# The Art of Artificial Reasoning for (Small) Language Models

## Yejin Choi

### Stanford & NVIDIA 💚

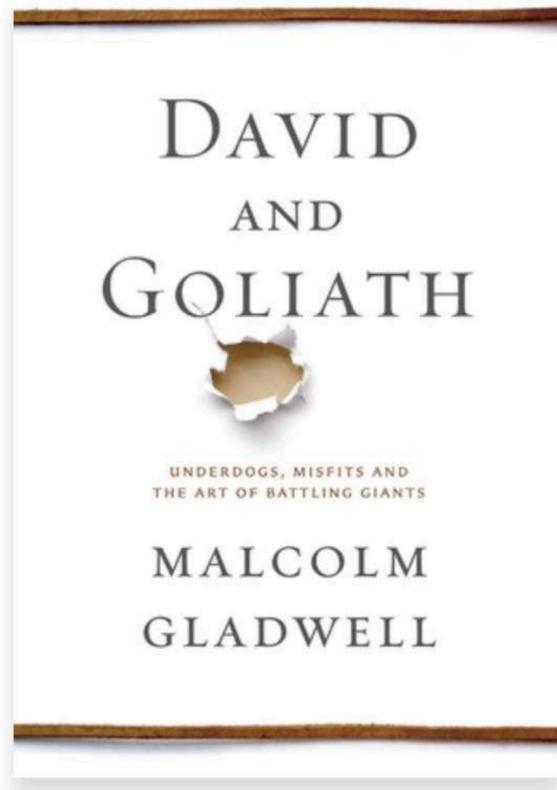# NVIDIA LEADS OPEN SOURCE AI MOMENTUM AS CHINESE LABS CLOSE IN



October 23, 2024

*Considering repos above 10 space likes and 500 downloads*

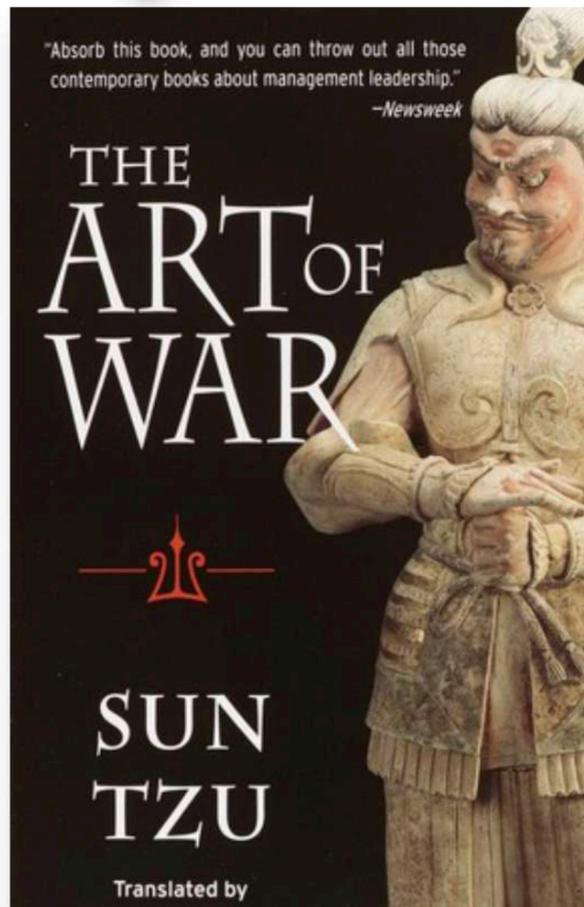# NVIDIA LEADS OPEN SOURCE AI MOMENTUM AS CHINESE LABS CLOSE IN



1. Small models are in high demand, with daily downloads hitting tens of millions (!!! 🔥) for top repositories
2. China's contribution to the open source ecosystem has skyrocketed, with Alibaba now surpassing Hugging Face 😮
3. NVIDIA has emerged as a rising star ⭐ (or a risen star? 😎), for open-source contributions 💚

**David and Goliath: Underdogs, Misfits, and the Art of Battling Giants**

by Malcolm Gladwell

**The Art of War**

by Sun Tzu, Thomas Cleary, Pulat Otkan (Translator)

*David vs. Goliath:*
*the Art of Scaling Intelligence*
*in the Era of Extreme-Scale Neural Models*

Three components to innovate:
Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏

# The Era of Brute-Force Scaling is Over
## The Era of Smart Scaling Begins



Pre-training as we know it will end

Compute is growing:

- Better hardware
- Better algorithms
- Larger clusters

Data is not growing:

- We have but one internet
- **The fossil fuel of AI**

-- Ilya Sutskever from his test time award talk at NeurIPS 2024

# The Era of Brute-Force Scaling is Over

## The Era of Smart Scaling Begins



THE NEW YORKER 100

Newsletter    My Acc

Fiction & Poetry    Humor & Cartoons    Magazine    Puzzles & Games    Video    Podcasts

OPEN QUESTIONS

WHAT IF A.I. DOESN'T GET MUCH BETTER THAN THIS?

*GPT-5, a new release from OpenAI, is the latest product to suggest that progress on large language models has stalled.*

By Cal Newport

August 12, 2025



Pre-training as we know it will end

Compute is growing:
- Better hardware
- Better algorithms
- Larger clusters

Data is not growing:
- We have but one internet
- **The fossil fuel of AI**

## NewScientist

Sign in    Enter search keywords

News    Features    Newsletters    Podcasts    Video    Comment    Culture    Games    |    This week's magazine

Health    Space    Physics    Technology    Environment    Mind    Humans    Life    Mathematics    Chemistry    Earth    Society

**Analysis** and **Technology**

# GPT-5's modest gains suggest AI progress is slowing down

OpenAI's latest large language model has achieved seemingly underwhelming improvements in performance, leading to questions about whether the AI industry can make significant advancements with its current designs

By Alex Wilkins

13 August 2025

# **The Era of Brute-Force Scaling is Over**

The Era of Smart Scaling Begins

Three broad ways to cope with data saturation:

1. Learn better & faster with limited data
   - Alternative architectures
   - Alternative training recipes
2. Synthesize new data
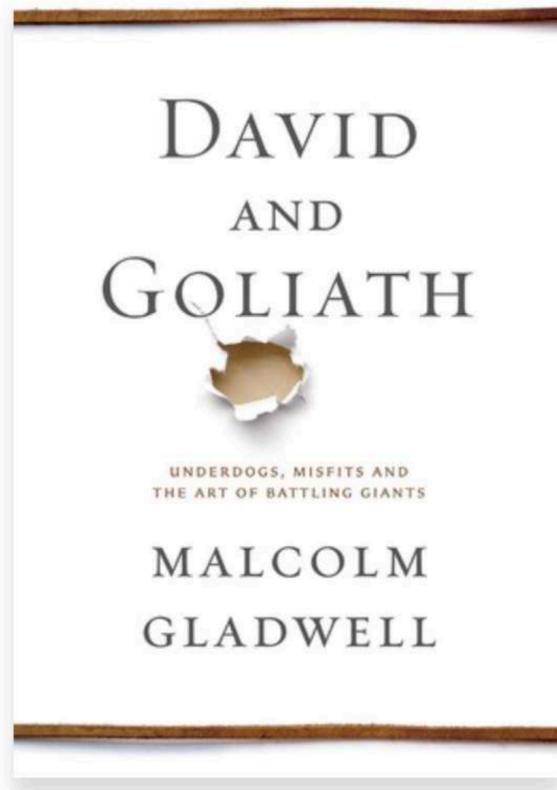   - Algorithms to generate the outer space of the internet data
3. Reason beyond what is in the data
   - Test-time reasoning algorithms
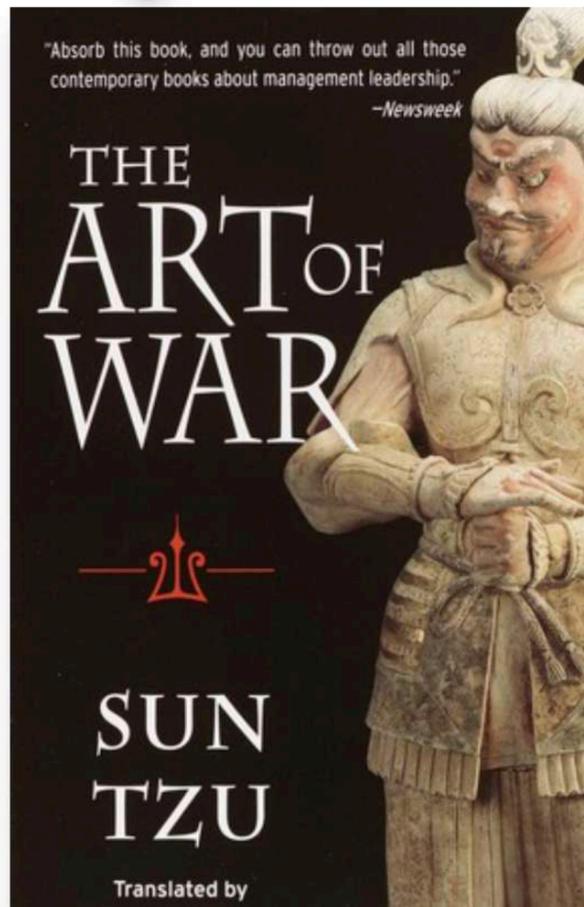   - Test-time training algorithms

=> Scaling of "intelligence" continues

**David and Goliath: Underdogs, Misfits, and the Art of Battling Giants**

by Malcolm Gladwell

**The Art of War**

by Sun Tzu, Thomas Cleary, Pulat Otkan (Translator)

*David vs. Goliath:*
*the Art of Scaling Intelligence*
*in the Era of Extreme-Scale Neural Models*

Three components to innovate:
Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏

# David vs. Goliath:

## the Art of Scaling Intelligence

### in the Era of Extreme-Scale Neural Models

h: Underdogs,
Art of Battling

In this talk:

===

ProRL: Prolonged RL
Prismatic Synthesis
RL as Pretraining

===

"Smaller but Better"
"Algorithms for the Win"

Three components to innovate:
Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏

# ProRL

Prolonged Reinforcement Learning Expands
Reasoning Boundaries in Large Language Models

**Mingjie Liu**   **Shizhe Diao**   **Ximing Lu**   **Jian Hu**   **Xin Dong**
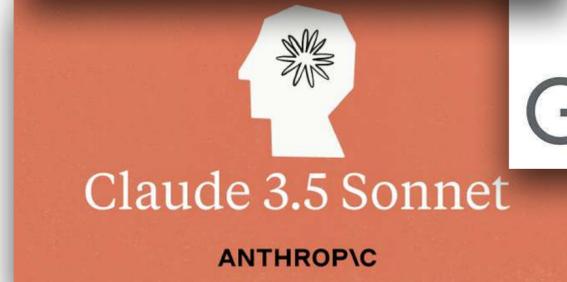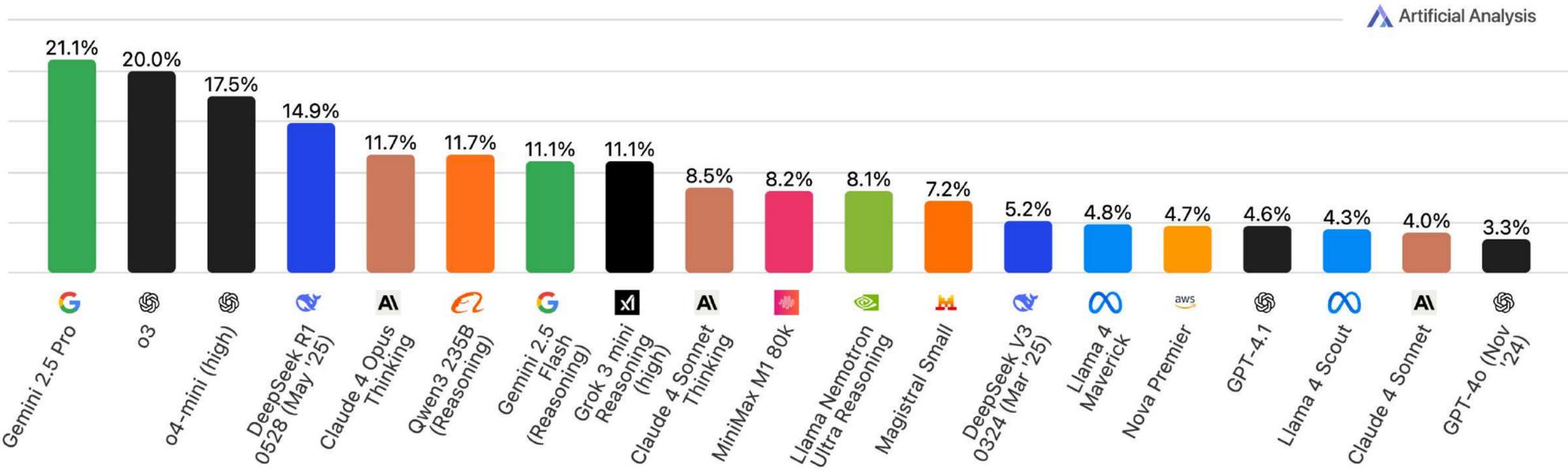**Yejin Choi**   **Jan Kautz**   **Yi Dong**

