# Evolution of System Designs from an AI Engineer Perspective

**Yangqing Jia**

**Lepton AI (now part of NVIDIA)**

# Who is this dude?

**Researcher -> Engineer -> Entrepreneur**

- Berkeley PhD 2009-2013

- Google, Facebook, Alibaba

- Did a bunch of open source work

- Lepton AI -> NVidia

0
# Demystify "LLM" and "AGI"

# "The Chinese Typewriter Problem"

For more comparative culture discussions, check
Thomas S Mullaney, "The Chinese Typewriter: a History"

# 中文打字機基本字盤表

| | 活用字(二) | 常用字 | 特用字 | 活用字(二) | 間用字 |
|---|---|---|---|---|---|



中文打字機基本字盤表（字盤排列圖）

# 中文打字機基本字盤表

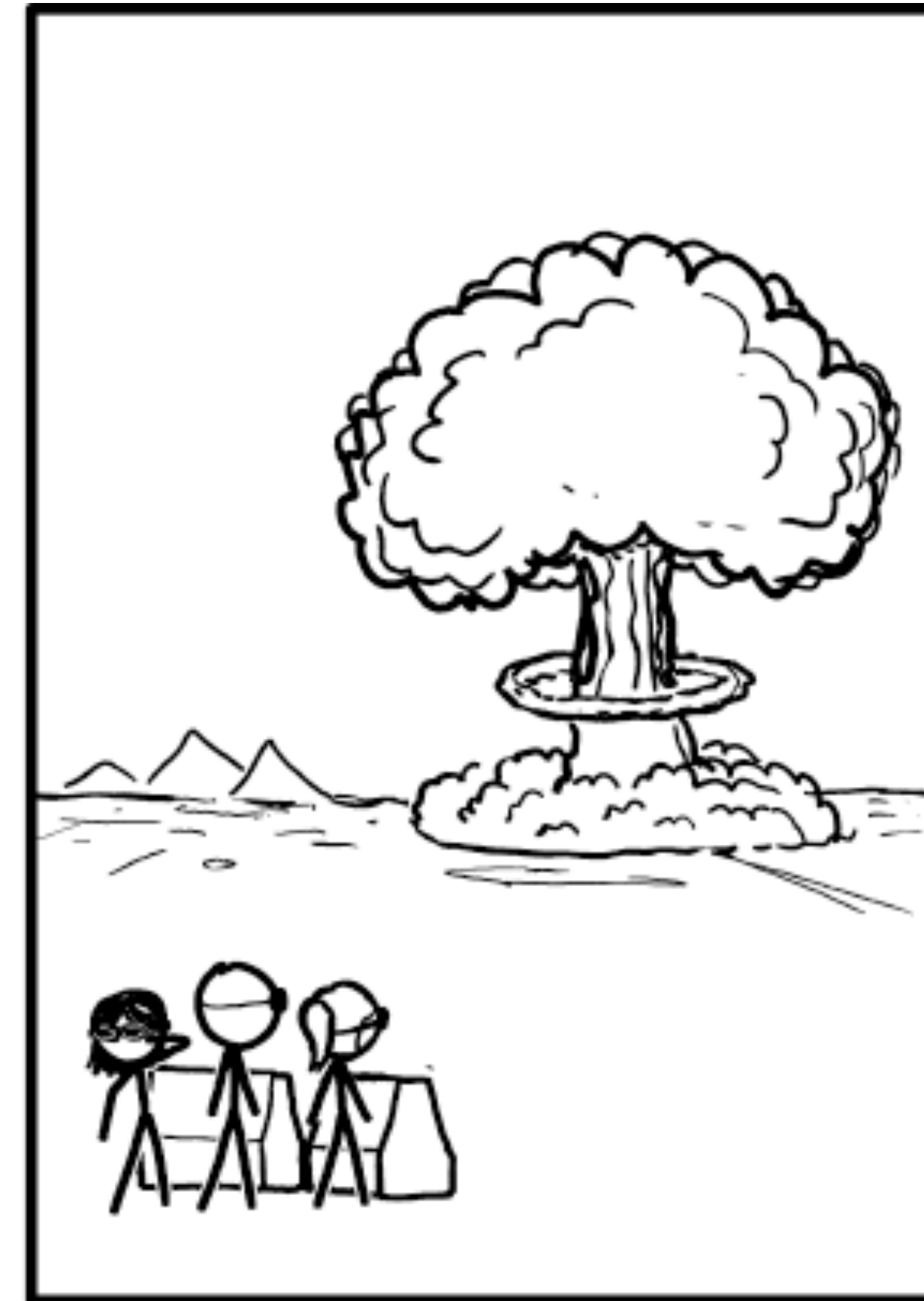| 間用字 | 活用字（二） | 常用字 | 特用字 | 活用字（一） | 間用字 |
|---|---|---|---|---|---|

講識認　談話說

國央中華

汽油漆漏洞

# What I would like to cover…

**My learnings along the road of AI models, apps, and infra**

- New algorithms continue to grow LLM algorithms.

- Application space thrives.

