_m \in \mathbf{R}^{d_v}$

$$K = [ \ k_1 \ \ \cdots \ \ k_m \ ] \in \mathbf{R}^{d_k \times m}, V = [ \ v_1 \ \ \cdots \ \ v_m \ ] \in \mathbf{R}^{d_v \times m}$$

- then

$$K^T Q / \sqrt{d_k} = \begin{bmatrix} & & \vdots \\ - & k_j^T q_i / \sqrt{d_k} & - \\ & & \vdots \end{bmatrix}$$

e.g., dependency between  $i$ th output token and  $j$ th input token is

$$a_{ij} = \exp \left( k_j^T q_i / \sqrt{d_k} \right) / \sum_{j=1}^m \exp \left( k_j^T q_i / \sqrt{d_k} \right)$$

- value obtained by  $i$ th query,  $q_i$  in  $\text{Attention}(Q, K, V)$

$$a_{i,1}v_1 + \cdots + a_{i,m}v_m$$

## Multi-head attention

- evaluate  $h$  single-head attentions (in parallel)
- $d_e$ : dimension for embeddings
- embeddings

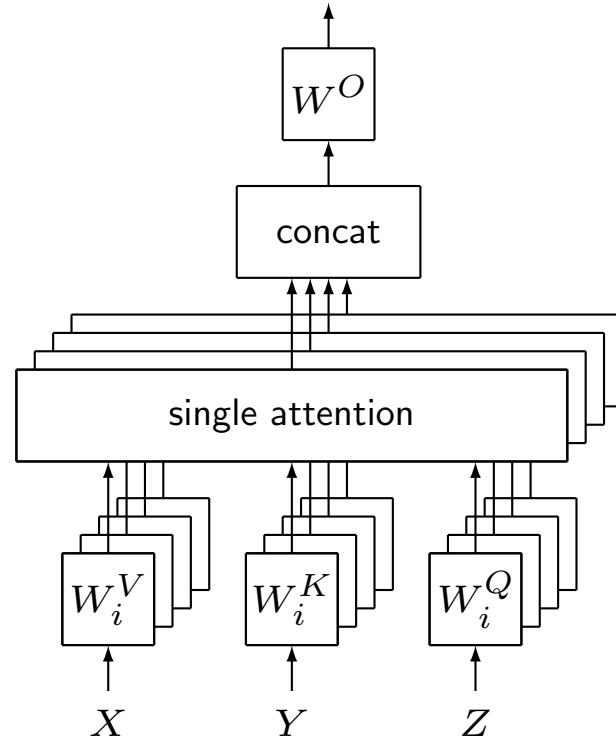
$$X \in \mathbb{R}^{d_e \times m}, Y \in \mathbb{R}^{d_e \times m}, Z \in \mathbb{R}^{d_e \times n}$$

e.g.,  $n$ : input sequence length &  $m$ : output sequence length in machine translation

- $h$  key/query/value weight matrices:  $W_i^K, W_i^Q \in \mathbb{R}^{d_k \times d_e}$ ,  $W_i^V \in \mathbb{R}^{d_v \times d_e}$  ( $i = 1, \dots, h$ )
- linear output layers:  $W^O \in \mathbb{R}^{d_e \times hdv}$
- *multi-head attention!*

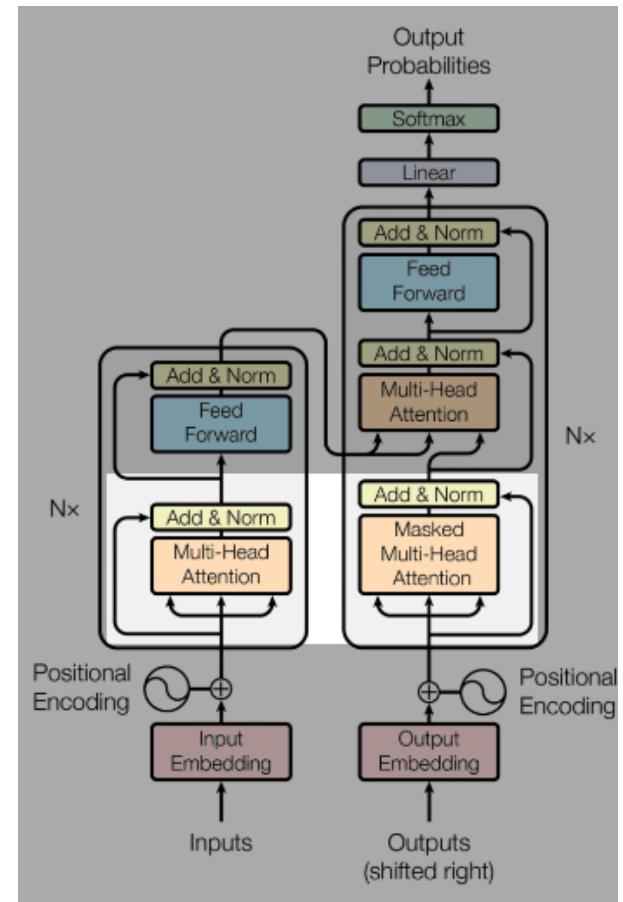
$$W^O \begin{bmatrix} A_1 \\ \vdots \\ A_h \end{bmatrix} \in \mathbb{R}^{d_e \times n},$$

$$A_i = \text{Attention}(W_i^Q Z, W_i^K Y, W_i^V X) \in \mathbb{R}^{d_v \times n}$$



