# Large Language Model Trends

Imanol Schlag

- **Introduction**

- **Past Trends**

- **Future Trends**

- **Swiss AI Initiative**
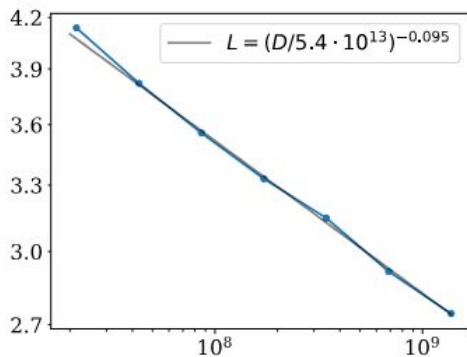
- **Q&A**

# Introduction

- Imanol Schlag

- Neural network research since 2016

- PhD at the Swiss AI Lab (IDSIA) under

  Jürgen Schmidhuber

- 20+ publications on ML/AI
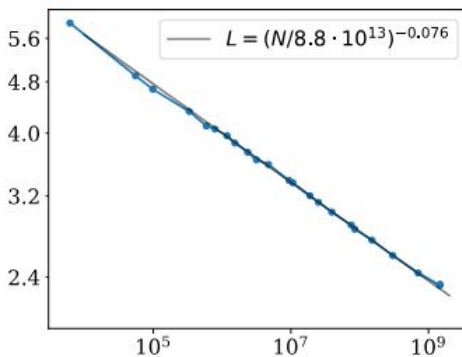
- Research Scientist at the ETH AI Center
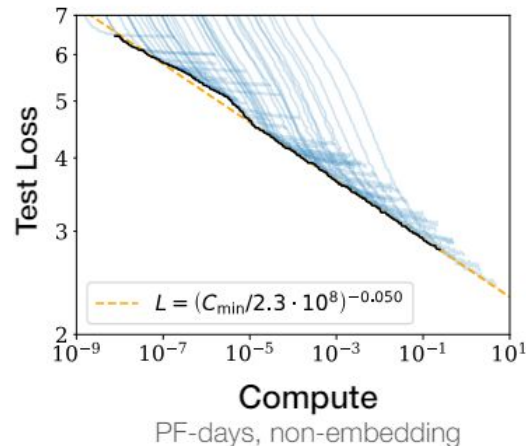