- AI Infra has become the 3rd pillar of enterprise IT strategy.

- … and interestingly, we see history repeat itself.

- If you would like - ups and downs of a startup.
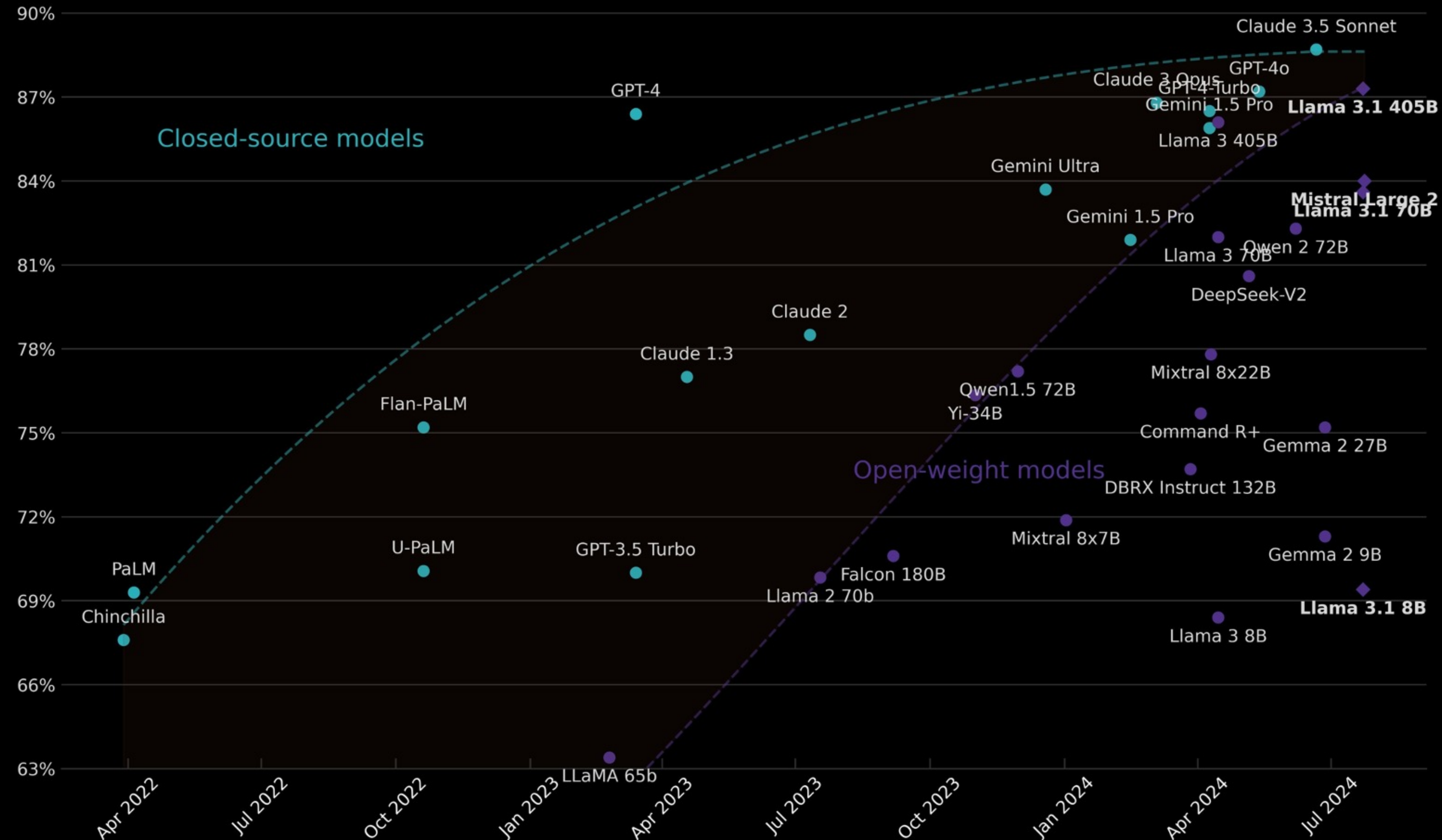
# 1
## New algorithms
continue to drive
LLM models.

**Closed-source vs. open-weight models**

Llama 3.1 405B closes the gap with closed-source models for the first time in history.