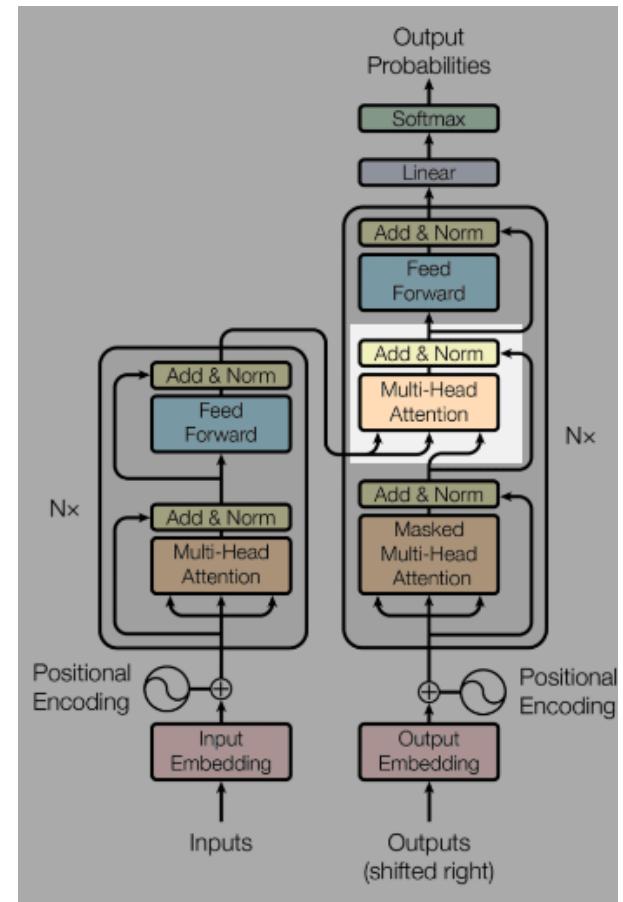
# Self attention

- $m = n$
- encoder
  - keys & values & queries ( $K, V, Q$ ) come from same place (from previous layer)
  - every token attends to every other token in input sequence
- decoder
  - keys & values & queries ( $K, V, Q$ ) come from same place (from previous layer)
  - every token attends to other tokens up to that position
  - prevent leftward information flow to right to preserve causality
  - assign  $-\infty$  for illegal connections in softmax (masking)



## Encoder-decoder attention

- $m$ : length of input sequence
- $n$ : length of output sequence
- $n$  queries ( $Q$ ) come from previous decoder layer
- $m$  keys /  $m$  values ( $K, V$ ) come from output of encoder
- every token in output sequence attends to every token in input sequence