NVIDIA

{mingjiel, sdiao, ximingl, jianh, xind, yejinc, jkautz, yidong}@nvidia.com

# 2025: The Rise of L**R**Ms (as opposed to L**L**Ms)

## Humanity's Last Exam Benchmark Leaderboard: Results
Independently conducted by Artificial Analysis

Artificial Analysis

| Model | Score |
|---|---|
| Gemini 2.5 Pro | 21.1% |
| o3 | 20.0% |
| o4-mini (high) | 17.5% |
| DeepSeek R1 0528 (May '25) | 14.9% |
| Claude 4 Opus Thinking | 11.7% |
| Qwen3 235B (Reasoning) | 11.7% |
| Gemini 2.5 Flash (Reasoning) | 11.1% |
| Grok 3 mini Reasoning (high) | 11.1% |
| Claude 4 Sonnet Thinking | 8.5% |
| MiniMax M1 80k | 8.2% |
| Llama Nemotron Ultra Reasoning | 8.1% |
| Magistral Small | 7.2% |
| DeepSeek V3 0324 (Mar '25) | 5.2% |
| Llama 4 Maverick | 4.8% |
| Nova Premier | 4.7% |
| GPT-4.1 | 4.6% |
| Llama 4 Scout | 4.3% |
| Claude 4 Sonnet | 4.0% |
| GPT-4o (Nov '24) | 3.3% |

- **Long thought** (chain-of-thought)
- The power of **Reinforcement Learning**
- Shift from **Imitation Learning** to **Exploration Learning**

# 2025: The Rise of LRMs (as opposed to LLMs)

## Does Math Reasoning Improve General LLM Capabilities? Understanding Transferability of LLM Reasoning

Maggie Huan[1,2,*], Yuetai Li[3,*], Tuney Zheng[4,*], Xiaoyu Xu[5], Seungone Kim[1], Minxin Du[5], Radha Poovendran[3], Graham Neubig[1], Xiang Yue[1,†]
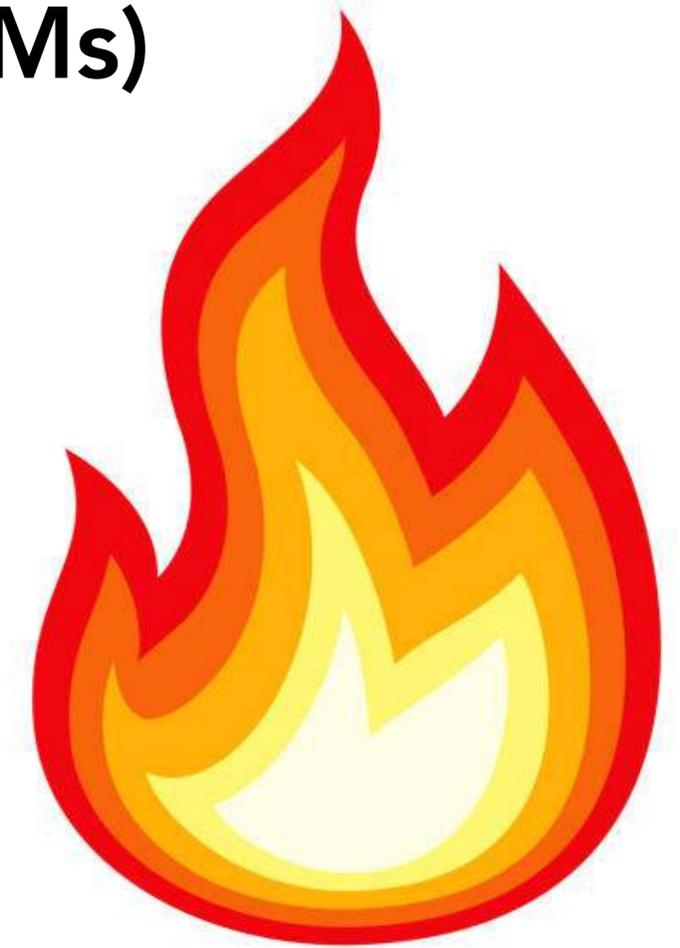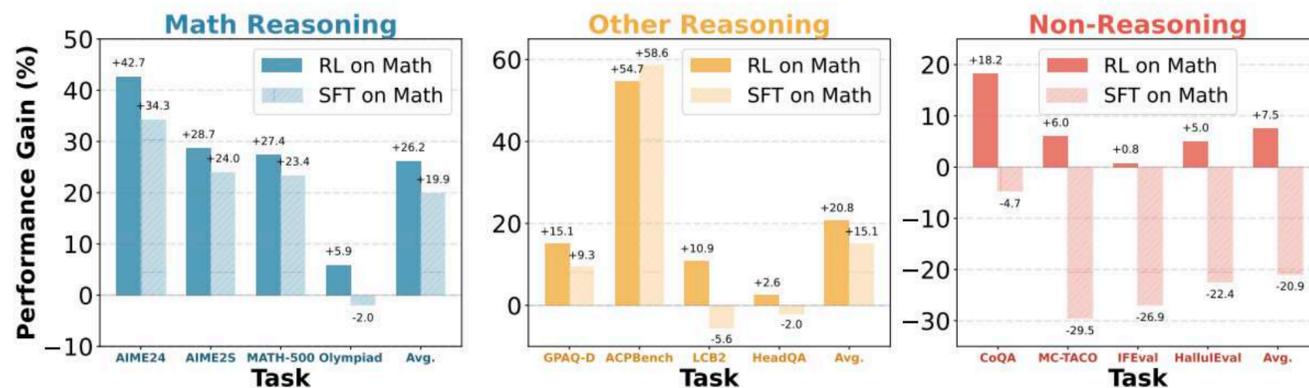
[1]Carnegie Mellon University  [2]University of Pennsylvania  [3]University of Washington
[4]M-A-P  [5]The Hong Kong Polytechnic University

