

国内外の生成AIの開発状況

第16回 自動チューニング技術の現状と応用に関するシンポジウム

2024/12/26

国内の基盤モデルに関する主要な取り組み

SB Intuitions

- Sarashina-8x70B

ELYZA

- Llama-3-ELYZA-JP-70B

Preferred Elements

- PLaMo-100B

CyberAgent

- CyberAgentLM3-22B

GENIAC 第1期

- ELYZA, Kotoba Tech, 富士通, ABEJA, Sakana AI, NII, Stockmark, Turing, 松尾研, PFE

NII LLM-jp

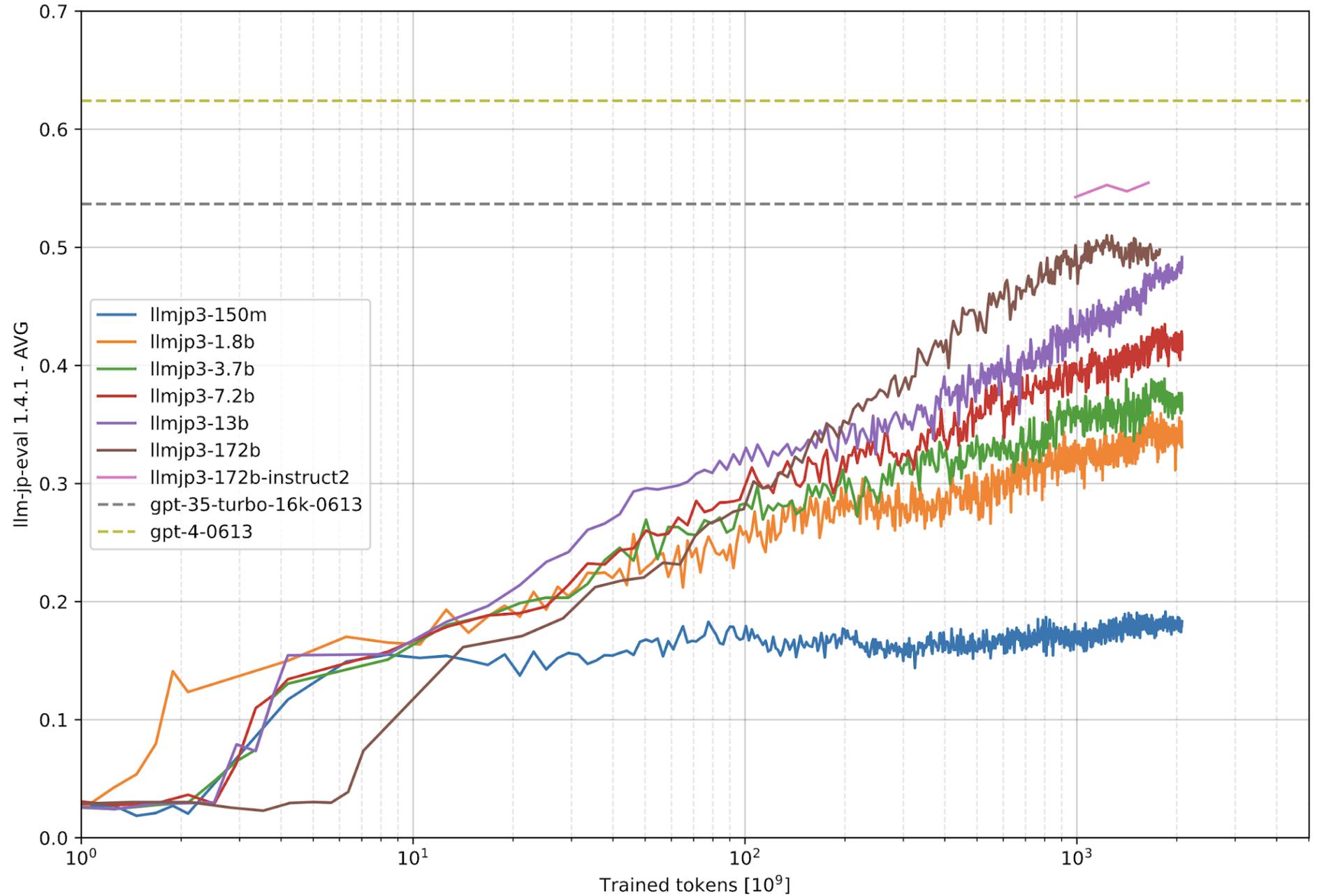