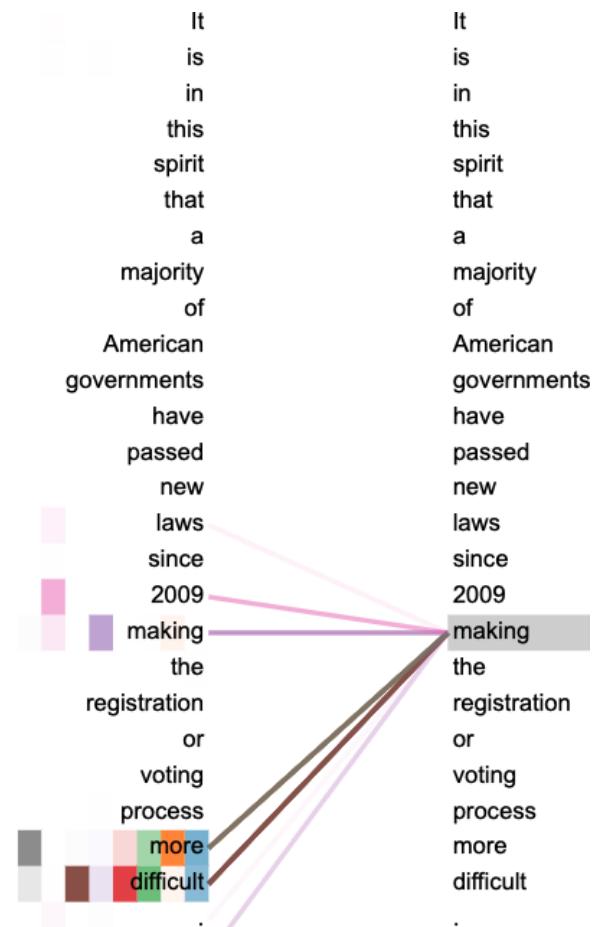


## Visualization of self attentions

example sentence

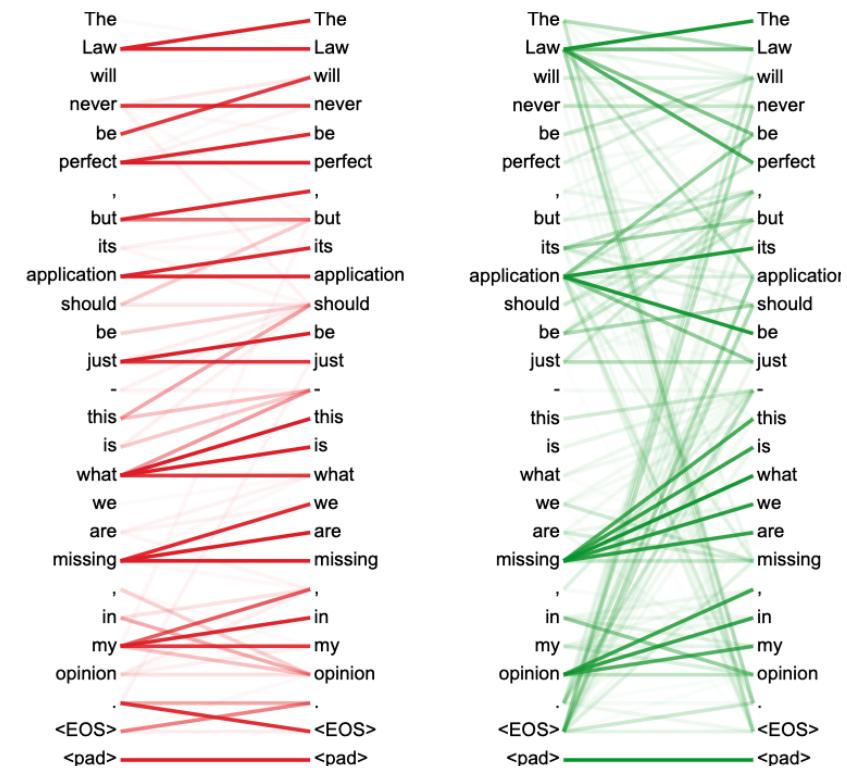
"It is in this spirit that a majority of American governments have passed new laws since 2009 making the registration or voting process more difficult."

- self attention of encoder (of a layer)
  - right figure
    - show dependencies between "making" and other words
    - different columns of colors represent different heads
  - "making" has strong dependency to "2009", "more", and "difficult"

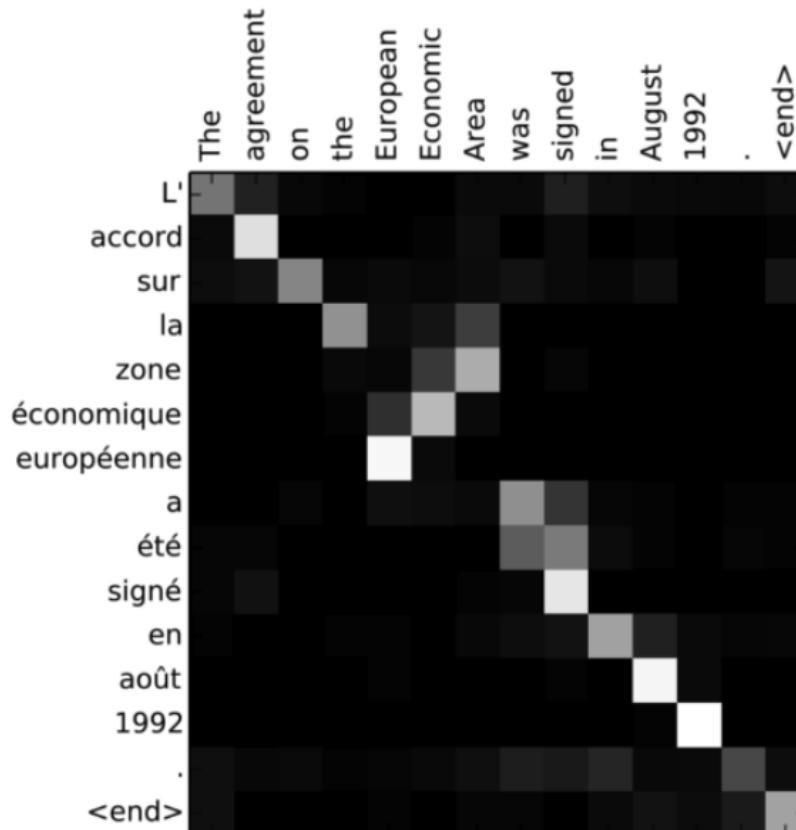


## Visualization of multi-head self attentions

- self attentions of encoder for two heads (of a layer)
  - different heads represent different structures  
→ advantages of multiple heads
  - multiple heads work together to collectively yield good results
  - dependencies *not* have absolute meanings (like embeddings in collaborative filtering)
  - randomness in resulting dependencies exists due to stochastic nature of ML training



## Visualization of encoder-decoder attentions



- machine translation: English → French
  - input sentence: “The agreement on the European Economic Area was signed in August 1992.”
  - output sentence: “L’ accord sur la zone économique européenne a été signé en août 1992.”
- encoder-decoder attention reveals relevance between
  - European ↔ européenne
  - Economic ↔ économique
  - Area ↔ zone