ziyuh@seas.upenn.edu   yuetaili@uw.edu   xyue2@andrew.cmu.edu

**Abstract:** Math reasoning has become the poster child of progress in large language models (LLMs), with new models rapidly surpassing human-level performance on benchmarks like MATH and AIME. But as math leaderboards improve week by week, it is worth asking: *do these gains reflect broader problem-solving ability or just narrow overfitting?* To answer this question, we evaluate over 20 open-weight reasoning-tuned models across a broad suite of tasks, including math, scientific QA, agent planning, coding, and standard instruction-following. We surprisingly find that most models that succeed in math fail to transfer their gains to other domains. To rigorously study this phenomenon, we conduct controlled experiments on Qwen3-14B models using math-only data but different tuning methods. We find that reinforcement learning (RL)-tuned models generalize well across domains, while supervised fine-tuning (SFT)-tuned models often forget general capabilities. Latent-space representation and token-space distribution shift analyses reveal that SFT induces substantial representation and output drift, while RL preserves general-domain structure. Our results suggest a need to rethink standard post-training recipes, particularly the reliance on SFT-distilled data for advancing reasoning models.

 github.com/ReasoningTransfer/Transferability-of-LLM-Reasoning
 huggingface.co/ReasoningTransferability



## SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training

Tianzhe Chu ♠* Yuexiang Zhai ♥♣* Jihan Yang ♦ Shengbang Tong ♦
Saining Xie ♠♦ Dale Schuurmans ♠♧ Quoc V. Le ♠ Sergey Levine ♥ Yi Ma ♠♥

# 2025: The Rise of L**R**Ms (as opposed to L**L**Ms)

## Does Math Reasoning Improve General LLM Capabilities? Understanding Transferability of LLM Reasoning

Maggie Huan[1,2,*], Yuetai Li[3,*], Tuney Zheng[4,*], Xiaoyu Xu[5], Seungone Kim[1], Minxin Du[5], Radha Poovendran[3], Graham Neubig[1], Xiang Yue[1,†]
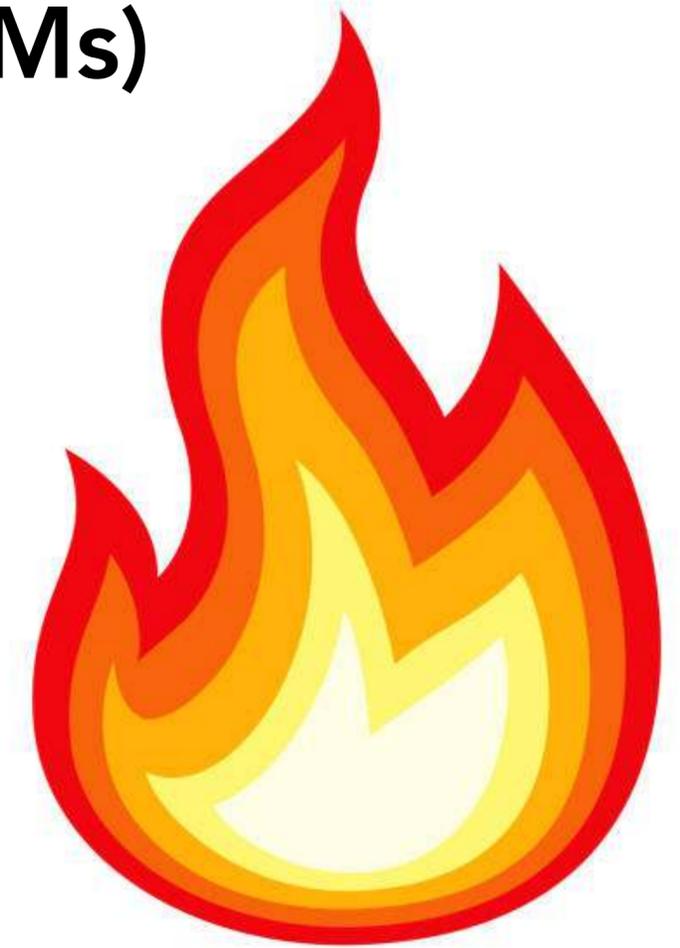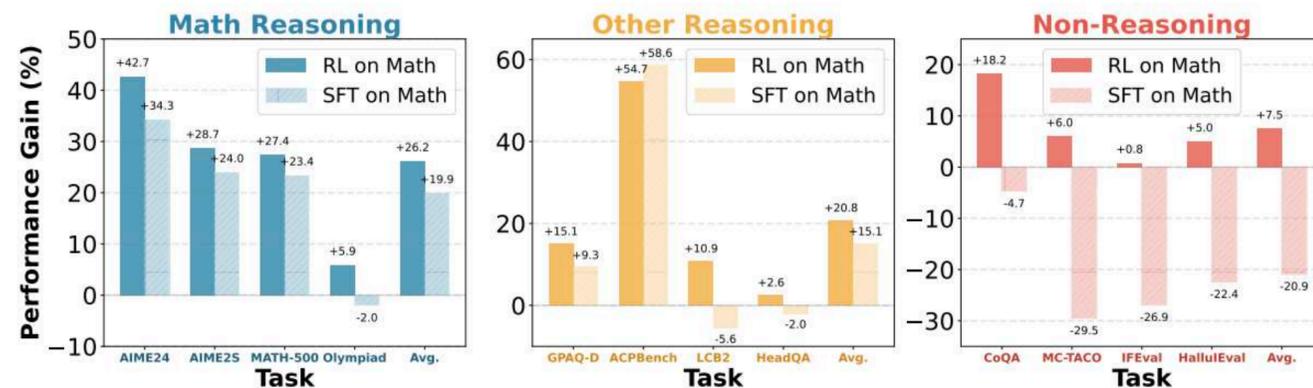
[1]Carnegie Mellon University   [2]University of Pennsylvania   [3]University of Washington
[4]M-A-P   [5]The Hong Kong Polytechnic University