SOLA-Stafette 2024

# 2020: GPT 3

- GPT 2, February 2019

- 175B Parameters, May 2020

- In-context learning -> excel at NLP tasks



$L = (D/5.4 \cdot 10^{13})^{-0.095}$

$L = (N/8.8 \cdot 10^{13})^{-0.076}$

**Dataset Size**
tokens

**Parameters**
non-embedding



$L = (C_{min}/2.3 \cdot 10^{8})^{-0.050}$

**Compute**
PF-days, non-embedding
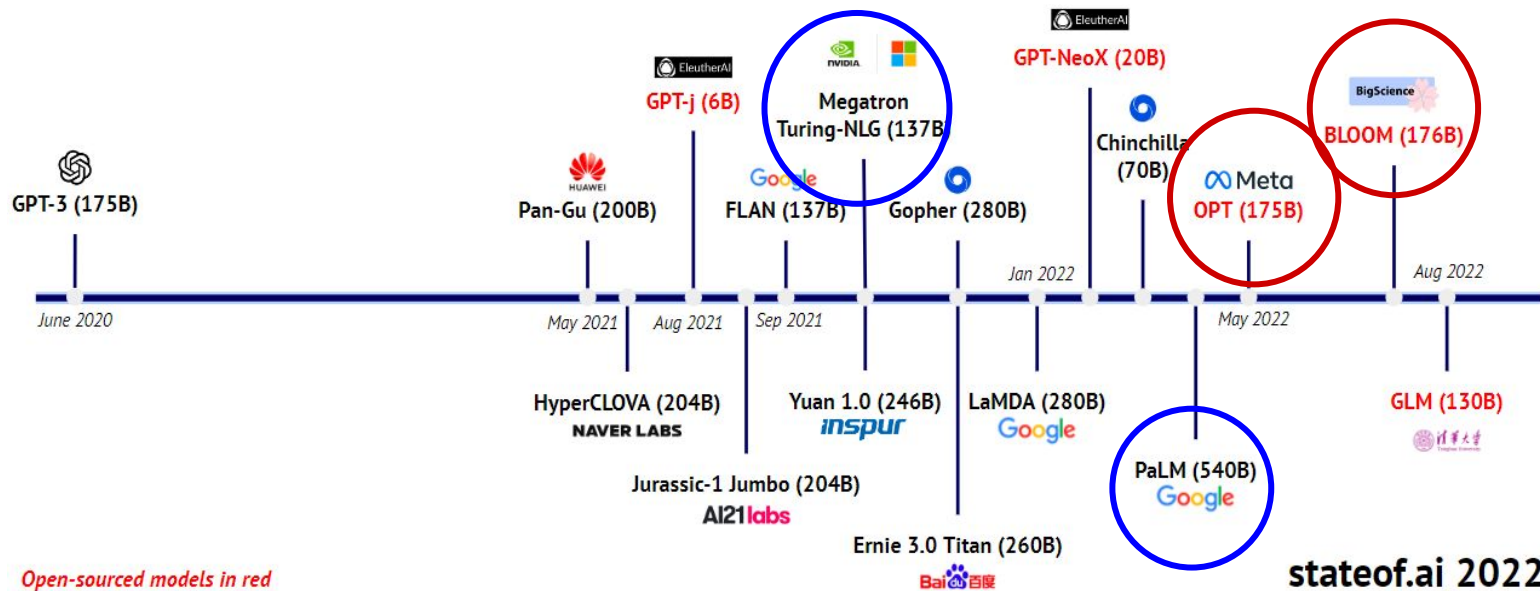
Kaplan et al. (January 2020)

# 2021: First Copycats and API-based Research

- Jurassic-1, Yuan 1.0, Ernie, Gopher

- Sparse methods at a Trillion parameters: GShard, GLaM

- Research explosion based on OpenAI API

  - Many new benchmarks

  - GPT-3 is not truthful / hallucinations

  - Prompting is brittle and challenging

  - Increased discussion around ethics & safety
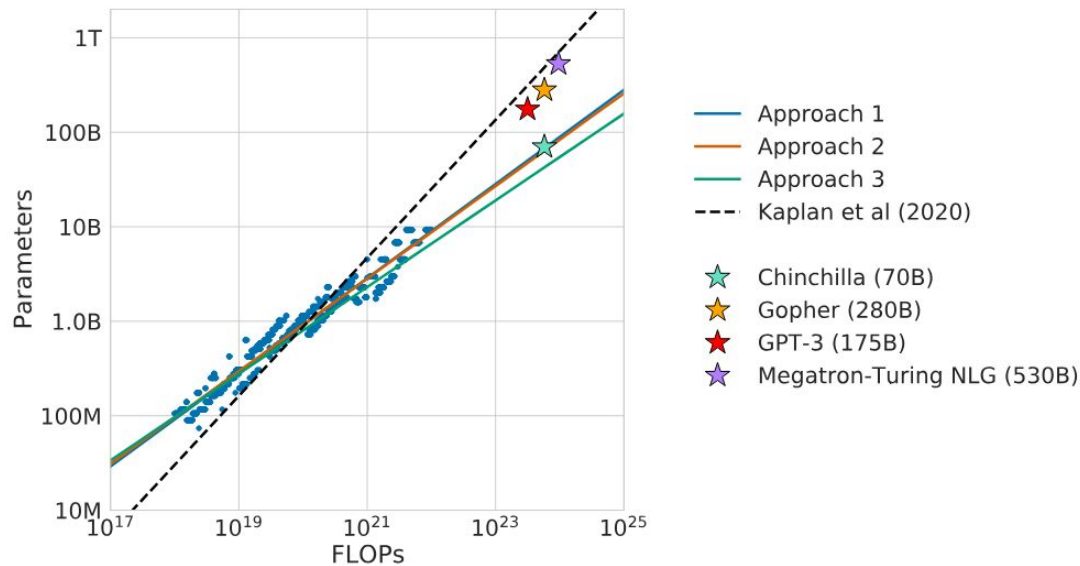
# 2022a: Bigger! Better?