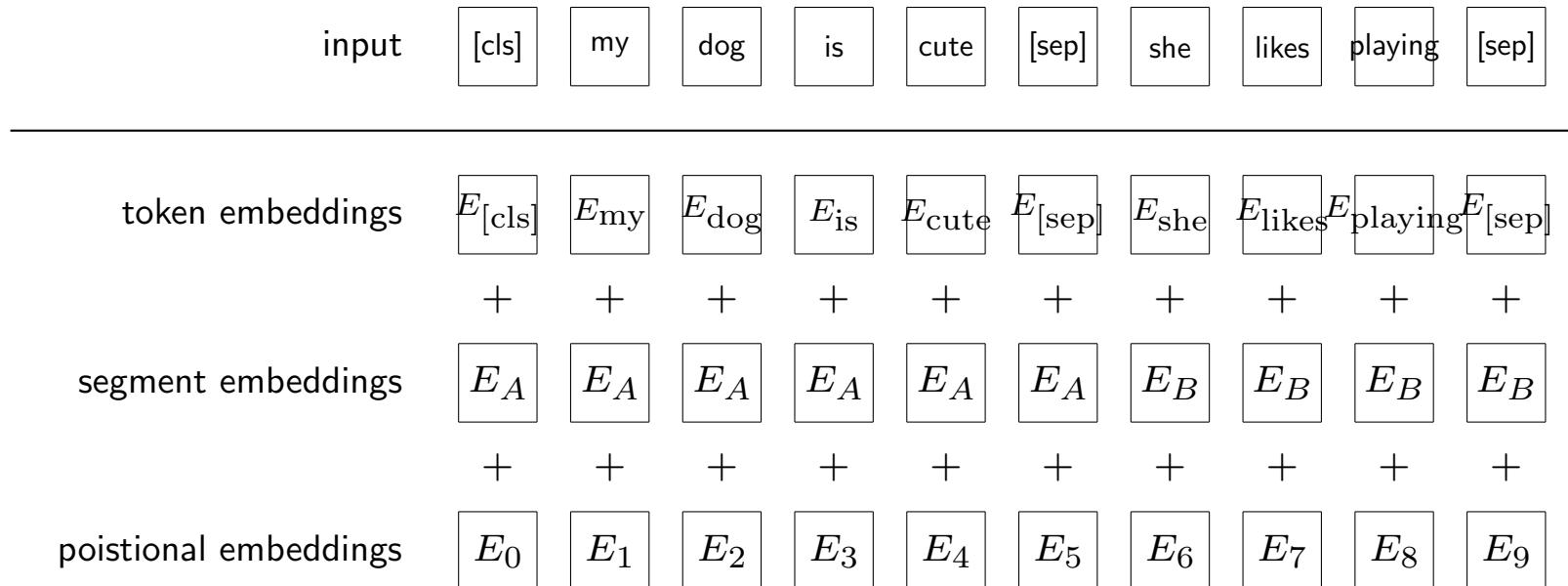
## Model complexity

- computational complexity
  - $n$ : sequence length,  $d$ : embedding dimension
  - complexity per layer - self-attention:  $\mathcal{O}(n^2d)$ , recurrent:  $\mathcal{O}(1)$
  - sequential operations - self-attention:  $\mathcal{O}(1)$ , recurrent:  $\mathcal{O}(n)$
  - maximum path length - self-attention:  $\mathcal{O}(1)$ , recurrent:  $\mathcal{O}(n)$
- *massive parallel processing, long context windows*
  - makes NVidia more competitive, hence profitable!
  - makes SK Hynix prevail HBM market!

## **Variants of Transformer**

## Bidirectional encoder representations from transformers (BERT)

- Bidirectional Encoder Representations from Transformers [DCLT19]
- pre-train deep bidirectional representations from unlabeled text
- fine-tunable for multiple purposes



## **Implications & Challenges**

## Multimodal learning

- understand information from multiple modalities, *e.g.*, text, images, audio, and video
- representation learning
  - language representation + image / video / text / audio representation
  - learn multimodal representations together
- outputs
  - captions for images, videos with narration, musics with lyrics
- collaboration among different modalities
  - understand image world (open system) using language (closed system)



## Implications of success of LLMs

- (very) many researchers change gears towards LLM
  - from computer vision (CV), speech, music, video, even reinforcement learning
- *LLM is not (only) about languages . . .*
  - humans have . . .
    - evolved and optimized (natural) language structures for eons
    - handed down knowledge using natural languages for thousands of years
  - natural language optimized (in human brains) through *thousands of generations by evolution*
  - *can connect non-linguistic world (open system) using language structures (closed system)*

## Challenges in LLMs

- *hallucination - can give entirely plausible outcome that is false*
- data poison attack
- unethical or illegal content generation
- huge resource necessary for both training & inference
- model size - need compact models
- outdated knowledge - can be couple of years old
- lack of reproducibility
- *biases - more on this later . . .*

do not, though, focus on downsides but on *infinite possibilities!*

- it evolves like internet / mobile / electricity
- only “tip of the iceberg” found & revealed