ziyuh@seas.upenn.edu   yuetaili@uw.edu   xyue2@andrew.cmu.edu

**Abstract:** Math reasoning has become the poster child of progress in large language models (LLMs), with new models rapidly surpassing human-level performance on benchmarks like MATH and AIME. But as math leaderboards improve week by week, it is worth asking: *do these gains reflect broader problem-solving ability or just narrow overfitting?* To answer this question, we evaluate over 20 open-weight reasoning-tuned models across a broad suite of tasks, including math, scientific QA, agent planning, coding, and standard instruction-following. We surprisingly find that most models that succeed in math fail to transfer their gains to other domains. To rigorously study this phenomenon, we conduct controlled experiments on Qwen3-14B models using math-only data but different tuning methods. We find that reinforcement learning (RL)-tuned models generalize well across domains, while supervised fine-tuning (SFT)-tuned models often forget general capabilities. Latent-space representation and token-space distribution shift analyses reveal that SFT induces substantial representation and output drift, while RL preserves general-domain structure. Our results suggest a need to rethink standard post-training recipes, particularly the reliance on SFT-distilled data for advancing reasoning models.

github.com/ReasoningTransfer/Transferability-of-LLM-Reasoning
huggingface.co/ReasoningTransferability



## SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training

Tianzhe Chu ♠* Yuexiang Zhai ♥♣* Jihan Yang ♦ Shengbang Tong ♦
Saining Xie ♠♦ Dale Schuurmans ♠♣ Quoc V. Le ♠ Sergey Levine ♥ Yi Ma ♠♥

# Striking findings against RL



Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?

Yang Yue[1][*][†], Zhiqi Chen[1][*], Rui Lu[1], Andrew Zhao[1], Zhaokai Wang[2], Yang Yue[1], Shiji Song[1], and Gao Huang[1][✉]

[1] LeapLab, Tsinghua University    [2] Shanghai Jiao Tong University

[*] Equal Contribution    [†] Project Lead    [✉] Corresponding Author

Pass@1 better after RLVR
Pass@K worse after RLVR
— compared to the base LLM —

# Striking findings against RL

## Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?

Yang Yue[1][*][†], Zhiqi Chen[1][*], Rui Lu[1], Andrew Zhao[1], Zhaokai Wang[2], Yang Yue[1], Shiji Song[1], and Gao Huang[1][✉]

[1] LeapLab, Tsinghua University    [2] Shanghai Jiao Tong University

[*] Equal Contribution    [†] Project Lead    [✉] Corresponding Author

Pass@1 better after RLVR
Pass@K worse after RLVR
— compared to the base LLM —

# Stoking findings against RL

## Does Reinforcement ... ...ing Capacity in LLMs

Yang Yue[1][*][†], Zhiqi Chen ...
Shiji Song[1], and Gao Huang ...

[1]LeapLab, Tsinghua University   [2]Shanghai Jiao Tong ...

[*] Equal Contribution   [†] Project Lead   [✉] Corresponding Autho...

Pas... ... better after RLVR
Pa... ... worse after RLVR
— c... ...d to the base LLM —

## Echo Chamber: RL Post-training Amplifies Behaviors ... in Pretraining
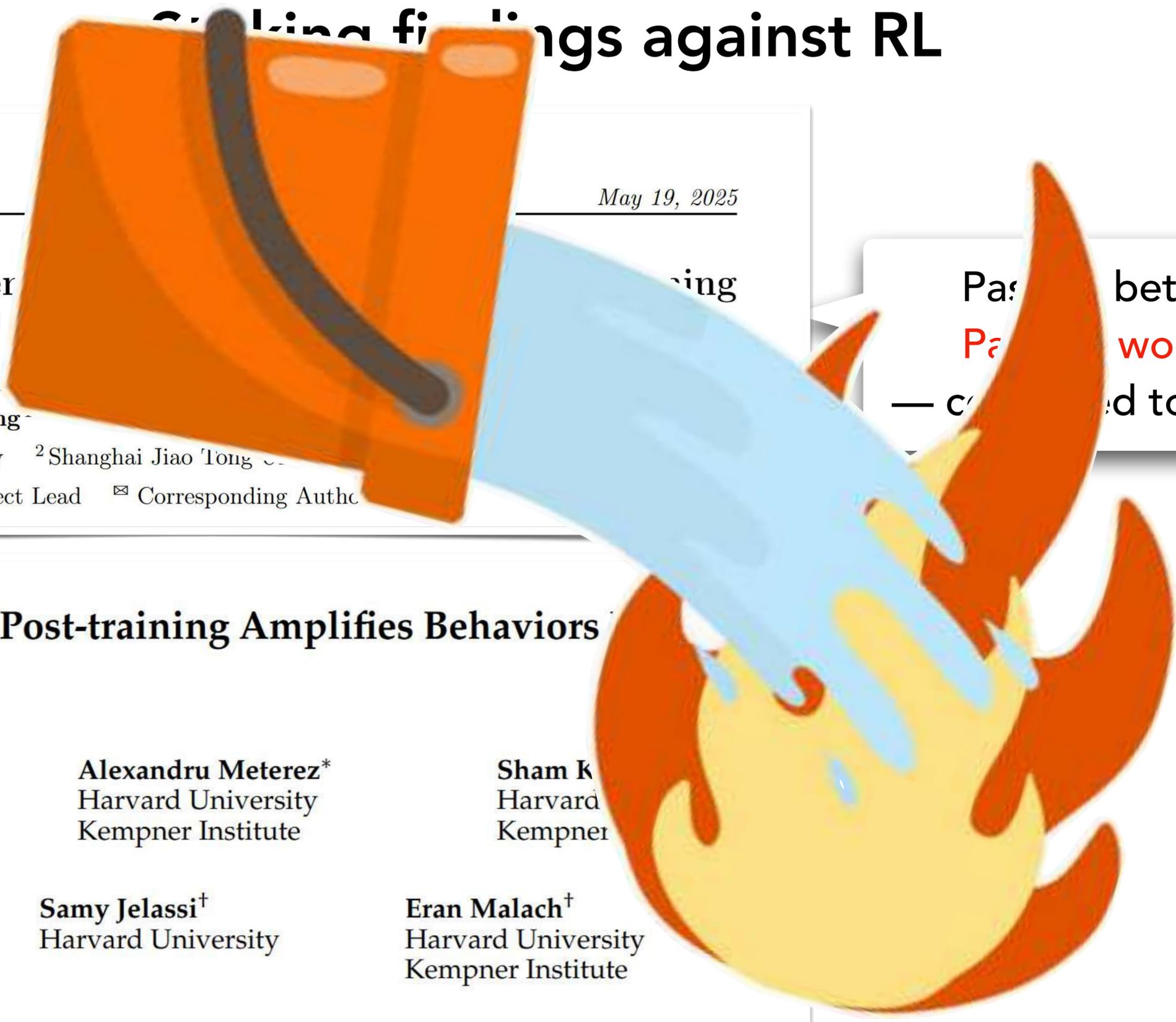
**Rosie Zhao**[*]
Harvard University
Kempner Institute