- llm-jp-172B

東京科学大学 + 産総研

- Llama-3.1-Swallow-70B

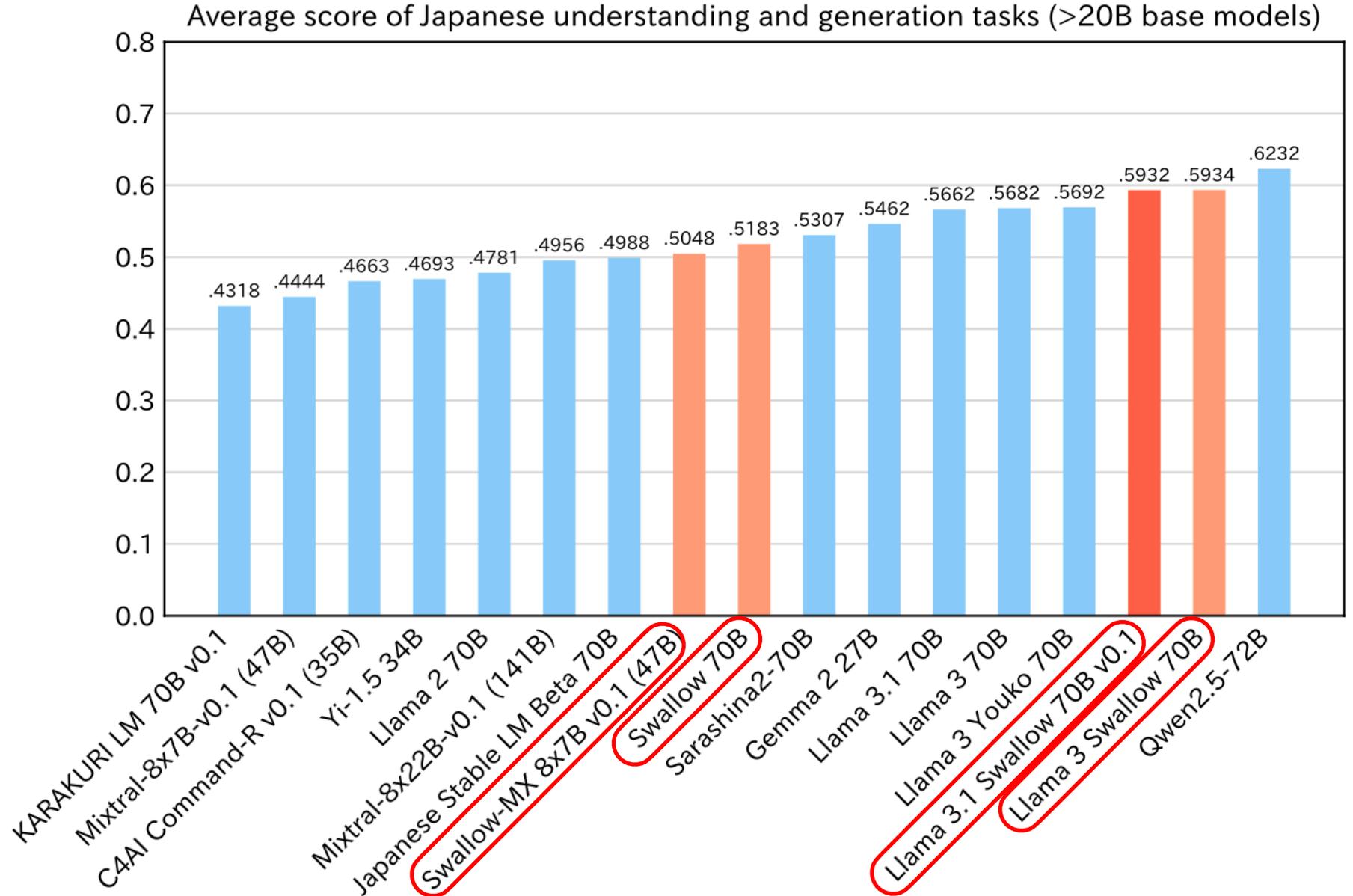
LLM-jp

- Largest effort to train an open LLM from scratch
- Model weights, dataset, checkpoints are open
- Constantly monitoring benchmark performance
- 172B benchmark scores are saturating
- Instruct model beats GPT-3.5