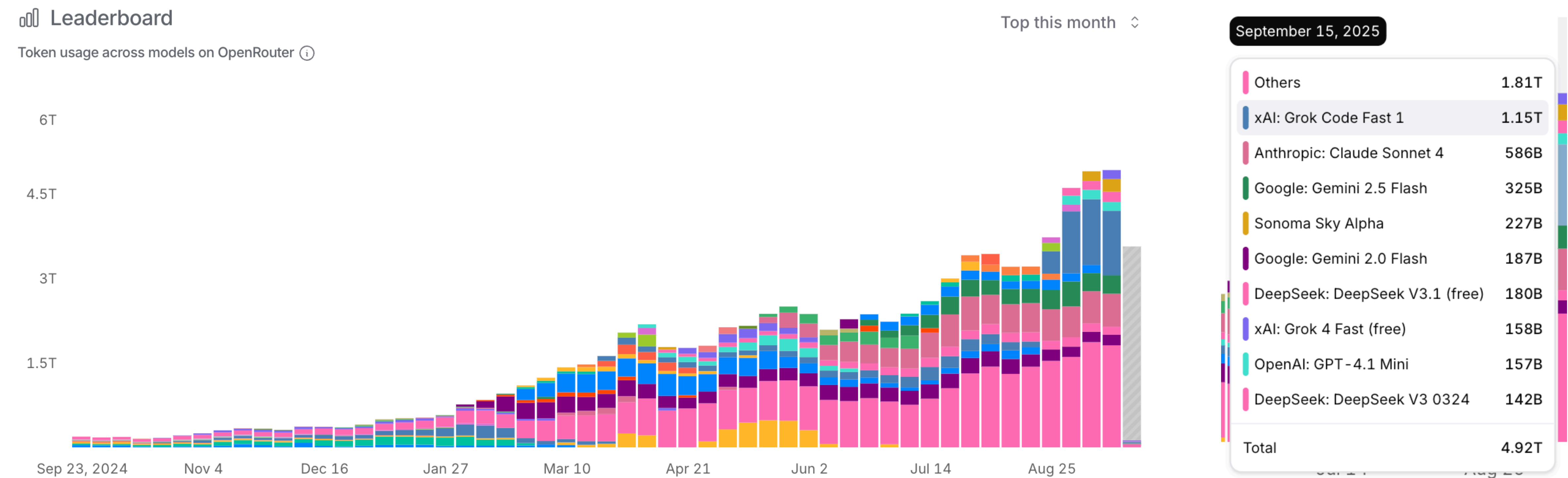@maximelabonne

- And the latest…
  - GPT-5
  - Grok
  - Gemini
  - Kimi
  - DeepSeek
  - Qwen
  - GPT-OSS

# There does not seems to be "a bubble"

## Consumption continue growing - not only model training.



**Leaderboard**

Token usage across models on OpenRouter ⓘ

Top this month ⇅

**September 15, 2025**

| | |
|---|---|
| Others | 1.81T |
| xAI: Grok Code Fast 1 | 1.15T |
| Anthropic: Claude Sonnet 4 | 586B |
| Google: Gemini 2.5 Flash | 325B |
| Sonoma Sky Alpha | 227B |
| Google: Gemini 2.0 Flash | 187B |
| DeepSeek: DeepSeek V3.1 (free) | 180B |
| xAI: Grok 4 Fast (free) | 158B |
| OpenAI: GPT-4.1 Mini | 157B |
| DeepSeek: DeepSeek V3 0324 | 142B |
| **Total** | **4.92T** |

Source: openrouter.ai

# New Algorithms Drive Continued Improvements
## My personal opinion and historical analogies...

| Date | Nov 2022 | Dec 2023 | Sep 2024 | Jan 2025 (and earlier) |
|------|----------|----------|----------|------------------------|
| Algorithm | GPT (3.5) | MoE (Mixtral 8x7B) | Test time Scaling | **Reinforcement Learning** |
| Analogies | AlexNet (structural innovation) | Ensemble Learning Inception/ResNet | Fully convolutional network Multi-instance learning | General RL GANs |

# Hype Cycle for Generative AI, 2025

Plateau will be reached:  ○ < 2 yrs.   ◐ 2–5 yrs.   ● 5–10 yrs.   ⚠ >10 yrs.

EXPECTATIONS

**Peak of Inflated Expectations**

Agentic AI

Multimodal Generative AI

Open-Source LLMs

GenAI Application
Orchestration Frameworks

AI-Augmented
Software Engineering

Prompt Engineering

Reinforcement Learning
From Human Feedback

Domain-Specific GenAI Models

GraphRAG

Agent Development Frameworks

Generative AI Workload Accelerators

Artificial General Intelligence

AI TRiSM

Model Ops

Generative AI-Enabled Applications

Retrieval-Augmented Generation

Innovation Trigger ⓘ   Peak of Inflated Expectations   Trough of Disillusionment   Slope of Enlightenment   Plateau of Productivity

**TIME**

Gartner®

**2**
# ToC apps thrive.
# ToB apps are hopeful & nascent.



Image source: xkcd

**Perfect** app experience is correlated, but independent from models.

https://elmo.chat/

# ⚡ Elmo is your AI companion to create summaries and insights

Featured on Chrome Web Store

Developed by Established Publisher

## 🌐 Install from Chrome Web Store

Free / No Account Required / Supports Multiple Languages

Work seamlessly across HTML, Youtube, PDF, and Google Docs.