500B scale models / first GPT-3 level open source copycats



Open-sourced models in red

stateof.ai 2022

ETH AI CENTER 6

# 2022b: Largest Models are Undertrained

March: Chinchilla & Scaling Laws



Hofmann et al. (2022)

# 2022c: Instruction Tuning and Alignment

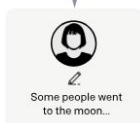- InstructGPT, Jan 2022

- RLHF, March 2022

# 2023a: Chatbots!

ChatGPT release: November 2022

## Road To 100 Million Users For Various Platforms

ChatGPT — (2 Months)

TikTok — (9 Months)

Youtube — (1.5 Years)

Instagram — (2.5 Years)

Facebook — (4.5 Years)

Twitter — (5 Years)

Spotify — (11 Years)

Netflix — (18 Years)

0 months    50 months    100 months    150 months    200 months    250 months

DEMANDSAGE

# 2023b: GPT-4

- Released in March

- Significant boost in performance veiled in secrecy

- Multimodal: image inputs + image generation through dall-e

- unofficial/leaked/rumoured details:

  - about 1.8T parameters, 120 layers, (x10 GPT-3)

  - 16 MLP-Experts each with ~111B parameters

  - 13T token training data

# 2023a: Chatbots!

- GPT-Turbo ChatGPT, March

- Claude 2, July

- Gemini, December

# 2023c: Competitive Open Weights LLMs

- Llama 1-65B, February

- MPT-30B, March

- Falcon-40B, June

- Llama 2-70B, July

- CodeLlama-70B, August

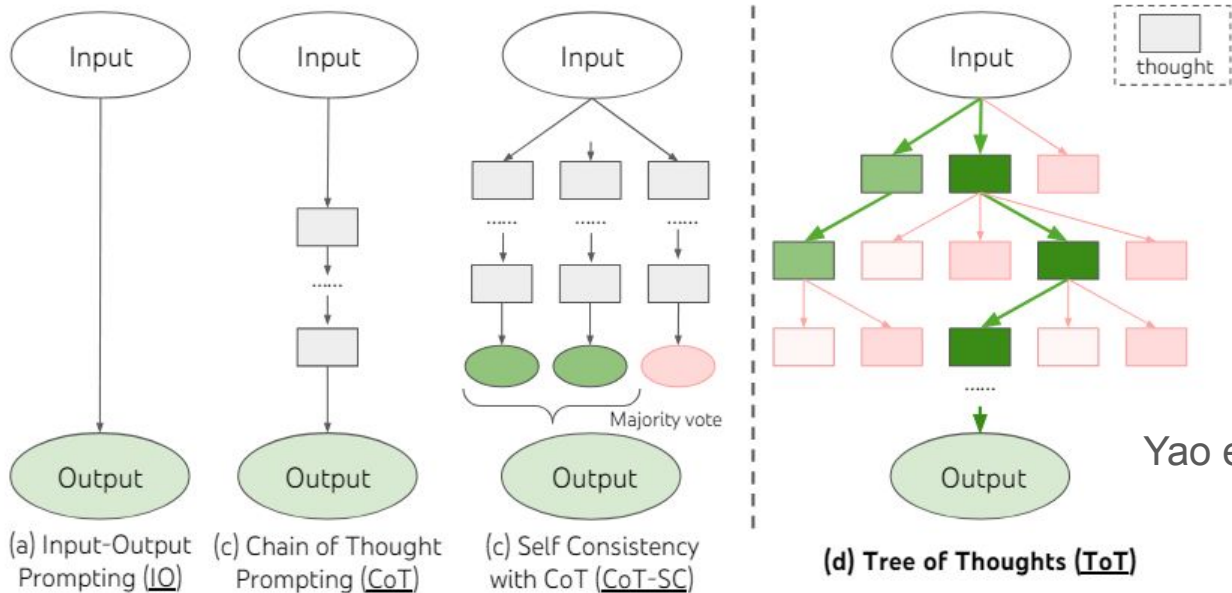|  |  | Humanities | STEM | Social Sciences | Other | Average |
|---|---|---|---|---|---|---|
| MPT | 7B | 26.7 | 25.3 | 27.1 | 28.2 | 26.8 |
|  | 30B | 44.5 | 39.0 | 52.8 | 52.9 | 46.9 |