Swallow

- Continual pre-training from open models
- Create original Japanese dataset from CommonCrawl
- Currently #1 among LLMs trained in Japan
- Qwen2.5 is a little better than Swallow



国外の基盤モデルに関する主要な取り組み

OpenAI

- ・ o3

Google

- ・ Gemini 2.0 Flash Thinking

X.ai

- ・ Grok 2

Anthropic

- ・ Claude 3.5 Sonnet

Meta

- ・ Llama 3.3

Mistral AI

- ・ Mistral Large 2

Alibaba

- ・ QwQ

NVIDIA

- ・ Nemotron 4

AWS

- ・ Nova Premier

Microsoft Azure

- ・ Phi 3

TTI

- ・ Falcon 3

Cohere

- ・ Command R+

DeepSeek

- ・ DeepSeek R1

0.1-AI

- ・ Yi Large

DataBricks

- ・ DBRX

Perplexity

- ・ Sonar 3.1 Large

Reka

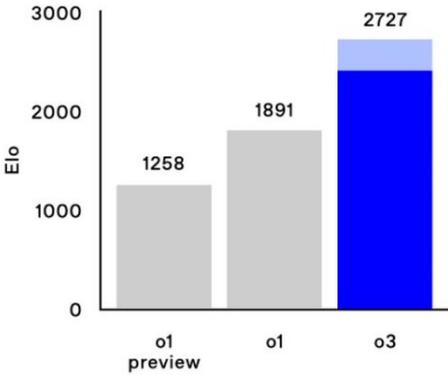
- ・ Reka Flash

AI21Labs

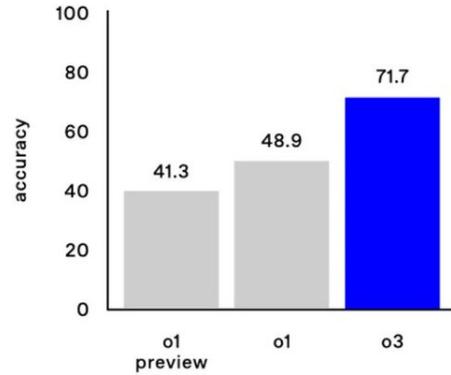
- ・ Jamba 1.5 Large

OpenAI o3

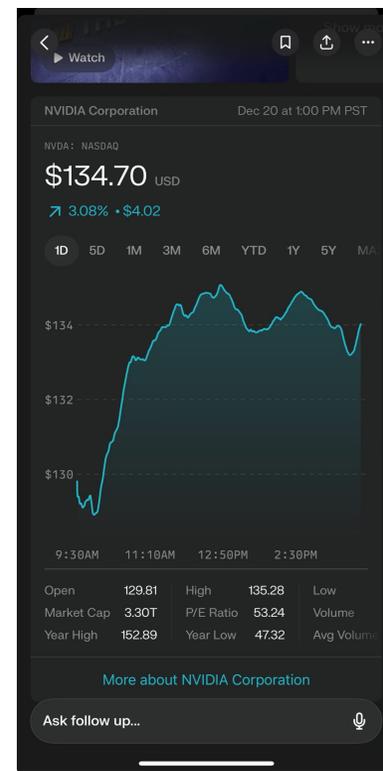
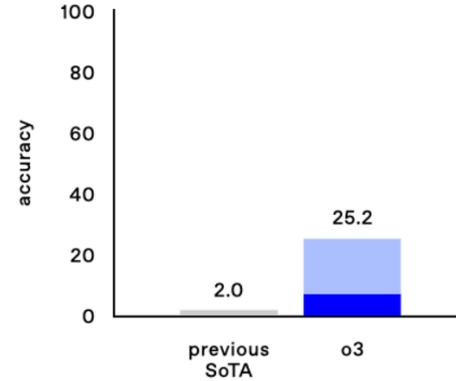
Competition Code (Codeforces)