---

**Guillermo Rauch** ✔

This looks really handy

---

**Bing Xu** ✔

Elmo is the co-pilot for Chrome.
Super helpful when reading new

---

**Diogo Santos** ✔

Elmo is your AI extension 🧠 for Chrome to create summaries, insights, and extended knowledge. What does Elmo offer? 🤔 ✓ Summaries and highlights; ✓ Keep asking questions; ✓ Dive deep into keywords; ✓ Chat with PDFs; ✓

---

**Alvin-GenAI** ✔

I've downloaded https://elmo.chat an AI Chrome extension to create summaries, insights and extended knowledge. It does - summarizes your websites right next to the page, - summarizes PDF files too, - summarizes

---

**Tulsi Prasad** ✔

I tried out this AI chrome extension yesterday! 1. No product hunt launches 2. Featured on Chrome Web Store like an OG 3. No signup needed, just plug and play! Works pretty accurate on blogs and textual content, but videos are

# Consumer App Landscape is Highly Fluid
## Due to the continued improvement of foundational models

## The Top 50 Gen AI Web Products, by Unique Monthly Visits

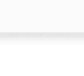| | | | | |
|---|---|---|---|---|
| 1. ChatGPT | 11. remove bg | 21. IIElevenLabs | 31. manus | 41. replit |
| 2. Gemini | 12. Doubao | 22. Hugging Face | 32. CIVITAI | 42. candy.ai |
| 3. deepseek | 13. NotebookLM | 23. Lovable | 33. KlingAI | 43. JuicyChat |
| 4. Grok | 14. SPICYCHAT.AI | 24. Crushon AI | 34. cutout.pro | 44. VEED |
| 5. character.ai | 15. SUNO | 25. GAMMA | 35. Adot | 45. Hailuo AI |
| 6. perplexity | 16. QuillBot | 26. CURSOR | 36. DeepAI | 46. Meta AI |
| 7. Claude | 17. Kimi | 27. Pixelcut | 37. Poe | 47. Remaker |
| 8. JanitorAI | 18. PolyBuzz | 28. Midjourney | 38. ZeroGPT | 48. JOI |
| 9. Quark | 19. SERART.AI | 29. TurboScribe | 39. Google Labs | 49. Monica |
| 10. Google AI Studio | 20. Qwen3 | 30. Photoroom | 40. Leonardo.Ai | 50. ourdream.ai |

a16z Consumer

# Consumer App Landscape is Highly Fluid
## Due to the continued improvement of foundational models

## The Top 50 Gen AI Web Products, by Unique Monthly Visits

| # | | # | | # | | # | | # | |
|---|---|---|---|---|---|---|---|---|---|
| 1. | ChatGPT | 11. | remove.bg | 21. | IIElevenLabs | 31. | manus | 41. | replit |
| 2. | Gemini | 12. | Doubao | 22. | Hugging Face | 32. | CIVITAI | 42. | candy.ai |
| 3. | deepseek | 13. | Notebook M | 23. | Lovable | 33. | Kling AI | 43. | JuicyChat |
| 4. | Grok | 14. | SPICYC | | | | | | |
| 5. | character.ai | 15. | SUNO | | | | | | |
| 6. | perplexity | 16. | QuillB | | | | | | |
| 7. | Claude | 17. | Kimi | | | | | | |
| 8. | JanitorAI | 18. | PolyBu | | | | | | |
| 9. | Quark | 19. | SEART. | | | | | | |
| 10. | Google AI Studio | 20. | Qwen | | | | | | |

## Top Gen AI Consumer Web Products: Newcomers

| # | | # | | # | | # | |
|---|---|---|---|---|---|---|---|
| 4. | Grok | 23. | Lovable | 41. | replit | 48. | JOI |
| 9. | Quark | 29. | TurboScribe | 43. | JuicyChat | 50. | ourdream.ai |
| 20. | Qwen3 | 31. | manus | 47. | Remaker | | |

# Prosumers: willingness to pay drives revenue

Cursor's Anysphere nabs $9.9B valuation, soars past $500M ARR

Runway, best known for its video-generating AI models, raises $308M

Plaud's $169 ChatGPT-powered NotePin has a permanent place in my travel bag

How Eleven Labs hit $200M revenue with a 291 person team in 2025.

Research Lab: Exploring new frontiers of voice generation. We are dedicated to researching and implementing innovative techniques in voice artificial intelligence (AI) to enhance the appeal of content across different languages and voices. Our goal is to reach new audiences and viewers by ensuring a more immersive and engaging experience.

| 2022 | $200M | 2122.56% | $381M |
|------|-------|----------|-------|
| Founded | 2025 Revenue | YOY | Funding |

# Still, there may be a big enterprise market
## Enterprises used to be slow. Now faster.



Years from $1 to $100M ARR
Public Enterprise Data, Productivity, and AI Companies

Notes: Assumes 24 months from founding to $1M ARR if not disclosed. Assumes exponential growth from founding date to most recent year company reached $100M ARR
Source: Historical filings and public company announcements

ICONIQ | Growth    glean

# Preliminary generative AI tech stack

Users

**Legend:**
- Apps
- Models
- Infrastructure

## Apps

*End-user facing B2B and B2C applications without proprietary models*

Examples: Jasper, Github Copilot

### End-to-End Apps

*End-user facing applications with proprietary models*

Examples: Midjourney, Runway

### Closed-Source Foundation Models

*Large-scale, pre-trained models exposed via APIs*

Examples: GPT-3 (OpenAI)