**genAI**

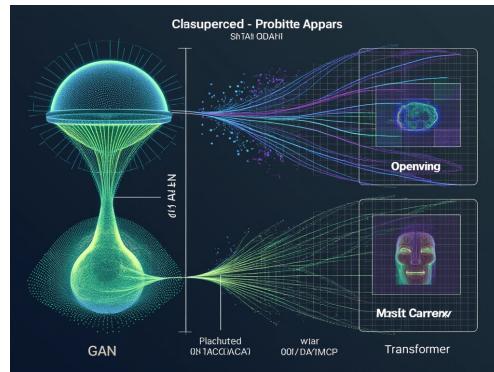
## **Definition of genAI**

## Generative AI

- genAI refers to systems capable of producing new (& original) contents based on patterns learned from training data (representation learning)
  - as opposed to discriminative models for, *e.g.*, classification, prediction & regression
  - here content can be text, images, audio, video, *etc.* - what about smell & taste?
- genAI model examples
  - generative adversarial networks (GANs), variational autoencoders (VAEs), diffusion models, Transformers



by Midjourney



by Grok 2 mini



by Generative AI Lab

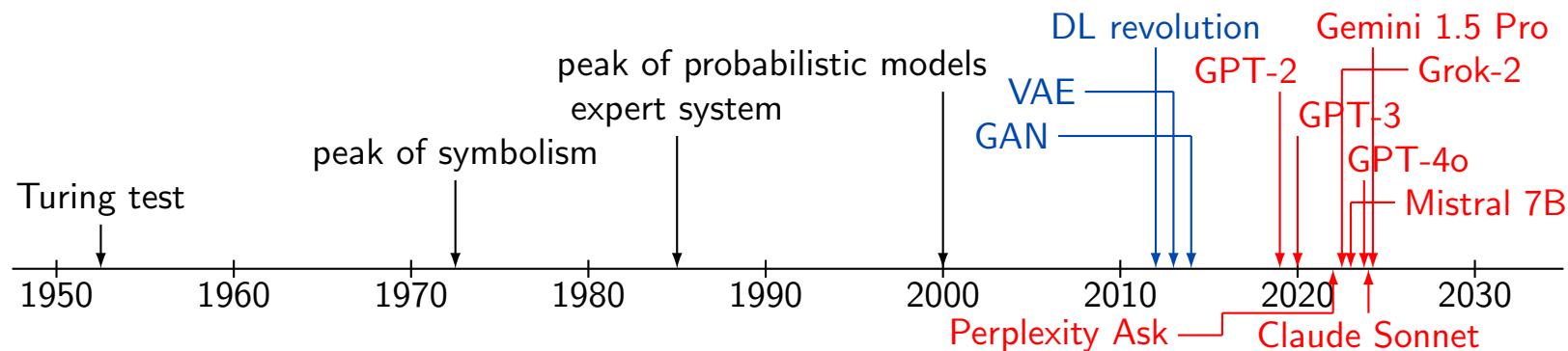
## Examples of genAI in action

- text generation
  - Claude, ChatGPT, Mistral, Perplexity, Gemini, Grok
  - conversational agent writing articles, code & even poetry
- image generation
  - DALL-E - creates images based on textual descriptions
  - Stable Diffusion - uses diffusion process to generate high-quality images from text prompts (by denoising random noise)
  - MidJourney - art and visual designs generated through deep learning
- music generation
  - Amper Music - generates unique music compositions
- code generation
  - GitHub Copilot - generates code snippets based on natural language prompts

# History of genAI

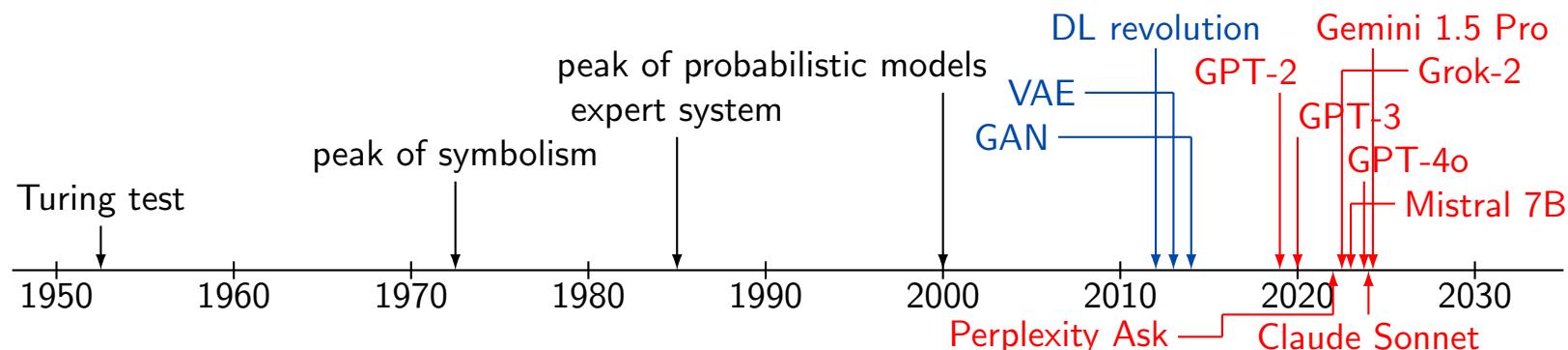
## Birth of AI - early foundations & precursor technologies