**Alexandru Meterez**[*]
Harvard University
Kempner Institute

**Sham K** ...
Harvard ...
Kempner ...

**Cengiz Pehlevan**
Harvard University
Kempner Institute

**Samy Jelassi**[†]
Harvard University

**Eran Malach**[†]
Harvard University
Kempner Institute

# Striking findings against RL

**Stella Li** ➡️ **CogSci2025**
@StellaLisy

🤯 We cracked RLVR with... Random Rewards?!
Training Qwen2.5-Math-7B with our Spurious Rewards improved MATH–500 by:
- Random rewards: +21%
- Incorrect rewards: +25%
- (FYI) Ground-truth rewards: + 28.8%
How could this even work⁉️ Here's why: 🧵
Blogpost: tinyurl.com/spurious-rewar...

## Spurious Rewards: Rethinking Training Signals in RLVR

Rulin Shao[1]* Shuyue Stella Li[1]* Rui Xin[1]* Scott Geng[1]* Yiping Wang[1]
Sewoong Oh[1] Simon Shaolei Du[1] Nathan Lambert[2] Sewon Min[3] Ranjay Krishna[1,2]
Yulia Tsvetkov[1] Hannaneh Hajishirzi[1,2] Pang Wei Koh[1,2] Luke Zettlemoyer[1]
[1]University of Washington    [2]Allen Institute for Artificial Intelligence
[3]University of California, Berkeley
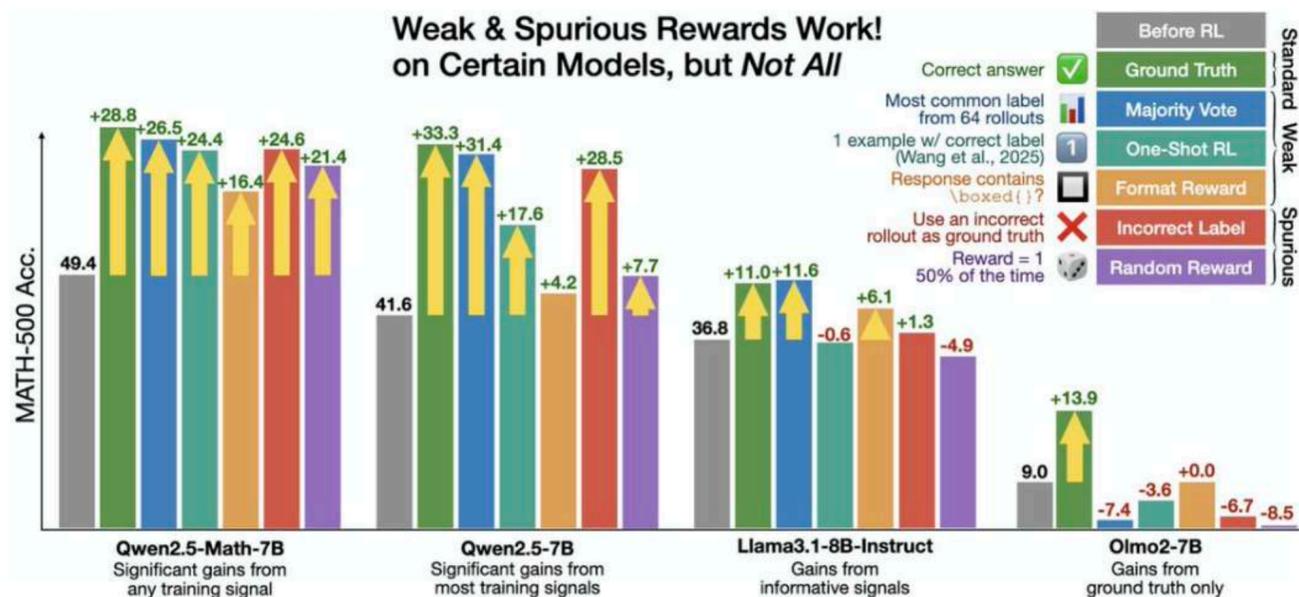{rulins,stelli,rx31,sgeng}@cs.washington.edu

Figure 1: MATH-500 accuracy after 150 steps of RLVR on various training signals. We show that

## Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?