### Model Hubs

*Platforms to share and host models*

Examples: Hugging Face, Replicate

### Open-Source Foundation Models

*Models released as trained weights*

Examples: Stable Diffusion (Stability)

# Preliminary generative AI tech stack

Users

| | |
|---|---|
| Apps | |
| Models | |
| Infrastructure | |

## Apps

*End-user facing B2B and B2C applications __without__ proprietary models*

Examples: Jasper, Github Copilot

### End-to-End Apps

*End-user facing applications __with__ proprietary models*

Examples: Midjourney, Runway

### Closed-Source Foundation Models

*Large-scale, pre-trained models exposed via APIs*

Examples: GPT-3 (OpenAI)

### Model Hubs

*Platforms to share and host models*

Examples: Hugging Face, Replicate

### Open-Source Foundation Models

*Models released as trained weights*

Examples: Stable Diffusion (Stability)

## Cloud Platforms

*Compute hardware exposed to developers in a cloud deployment model*

Examples: AWS, GCP, Azure, Coreweave

## Compute Hardware

*Accelerator chips optimized for model training and inference workloads*

Examples: GPUs (Nvidia), TPUs (Google)

a16z Enterprise

**3**
# AI infra is the 3rd pillar in IT strategy.

"The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin."

- Professor Richard Sutton, "The Bitter Lesson"

# The Third Pillar

Computers are used for large-scale physics / weather simulations. Large clusters of scientific computing machines.

**Scientific Computing**

Spearheaded by Amazon Web Services, the public cloud does a great job in web serving: moving data around, like webpages / images / videos.

**"Web Service Cloud"**

Modern AI applications call for exaflops computation power, over high performance, heterogeneous and cloud native infra.
This has never been seen in the history of cloud computing.

**Need for AI Cloud**

**1990s**

**2010s**

**1970s**

**2000s**

**2020s**

**Virtual Private Servers**

First application of managed machines - but still limited offering of software and applications on top of the machines.

**"Data Cloud"**

E-commerce has called for the processing power of exabyte scale data. Snowflake and Databricks emerged and grew into unicorns eventually.