Software Engineering (SWE-bench Verified)

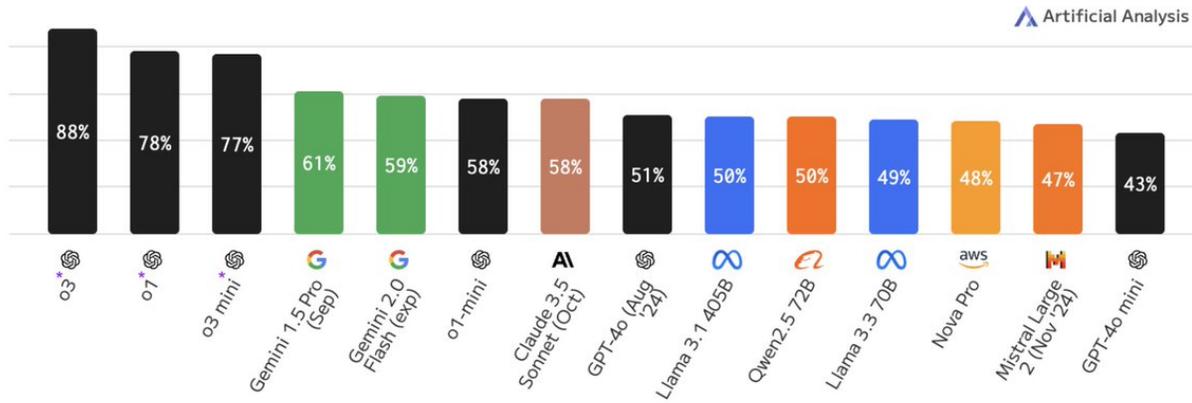


Research Math (EpochAI Frontier Math)