| Falcon | 7B | 26.4 | 26.2 | 24.7 | 27.4 | 26.2 |
|  | 40B | 49.3 | 45.5 | 65.4 | 65.0 | 55.4 |
| LLAMA 1 | 7B | 34.0 | 30.5 | 38.3 | 38.1 | 35.1 |
|  | 13B | 45.0 | 35.8 | 53.8 | 53.3 | 46.9 |
|  | 33B | 55.8 | 46.0 | 66.7 | 63.4 | 57.8 |
|  | 65B | 61.8 | 51.7 | 72.9 | 67.4 | 63.4 |
| LLAMA 2 | 7B | 42.9 | 36.4 | 51.2 | 52.2 | 45.3 |
|  | 13B | 52.8 | 44.1 | 62.6 | 61.1 | 54.8 |
|  | 34B | 59.4 | 52.1 | 71.8 | 69.2 | 62.6 |
|  | 70B | **65.0** | **58.0** | **80.3** | **74.6** | **68.9** |

# 2023c: Competitive Open Weights LLMs

- Llama 1, February 2023

    ~5,500 citations

- Llama 2, July 2023

    ~4,300 citations

    7B: ~1.3M downloads

    70B: ~380k downloads

- 7,000 GitHub projects mentioning LLama

# 2023d: Explosion in LLM Projects/Research

- Agents using software tools

- Increasingly sophisticated prompting techniques, e.g. tree of thoughts
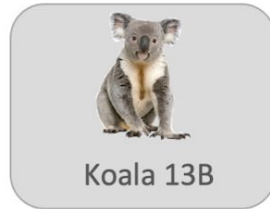


Yao et al. (2023)

# 2023d: Explosion in LLM Projects/Research

- Agents using software tools

- Increasingly sophisticated prompting techniques, e.g. tree of thoughts

- New alignment methods leads to many Llama derivatives:

Stanford Alpaca

Koala 13B

Vicuna

Orca

WizardLM

Llava

# 2024a: Next Generation General Purpose Assistants

- Multimodal (Voice, Video, Image)

- Long context (100k-1M tokens)

- Memory

- Strong coding and multilingual

- Tool use / execution environment

- Websearch / document upload



November 2023



February 2024



March 2024



yesterday



literally now?