# Why? AI is different from conventional compute

## Data Compute
- IO >> compute
- Simple abstraction
- Very Distributed Systems

Easy to use
Hard for infra

## Web Services
- IO > compute
- Arbitrary code
- Embarrassingly Parallel system

(Kinda) easy to use
(Kinda) easy for infra

## AI Compute
- Compute >> IO
- Arbitrary code
- Very Distributed Systems

(Pretty) hard to use
(Pretty) hard for infra

Lepton AI

# Conventional cloud value proposition no longer holds…

**SOFTWARE**

**SUPPLY CHAIN**

RETURN OF THE MPI

# Conventional cloud value proposition no longer holds…

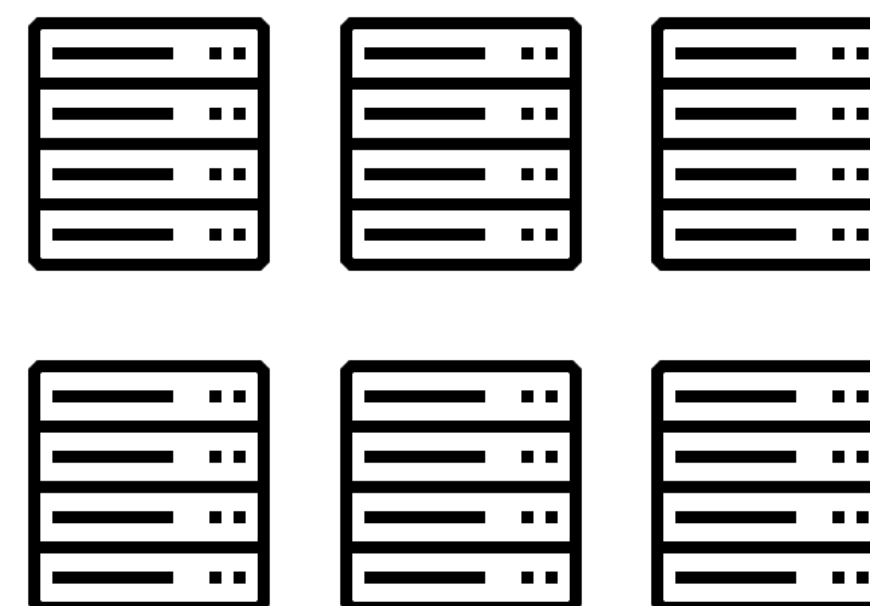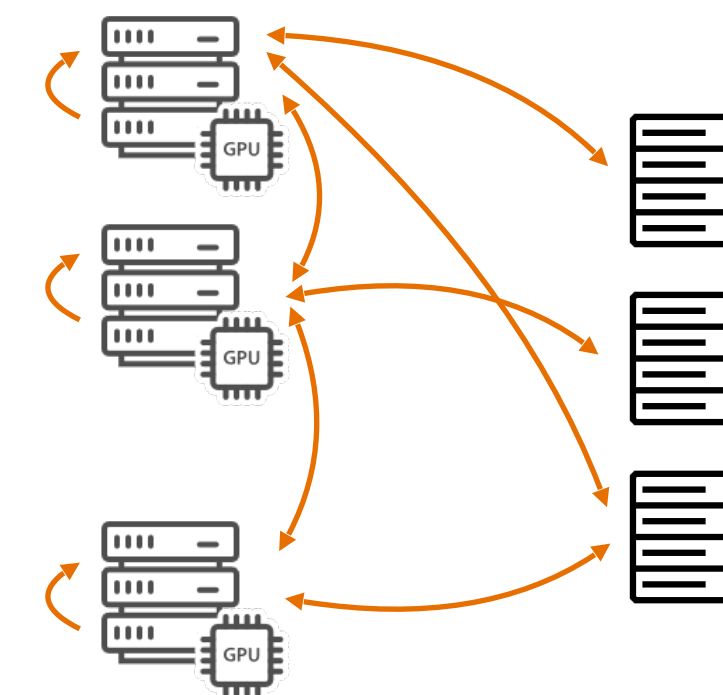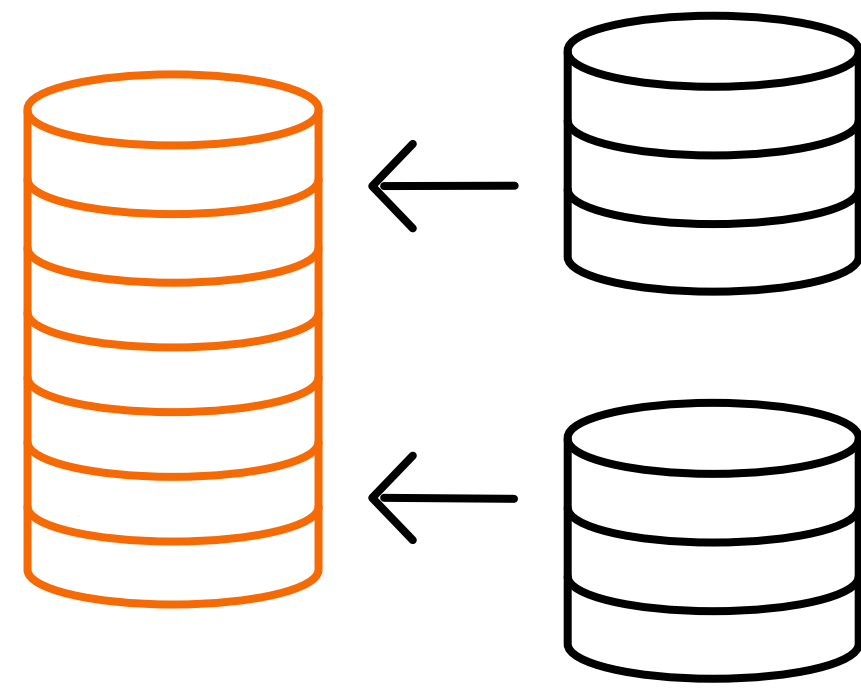| | Conventional Cloud | AI Cloud |
|---|---|---|
| **Software: Variety** | Complicated: Many applications Middleware | Simple: "AI frameworks" |
| **Software: Workload** | Varied: Compute, storage, network, big data, database, etc. | Unified: Numerical Computation |
| **Supply Chain: Flexibility** | High: CPU based Virtualization/Migration | Low: Large training Hard to live migrate |
| **Supply Chain: Interchangeability** | High: VMs can do many different jobs | Low: Really just doing AI compute |

# You shouldn't run on baremetal. Also, K8s is wrong.

# What you want to care about…

- Developer efficiency
  - Developer's time is precious
- Infra efficiency
  - GPUs do die
  - (More frequently than you think)