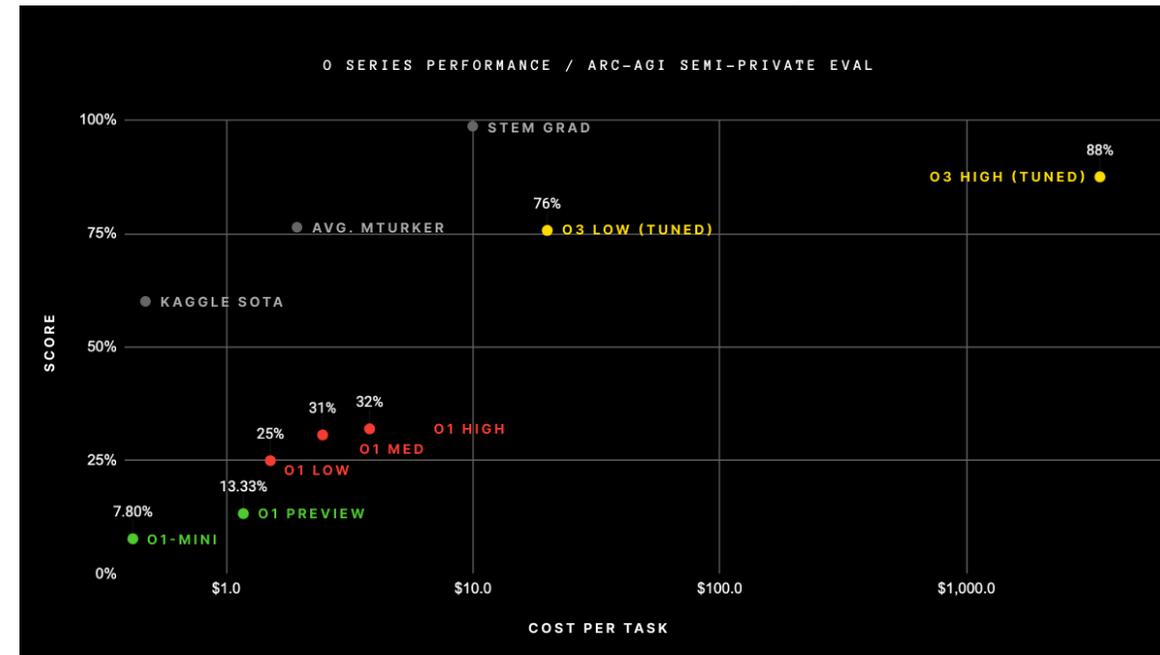


GPQA Diamond (Scientific Reasoning & Knowledge)

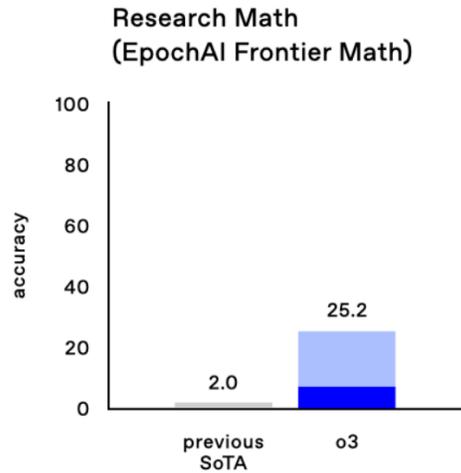
GPQA Diamond scores independently evaluated by Artificial Analysis; Higher is better



* Results claimed by OpenAI, not yet independently benchmarked by Artificial Analysis



Frontier Math



Prime field continuous extensions

Problem Solution

Let a_n for $n \in \mathbb{Z}$ be the sequence of integers satisfying the recurrence formula

$$a_n = 198130309625a_{n-1} + 354973292077a_{n-2} - 427761277677a_{n-3} + 370639957a_{n-4}$$

with initial conditions $a_i = i$ for $0 \leq i \leq 3$. Find the smallest prime $p \equiv 4 \pmod{7}$ for which the function $\mathbb{Z} \rightarrow \mathbb{Z}$ given by $n \mapsto a_n$ can be extended to a continuous function on \mathbb{Z}_p .

Find the degree 19 polynomial

Problem Solution

Construct a degree 19 polynomial $p(x) \in \mathbb{C}[x]$ such that $X := \{p(x) = p(y)\} \subset \mathbb{P}^1 \times \mathbb{P}^1$ has at least 3 (but not all linear) irreducible components over \mathbb{C} . Choose $p(x)$ to be odd, monic, have real coefficients and linear coefficient -19 and calculate $p(19)$.

ARC-AGI

ARC-AGI特化モデル

2024 HIGH SCORE WINNERS

Rank	Team	Code	Paper	Score	Prize
1st	the ARCHitects	CODE	PAPER	53.5%	\$25k
2nd	G. Barbadillo	CODE	PAPER	40%	\$10k
3rd	alijs	CODE		40%	\$5k
4th	William Wu	CODE		37%	\$5k
5th	PooHAI	CODE	PAPER	37%	\$5k