- 1950s ~ 1970s
  - Alan Turing - concept of “*thinking machine*” & *Turing test* to evaluate machine intelligence (1950s)
  - *symbolists* (as opposed to connectionists) - early AI focused on symbolic reasoning, logic & problem-solving - Dartmouth Conference in 1956 by *John McCarthy, Marvin Minsky, Allen Newell & Herbert A. Simon*
  - precursor technologies - genetic algorithms (GAs), Markov chains & *hidden Markov models (HMMs)* - laying foundation for generative processes (1970s ~)



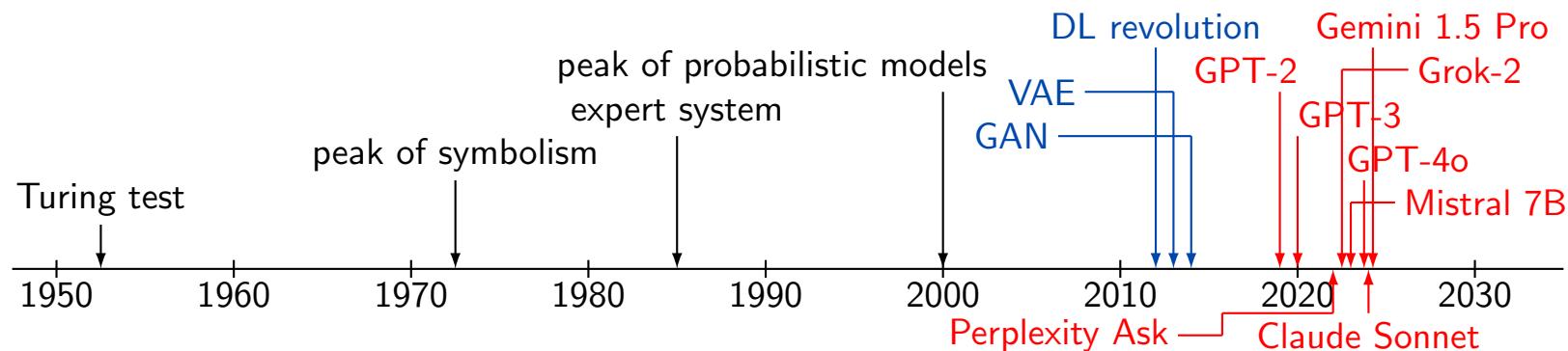
## Rule-based systems & probabilistic models

- 1980s ~ early 2000s
  - *expert systems* (1980s) - AI systems designed to mimic human decision-making in specific domains
  - development of neural networks (NN) w/ backpropagation *training multi-layered networks* - setting stage for way more complex generative models
  - *probabilistic models* (including network models, *i.e.*, Bayesian networks) & Markov models - laying groundwork for data generation & pattern prediction



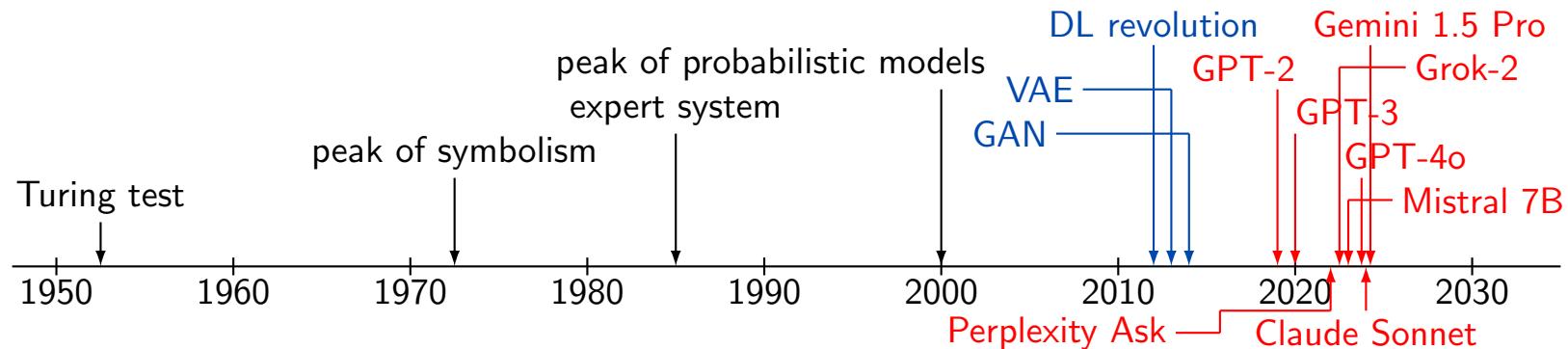
## Rise of deep learning & generative models

- 2010s - breakthrough in genAI
  - *deep learning (DL) revolution* - advances in GPU computing and data availability led to the rapid development of deep neural networks.
  - *variational autoencoder (VAE)* (2013) - by Kingma and Welling - learns mappings between input and latent spaces
  - *generative adversarial network (GAN)* (2014) - by Ian Goodfellow - game-changer in generative modeling where two NNs compete each other to create realistic data
    - widely used in image generation & creative tasks