Root-cause of Interruptions During a 54-day Period of Llama 3 405B Pre-training

- NCCL Watchdog Timeouts 1.7%
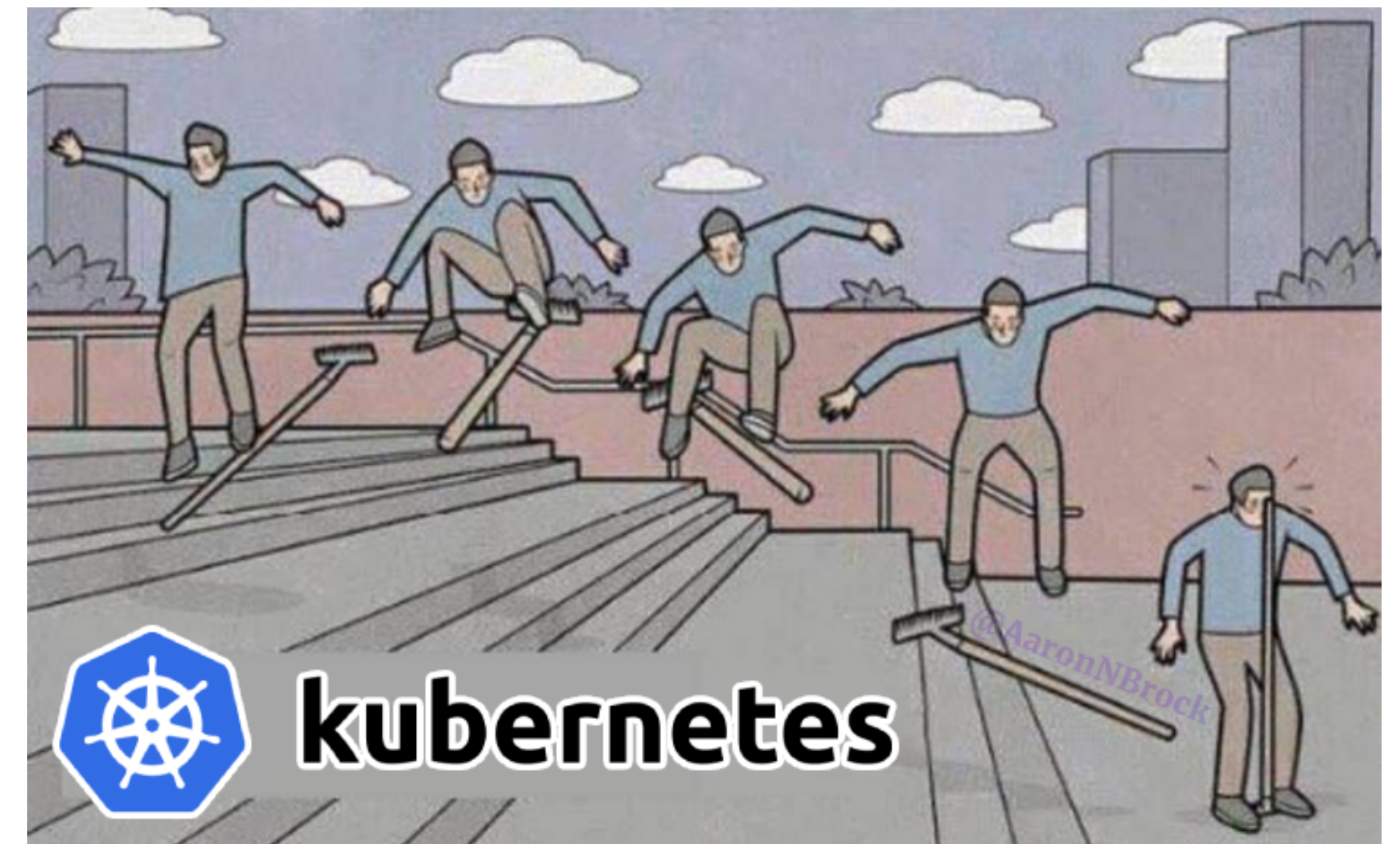- Network Switch/Cable 8.4%
- Host Maintenance Unplanned 7.6%
- System Memory 0.5%
- CPU 0.5%
- Power Supply 0.7%
- SSD 0.7%
- NIC 1.7%
- GPU Thermal Interface + Sensor 1.4%
- Silent Data Corruption 1.4%
- GPU System Processor 4.1%
- GPU SRAM Memory 4.5%
- GPU HBM3 Memory 17.2%
- Software Bug 12.9%
- Dependency 0.5%
- Faulty GPU 35.3%

# Best Practices?

- Multi-cloud supply chain management

- Elasticity and utilization management

- AI native platform to unify dev, training, and inference

- Build your own team around model and applications

| Serverless Cloud | Lepton API Services | Enterprise Deployment |
|---|---|---|

**Lepton AI Cloud Architecture**

| Deployments (Inference) | Jobs (Training) | Pods (Development) | Fast Runtimes (LLM, SD, etc) |
|---|---|---|---|

| Global Overlay Network | Infra Health Management | Lepton Optimized Kubernetes |
|---|---|---|

| Bare Metal & VM | High Throughput Storage | Cloud native middleware |
|---|---|---|

| Multi Cloud Hardware Sources |
|---|

AI infra is different, but also the same.

# 4
# HW and SW design: back to the future?

**Model Parameter Size (Log Scale)** vs. year (2005–2025)
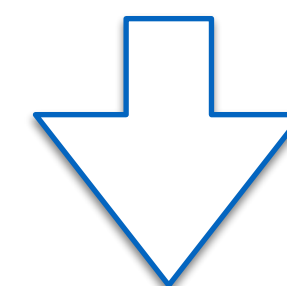
Y-axis labels: 10T, 1T, 100B, 10B, 1B, 100M, 10M, 1M, 100K, 10K, 1K, 100, 10

Annotated models: AlexNet, ResNet-V2, YOLOv2, GPT-2 (1.5B), GPT-3 175B, ChatGPT, GPT-4, Grok-1, DeepSeek-R1, Llama 2-70B, Llama 3.1-405B, Llama 4 Maverick, gpt-oss-120b