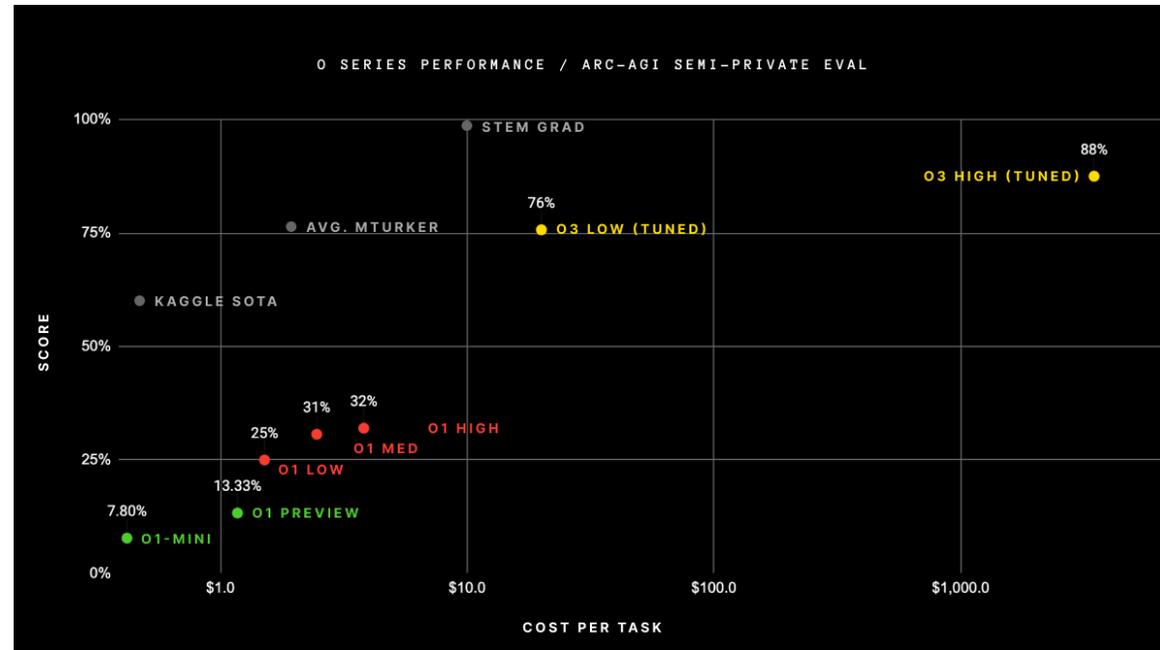
private

汎用モデル

2024 ARC-AGI-PUB HIGH SCORES

Model	Score	Prize	Code
o3 (coming soon)	75.7%	\$2.6k	
o1-preview	18%	21%	CODE
Claude 3.5 Sonnet	14%	21%	CODE
GPT-4o	5%	9%	CODE
Gemini 1.5	4.5%	8%	CODE

semi-private public



o1-like models from China

O1-CODER: AN O1 REPLICATION FOR CODING

Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiangming Shu, Jinlin Xiao, Chao Kong & Jitao Sang*
School of Computer Science and Technology
Beijing Jiaotong University
Beijing, China
{yuxiangzhang, wushangxi, yqyang, jiangmingshu, jinlinx, 23120361, jtsang}@bjtu.edu.cn

 Generative AI Research

O1 Replication Journey: A Strategic Progress Report – Part 1

Yiwei Qin^{1,4*} Xuefeng Li^{1,4*} Haoyang Zou^{4*} Yixiu Liu^{1,4*} Shijie Xia^{1,4*}
Zhen Huang⁴ Yixin Ye^{1,4} Weizhe Yuan² Hector Liu³ Yuanzhi Li³ Pengfei Liu^{1,4†}
¹Shanghai Jiao Tong University, ²New York University,
³MBZUAI, ⁴Generative AI Research Lab (GAIR)

Alibaba
International

2024-11-26

Marco-o1: Towards Open Reasoning Models for Open-Ended Solutions

Yu Zhao*, Huifeng Yin*, Bo Zeng, Hao Wang, Tianqi Shi, Chenyang Lyu, Longyue Wang, Weihua Luo, Kaifu Zhang

MarcoPolo Team, Alibaba International Digital Commerce

 Generative AI Research

O1 Replication Journey – Part 2: Surpassing O1-preview through Simple Distillation Big Progress or Bitter Lesson?