## Transformer models & multimodal AI

- late 2010s ~ Present
  - Transformer architecture (2017) - by Vaswani et al.
    - *revolutionized NLP*, e.g., LLM & various genAI models
  - GPT series - generative pre-trained transformer
    - GPT-2 (2019) - generating human-like texts - *marking leap in language models*
    - GPT-3 (2020) - 175B params - set *new standards for LLM*
  - multimodal systems - DALL-E & CLIP (2021) - *linking text and visual data*
  - emergence of diffusion models (2020s) - new approach for generating high-quality images - progressively “denoising” random noise (DALL-E 2 & Stable Diffusion)



# **Mathy Views on genAI**

## genAI models

- definition of generative model

$$\boxed{\mathcal{Z}} \xrightarrow{g_{\theta}(z)} \boxed{\mathcal{X}}$$

- *generate samples in original space,  $\mathcal{X}$ , from samples in latent space,  $\mathcal{Z}$*
- $g_{\theta}$  is parameterized model e.g., CNN / RNN / Transformer / diffuction-based model
- training
  - finding  $\theta$  that minimizes/maximizes some (statistical) loss/merit function so that  $\{g_{\theta}(z)\}_{z \in \mathcal{Z}}$  generates plausible point in  $\mathcal{X}$
- inference
  - random samples  $z$  to generated target samples  $x = g_{\theta}(z)$
  - e.g., image, text, voice, music, video

## VAE - early genAI model

- variational auto-encoder (VAE) [KW19]

$$\boxed{\mathcal{X}} \xrightarrow{q_\phi(z|x)} \boxed{\mathcal{Z}_0} \xrightarrow{p_\theta(x|z)} \boxed{\mathcal{X}}$$

- log-likelihood & ELBO - for any  $q_\phi(z|x)$

$$\begin{aligned} \log p_\theta(x) &= \mathbf{E}_{z \sim q_\phi(z|x)} \log p_\theta(x) = \mathbf{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \cdot \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ &= \mathcal{L}(\theta, \phi; x) + D_{KL}(q_\phi(z|x) \| p_\theta(z|x)) \geq \mathcal{L}(\theta, \phi; x) \end{aligned}$$

- (indirectly) maximize likelihood by maximizing evidence lower bound (ELBO)

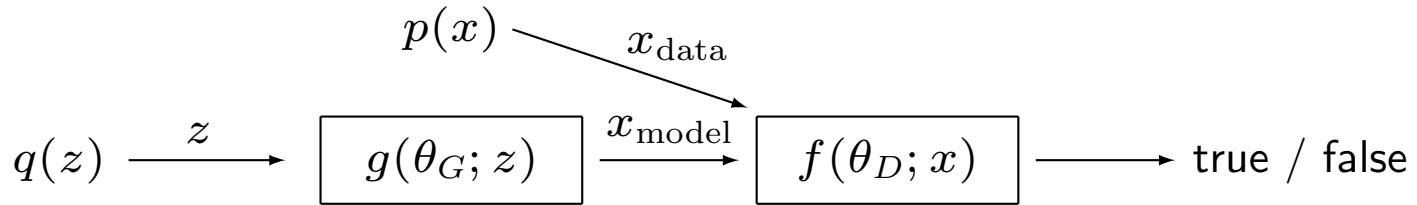
$$\mathcal{L}(\theta, \phi; x) = \mathbf{E}_{z \sim q_\phi(z|x)} \log \frac{p_\theta(x, z)}{q_\phi(z|x)}$$

- generative model

$$p_\theta(x|z)$$

## GAN - early genAI model

- generative adversarial networks (GAN) [GPAM<sup>+</sup>14]



- value function

$$V(\theta_D, \theta_G) = \mathbf{E}_{x \sim p(x)} \log f(\theta_D; x) + \mathbf{E}_{z \sim q(z)} \log(1 - f(\theta_D; g(\theta_G; z)))$$

- modeling via playing min-max game

$$\min_{\theta_G} \max_{\theta_D} V(\theta_D, \theta_G)$$

- generative model

$$g(\theta_G; z)$$

- variants: conditional / cycle / style / Wasserstein GAN

## genAI - LLM

- *maximize conditional probability*

$$\underset{\theta}{\text{maximize}} \ d(p_{\theta}(x_t|x_{t-1}, x_{t-2}, \dots), p_{\text{data}}(x_t|x_{t-1}, x_{t-2}, \dots))$$