ETH AI CENTER 16

# 2024b: Open Source/Weights Catching Up

- Grok-1, 314B, March

- DBRX, 132B, March

- Mixtral 8x22B, April

- Llama 3, 70B, April



Closed-source vs. Open-weight models (MMLU, 5-shot)

Labonne (2024)

# 2024b: Open Source/Weights Catching Up

- Grok-1, 314B, March

- DBRX, 132B, March

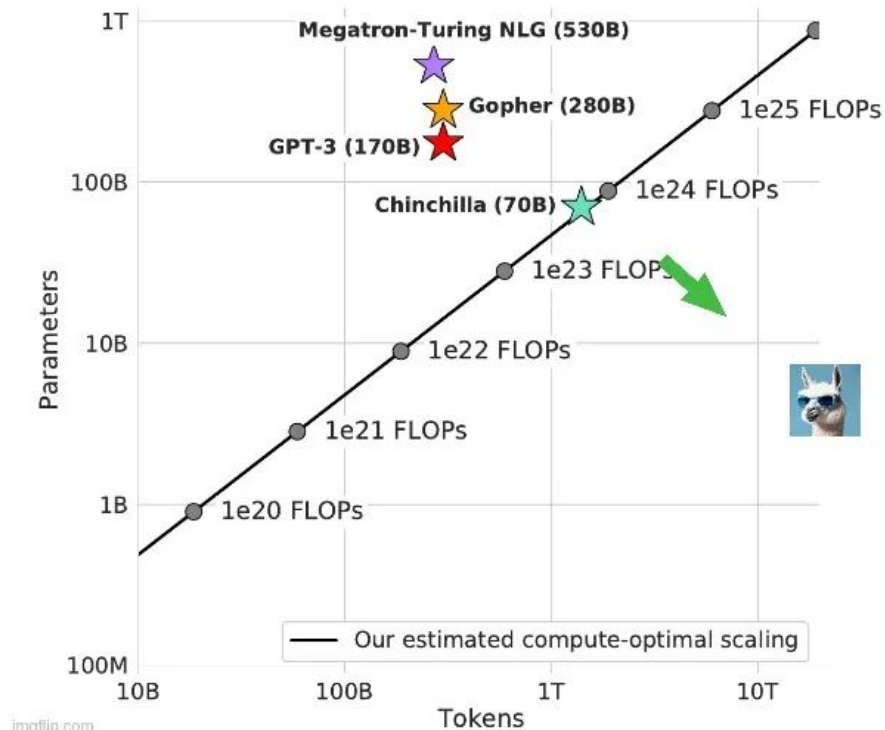- Mixtral 8x22B, April

- Llama 3, 70B, April

Inference cost matters!

1. 15T (!) tokens
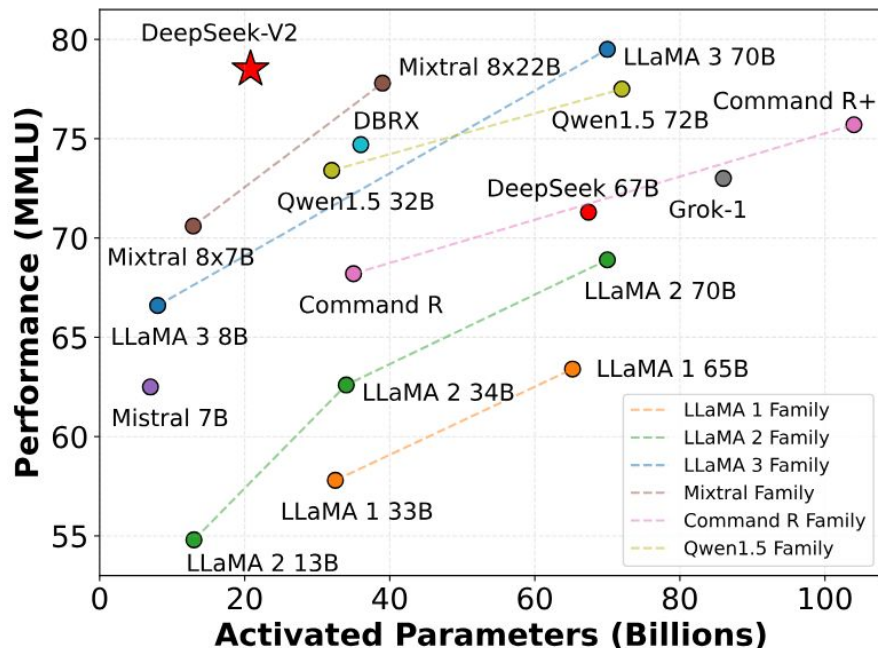2. Beats/competitive with Gemini Pro 1.5 and Claude 3 Sonnet