Yang Yue[1]*[†], Zhiqi Chen[1]*, Rui Lu[1], Andrew Zhao[1], Zhaokai Wang[2], Yang Yue[1],
Shiji Song[1], and Gao Huang[1]✉

[1]LeapLab, Tsinghua University    [2]Shanghai Jiao Tong University

* Equal Contribution    † Project Lead    ✉ Corresponding Author

## Echo Chamber: RL Post-training Amplifies Behaviors Learned in Pretraining

**Rosie Zhao*** 
Harvard University
Kempner Institute

**Alexandru Meterez*** 
Harvard University
Kempner Institute

**Sham Kakade** 
Harvard University
Kempner Institute

**Cengiz Pehlevan** 
Harvard University
Kempner Institute

**Samy Jelassi†** 
Harvard University

**Eran Malach†** 
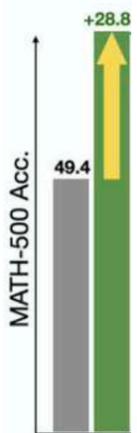Harvard University
Kempner Institute

st RL

*May 19, 2025*

Reasoning

[2], Yang Yue[1],

ors Learned

n Kakade
vard University
pner Institute

h[†]
University
stitute

**Reinforcement Learning is terrible.**

**Andrej Karpathy — "We're summoning ghosts, not building animals"**

**Striking findings against RL**

Stella Li ➡️ CogSci2025
@StellaLisy

🤯 We cracked RLVR with... Random Rewards?!
Training Qwen2.5-Math-7B with our Spurious Rewards improved MATH-500 by:
- Random rewards: +21%
- Incorrect rewards: +25%
- (FYI) Ground-truth rewards: + 28.8%
How could this even work⁉️ Here's why: 🧵
Blogpost: tinyurl.com/spurious-rewar...

**Chemistry between the base LLM and RL matters**
**Conclusions from effortless RL != effortful RL**

Weak & Spurious Rewards Work!
on Certain Models, but *Not All*

| | Standard |
| Correct answer ✅ | |
| Ground Truth | |
| Most common label from 64 rollouts 📊 | Majority Vote | Weak |
| 1 example w/ correct label (Wang et al., 2025) 1️⃣ | One-Shot RL | |
| Response contains \boxed{}? ⬜ | Format Reward | |
| Use an incorrect rollout as ground truth ❌ | Incorrect Label | Spurious |
| Reward = 1 50% of the time 🎲 | Random Reward | |

Figure 1: MATH-500 accuracy after 150 steps of RLVR on various training signals. We show that

{rulins,stelli,rx31,sgeng}@cs.washington.edu

[1] LeapLab, Tsinghua University   [2] Shanghai Jiao Tong University

* Equal Contribution   † Project Lead   ✉ Corresponding Author

**Echo Chamber: RL Post-training Amplifies Behaviors Learned in Pretraining**

**Rosie Zhao***
Harvard University
Kempner Institute

**Alexandru Meterez***
Harvard University
Kempner Institute

**Sham Kakade**
Harvard University
Kempner Institute

**Cengiz Pehlevan**
Harvard University
Kempner Institute

**Samy Jelassi**[†]
Harvard University

**Eran Malach**[†]
Harvard University
Kempner Institute

# Prolonged Reinforcement Learning

*ByteDance*

*Deepseek*

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).

# Prolonged Reinforcement Learning

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).

$$\mathscr{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \min \left( r_\theta(\tau) A(\tau), \ \text{clip}(r_\theta(\tau), 1 - \epsilon, 1 + \epsilon) A(\tau) \right) \right]$$

- **Dynamic Sampling** - targets mid-difficulty examples to maintain diverse learning signals
- **Decoupled Clipping** - 'clip-higher' uplifts low-probability tokens, encourages broader exploration.

$$\text{clip}(r_\theta(\tau), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}})$$

# Prolonged Reinforcement Learning

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).
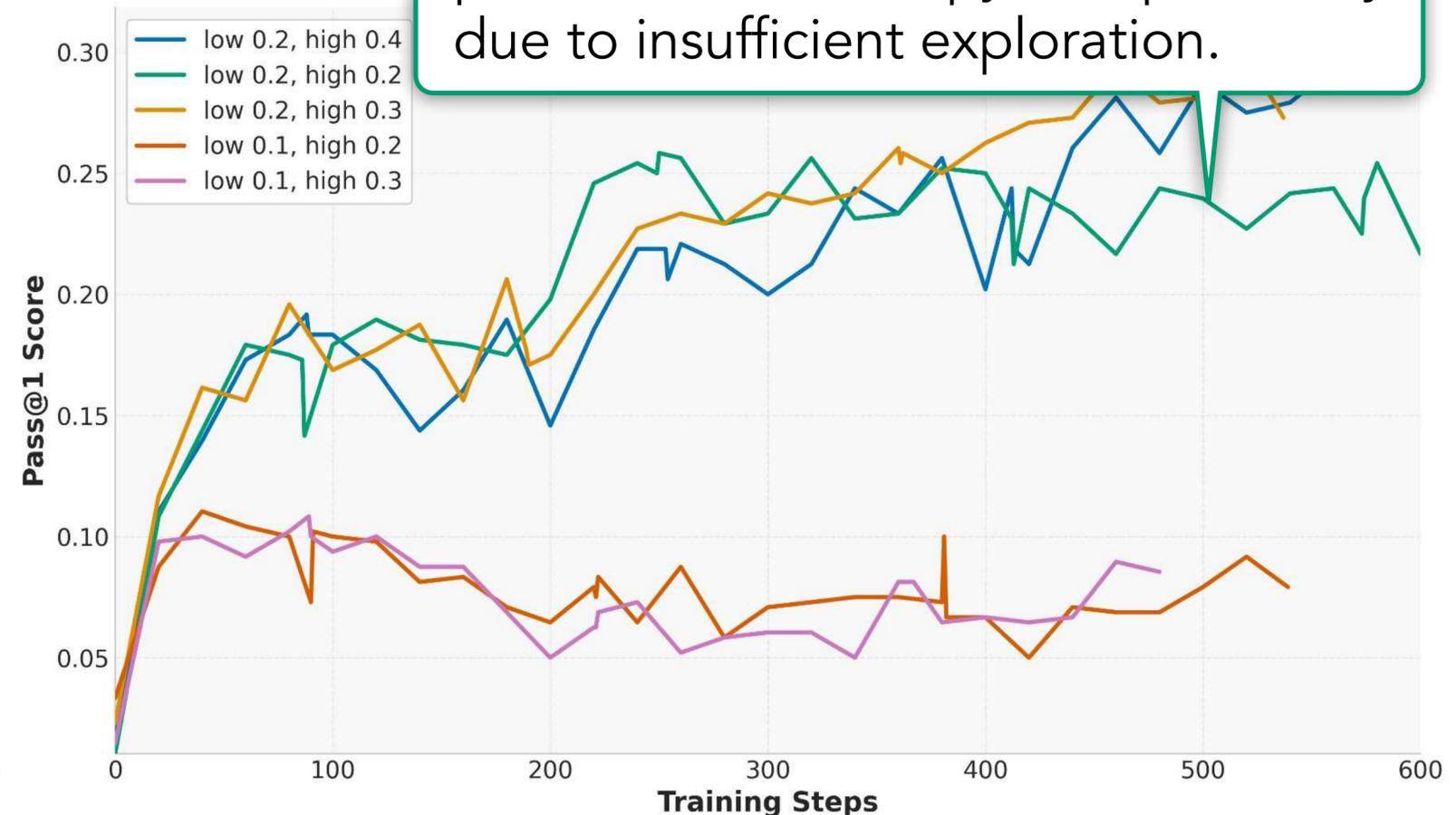
**ProRL**