where  $d(\cdot, \cdot)$  distance measure between probability distributions

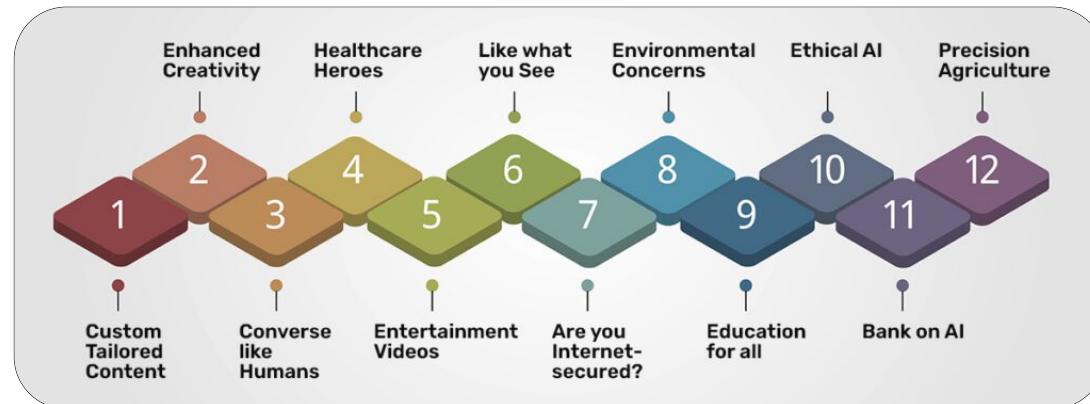
- previous sequence:  $x_{t-1}, x_{t-2}, \dots$
- next token:  $x_t$
- $p_{\theta}$  represented by (extremely) complicated model
  - e.g., containing multi-head & multi-layer Transformer architecture inside
- model parameters, e.g., for Llama2

$$\theta \in \mathbf{R}^{70,000,000,000}$$

## **Current Trend & Future Perspectives**

## Current trend of genAI

- rapid advancement in language models & multimodal AI capabilities
- rise of AI-assisted creativity & productivity tools
- growing adoption across industries
  - creative industries - design, entertainment, marketing, software development
  - life sciences - healthcare, medical, biotech
- infrastructure & accessibility, *e.g.*, Hugging Face democratizes AI development
- integration with cloud platforms & enterprise-level tools
- increased focus on AI ethics & responsible development



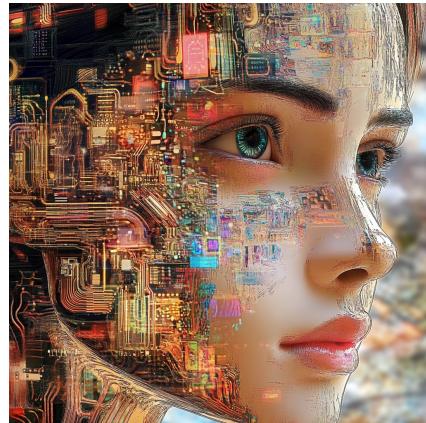
## Industry & business impacts

- how genAI is transforming industries
  - creative industries - content creation - advertising, gaming, film
  - life science - enhance research, drug discovery & personalized treatments
  - finance - automating document generation, risk modeling & fraud detection
  - manufacturing & Design - rapid prototyping, 3D modeling & optimization
  - business operations - automate routine tasks to boost productivity



## Future perspectives of genAI

- hyper-personalization - highly personalized content for individual users - music, products & services
- AI ethics & governance - concerns over deepfakes, misinformation & bias
- interdisciplinary synergies - integration with other fields such as quantum computing, neuroscience & robotics
- human-AI collaboration - augment human creativity rather than replace it
- energy efficiency - have to figure out how to dramatically reduce power consumption



# **Selected References & Sources**

## Selected references & sources

- Chris Miller “Chip War: The Fight for the World’s Most Critical Technology” (2022)
- Daniel Kahneman “Thinking, Fast and Slow” (2011)
- M. Shanahan “Talking About Large Language Models” (2022)
- A.Y. Halevy, P. Norvig, and F. Pereira “Unreasonable Effectiveness of Data” (2009)
- A. Vaswani, et al. “Attention is all you need” @ NeurIPS (2017)
- S. Yin, et. al. “A Survey on Multimodal LLMs” (2023)
- I.J. Goodfellow, ..., Y. Bengio “Generative adversarial networks (GAN)” (2014)
- T. Kuiken “Artificial Intelligence in the Biological Sciences: Uses, Safety, Security, and Oversight” (2023)
- Stanford Venture Investment Groups
- CEOs & CTOs @ startup companies in Silicon Valley
- VCs on Sand Hill Road - Palo Alto, Menlo Park, Woodside in California, USA

# **References**

## References

- [DCLT19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [GPAM<sup>+</sup>14] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [HGH<sup>+</sup>22] Sue Ellen Haupt, David John Gagne, William W. Hsieh, Vladimir Krasnopolksy, Amy McGovern, Caren Marzban, William Moninger, Valliappa Lakshmanan, Philippe Tissot, and John K. Williams. The history and practice of AI in the environmental sciences. *Bulletin of the American Meteorological Society*, 103(5):E1351 – E1370, 2022.
- [KW19] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends in Machine Learning*, 12(4):307–392, 2019.

- [Mil22] Chris Miller. *Chip war: fight for the world's most critical technology*. New York: Scribner, 2022.
- [MLZ22] Louis-Philippe Morency, Paul Pu Liang, and Amir Zadeh. Tutorial on multimodal machine learning. In Miguel Ballesteros, Yulia Tsvetkov, and Cecilia O. Alm, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Tutorial Abstracts*, pages 33–38, Seattle, United States, July 2022. Association for Computational Linguistics.
- [VSP<sup>+</sup>17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of 31st Conference on Neural Information Processing Systems (NIPS)*, 2017.
- [YFZ<sup>+</sup>24] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models, 2024.

**Thank You**