# 2024b: Open Source/Weights Catching Up

- Grok-1, 314B, March

- DBRX, 132B, March

- Mixtral 8x22B, April

- Llama 3, 70B, April

# 2024b: Open Source/Weights Catching Up
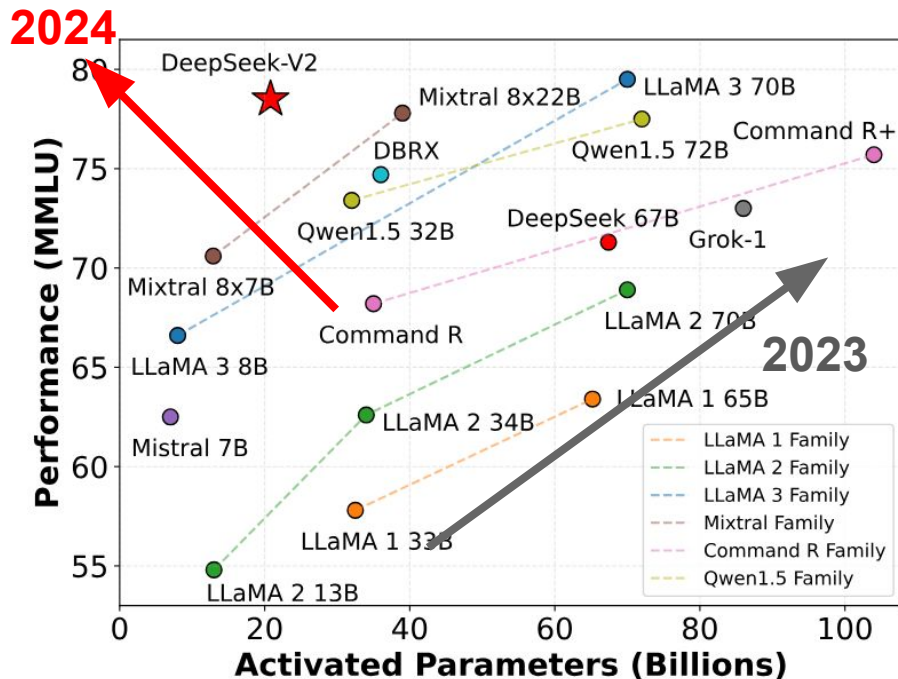
- Grok-1, 314B, March

- DBRX, 132B, March

- Mixtral 8x22B, April

- Llama 3, 70B, April

- DeepSeek-V2 67B, May
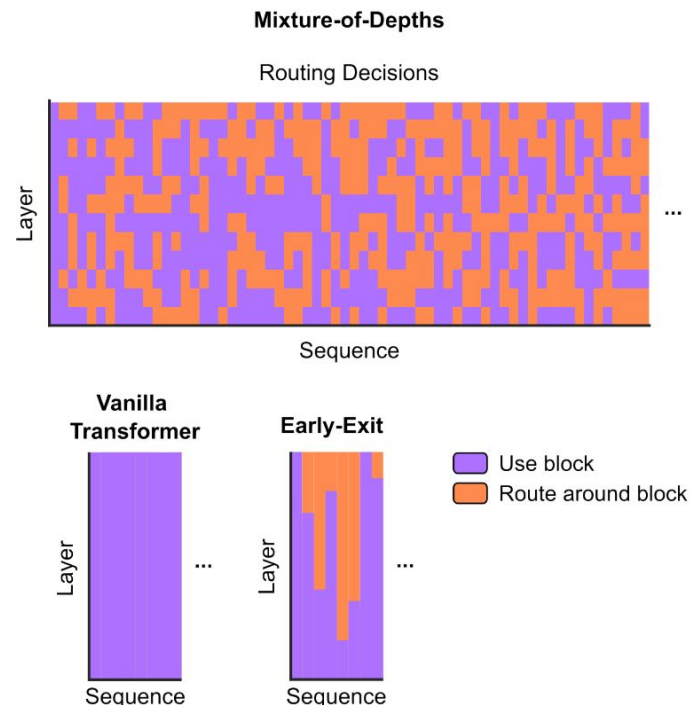


DeepSeek-AI (2024)

# 2024b: Open Source/Weights Catching Up!

- Grok-1, 314B, March

- DBRX, 132B, March

- Mixtral 8x22B, April

- Llama 3, 70B, April

- DeepSeek-V2 67B, May

- Llama 3 405B, June?



DeepSeek-AI (2024)

# Exciting Trend 1/6: MoE & Adaptive Computation

# Exciting Trend 2/6: The Return of Recurrent Nets

Associative RNN cell allows us to "scan" the sequence