Data Source: Epoch (2025) – with major processing by Our World in Data
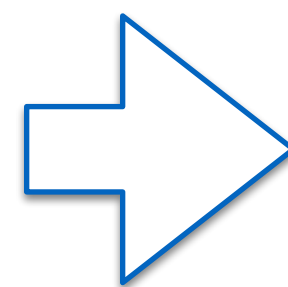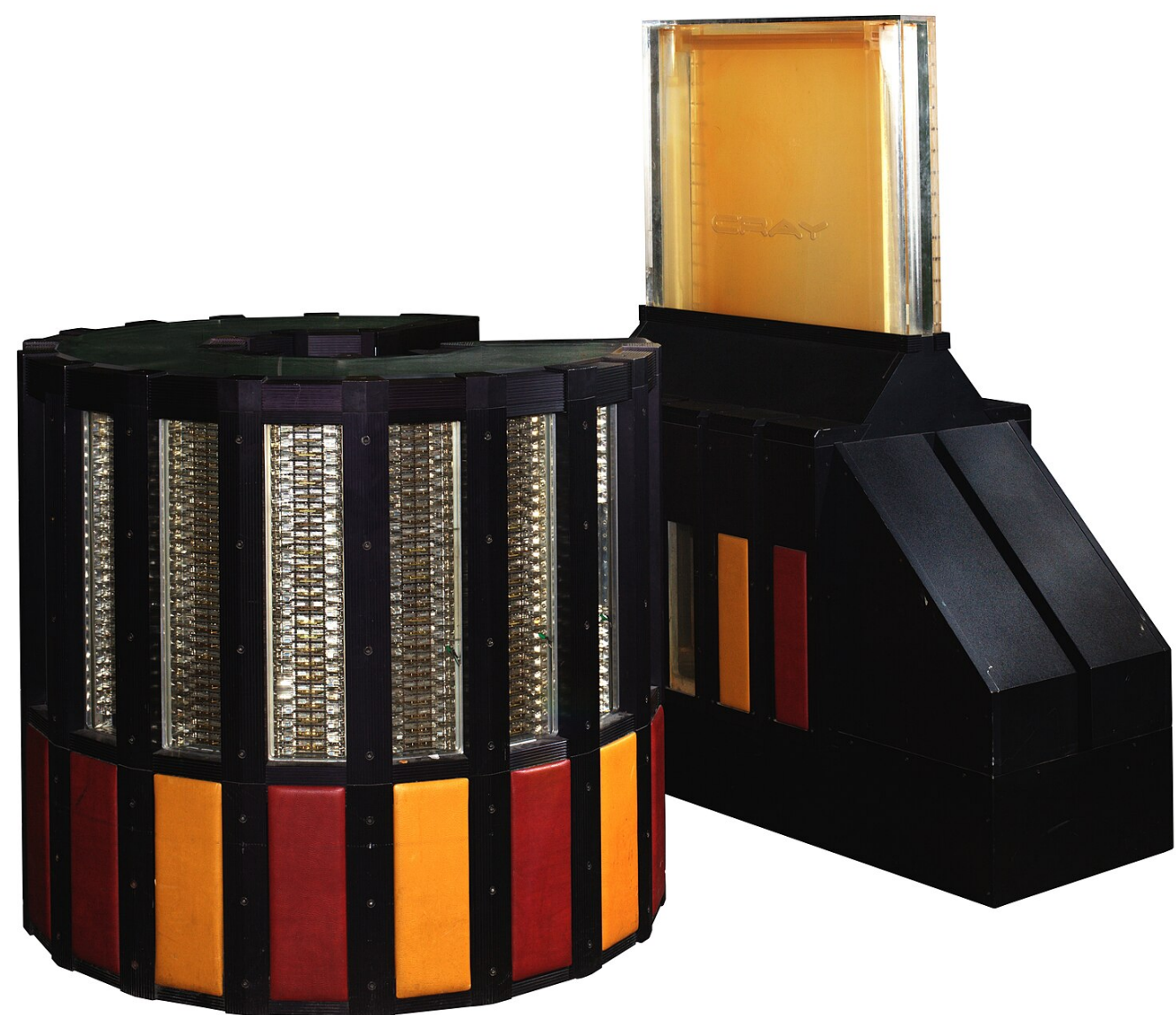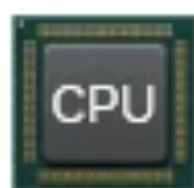
**Early AI Era**
≤100M parameters.
*Follows Moore's Law – Model size doubling every 20 months*
*Inference runs on CPUs*

**GPU AI Era**
~100M to ~1B parameters
*Doubling every 6 months*
*Inference runs on 1 GPU*

**Multi-GPU AI Era**
~1B to multi-trillion parameters
*Commercial AI Driven*
*Doubling every 10 months*
*Inference running up to 8 GPUs*

**Age of AI Reasoning at Scale**
*Drastic Increase in Compute for Reasoning*
*Expansion of Distributed Parallelism Techniques*
*Large Scale Mixture of Experts*
*Inference running up to 72 GPUs*

# Evolution of System Designs from an AI Engineer Perspective

**Yangqing Jia**

**Lepton AI (now part of NVIDIA)**