Zhen Huang^{4*} Haoyang Zou^{4*} Xuefeng Li^{1,4*} Yixiu Liu^{1,4*} Yuxiang Zheng^{1,4*}
Ethan Chern^{1,4*} Shijie Xia^{1,2,4*} Yiwei Qin⁴ Weizhe Yuan³ Pengfei Liu^{1,2,4†}
¹Shanghai Jiao Tong University, ²SII, ³NYU,
⁴Generative AI Research Lab (GAIR)

Benchmark	QwQ 32B-preview	OpenAI o1-preview	OpenAI o1-mini	Claude3.5 Sonnet	Qwen2.5-72B Instruct	GPT-4o
GPQA	65.2	72.3	60.0	65.0	49.0	53.6
AIME	50.0	44.6	56.7	16.0	23.3	9.3
MATH-500	90.6	85.5	90.0	78.3	82.6	76.6
LiveCodeBench	50.0	53.6	58.0	36.3	30.4	33.4

AI for AI を AI for Scienceにつなげる

Trillion Parameter Consortiumの大幅な方針転換

- OpenAIやAnthropicと協議し、AI for Scienceの大幅な方針転換を決めた
- Domain specific foundation model → Domain specific benchmarks
- 方針転換後はOpenAIが協力的でモデルを共有してくれる予定

古典的な深層学習の制約を排除

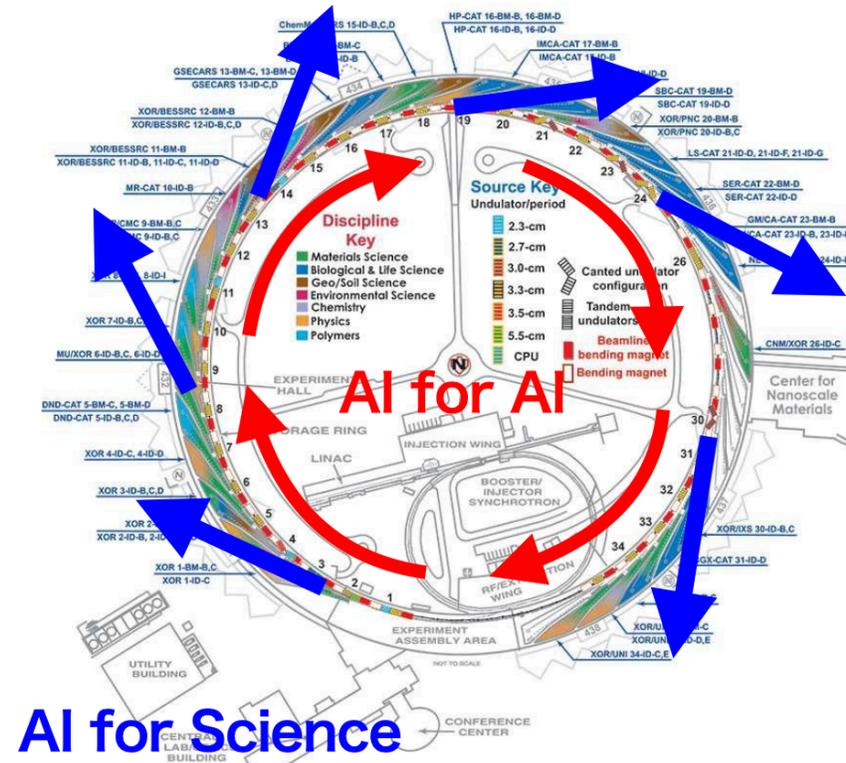
- 分類・要約 → 生成 → **インタラクション**

モデルにデータを入れるのが最善の方法か？

- 事前学習 → 事後学習 → **コンテキスト内学習**

データがAIに入っていくのではなく、AIがデータに入っていく

- AIを賢くするためのデータ vs それを使って分析したいデータ
- 前者はBigTechと競争してもリソースの無駄
- 後者を探索するためのAI → 知的な検索・探索



AMD GPUにおけるLLMの性能

- MI300Xは理論性能はH200より高い
- LLMの学習の性能は遥かに低い
- モデルによってはさらに低い
- この性能を出すだけでも大変

H100 vs H200 vs MI300X Basic Specifications			
	H100	H200	MI300X
Watts Per GPU (TDP)	700	700	750
All-in System Watts Per GPU	1,275	1,275	1,275
Memory Capacity (GB)	80GB	141GB	192GB
Memory Bandwidth (GB/s)	3,352	4,800	5,300
FP16/BF16 TFLOPS ¹	989	989	1,307
FP8 / FP6 / Int8 TFLOPS ¹	1,979	1,979	2,615

1. All FLOPS are dense