(a) Up: parent combines values of its children

(b) Down: right child combines statistics of its parent with the left sibling.

Blelloch (1990)

# Exciting Trend 2/6: The Return of RNNs

Associative RNNs increasingly competitive with Transformers



(a) Scaling curve during training

(b) Maximum throughput at 1B parameter scale.

De et al. (2024)

# Exciting Trend 3/6: Low Precision & New Hardware

**Precision**

- fp16 (standard) -> fp8 training (ongoing) -> fp4 (soon)

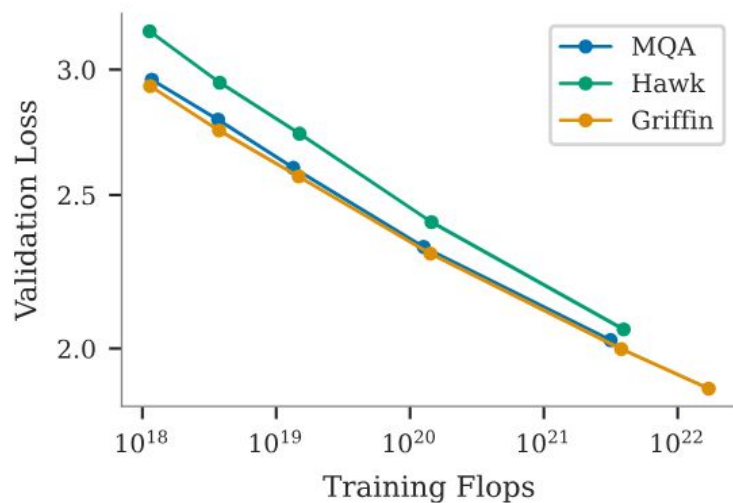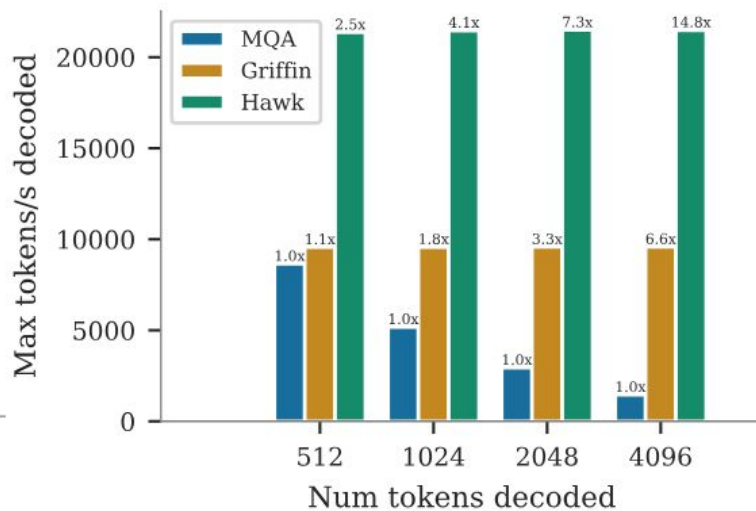- Inference with 4 bits, 2 bits, teneray (-1, 0, +1)

**Hardware**

- Specialised hardware for cheaper inference (groq.com)

- custom silicon (Apple, Meta, Google, Amazon, …)

- Nvidia Blackwell (FP4)

- Local LLMs -> embedded devices

**Transformer LLMs**

**16-bit Float (FP16/BF16)**

$$W = \begin{pmatrix} 0.2961 & -0.0495 & \dots & -0.4765 \\ 0.0413 & \dots & 0.2812 & 0.2403 \\ -0.1808 & 0.1304 & \dots & -0.1771 \\ -0.4809 & \dots & -0.1741 & -0.3853 \end{pmatrix}$$

**BitNet b1.58 (This Work)**

**{-1, 0, 1}**

$$W = \begin{pmatrix} 1 & -1 & \dots & 1 \\ 0 & \dots & -1 & -1 \\ -1 & 1 & \dots & 0 \\ -1 & \dots & 0 & -1 \end{pmatrix}$$

Ma et al. (2024)

# Exciting Trend 4/6: Quality Data & Synthetic Data

- High-quality public data for pretrain: FineWeb (15T tokens; webcrawl)



- Synthetic data for reasoning and alignment

1. Generate (E-step): The language model generates multiple output samples for each input context. Then, we filter these samples using a binary reward to collect the training dataset.
2. Improve (M-step): The original language model is supervised fine-tuned on the training dataset from the previous Generate step. The fine-tuned model is then used in the next Generate step.

Singh et al. (2024)

# Exciting Trend 5/6: Better Alignment

- RLHF: aligning LLMs via reinforcement learning with human feedback



- Alignment across modalities: images, videos, audio

# Exciting Trend 6/6: System Interactions & Agents

- Agents interacting / Tool use

- Web/Database Search & RAG: Retrieval Augmented Generation

- Coding & execution environment



Lilian Weng (2023)

# Swiss AI Initiative: swiss-ai.org

- National Research Initiative jointly lead by ETHZ and EPFL

- Scientific Council: 26 professors / researchers

- Assembly: >100 researchers

- 10M GPU hour commitment on Alps

# Swiss AI Initiative: swiss-ai.org



Alps Supercomputer: 10'000 GH200 GPUs

| Rank | System | Cores | Rmax (PFlop/s) | Rpeak (PFlop/s) | Power (kW) |
|------|--------|-------|----------------|-----------------|------------|
| 6 | **Alps** - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE Swiss National Supercomputing Centre (CSCS) Switzerland | 1,305,600 | 270.00 | 353.75 | 5,194 |

# Swiss AI Initiative: swiss-ai.org



| Foundation model for sustainability / climate | Foundation model for health | Foundation model for vision/robotics | Foundations model for education | Foundations models for sciences | Other areas |
|---|---|---|---|---|---|

**Conversational & alignment R&D**

| LLMs (focus on text) | Large multi-modal models | LLM security, red teaming & privacy | Other aspects of large models incl. "trustworthiness", learn new abilities from scale (of private data), ... |
|---|---|---|---|

**Infra & tools for scaling**

# Swiss AI Initiative: swiss-ai.org

**LLM Area:**

- An LLM for Switzerland
- Trustworthy and Responsible
- Transparent and compliant (open source / open weights)
- Multilingual with Swiss societal values

- Attract and develop talent
- Startup fuel
- Teaching and sharing lessons, code, models, …

- Collaborations: users, developers, legal, ….

# Questions?

Thank you for your attention.

Feel free to get in touch:

Imanol Schlag
ischlag@ethz.ch