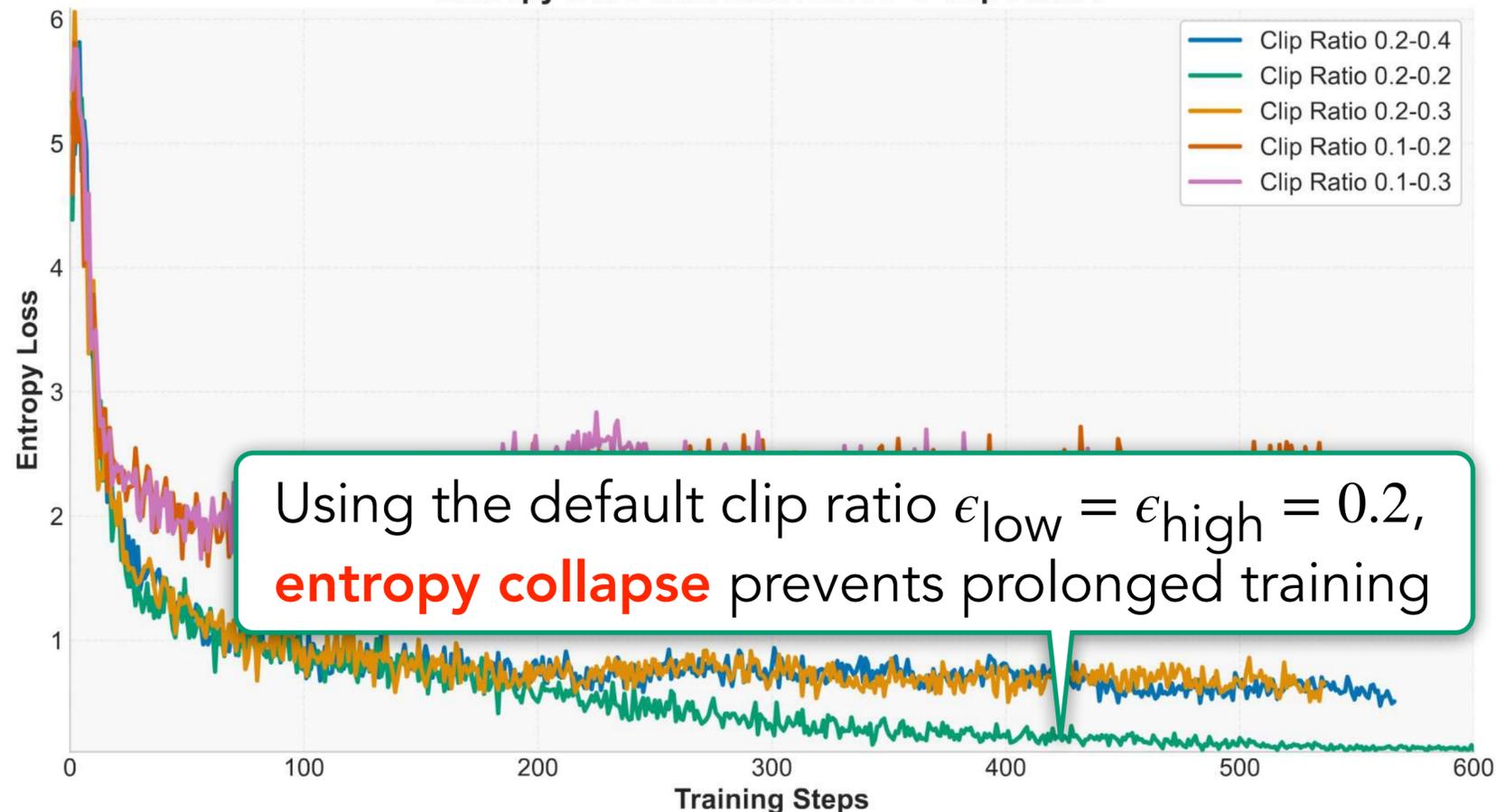Lesson I: Maintain **sustainable entropy** via clip ratio, balance exploration and exploitation.

$$\text{clip}(r_\theta(\tau),\ 1 - \epsilon_{\text{low}},\ 1 + \epsilon_{\text{high}})$$

Performance initially increases but plateaus after entropy collapse, likely due to insufficient exploration.

Using the default clip ratio $\epsilon_{\text{low}} = \epsilon_{\text{high}} = 0.2$, **entropy collapse** prevents prolonged training



Entropy Loss with Different PPO Clip Ratios

- Clip Ratio 0.2-0.4
- Clip Ratio 0.2-0.2
- Clip Ratio 0.2-0.3
- Clip Ratio 0.1-0.2
- Clip Ratio 0.1-0.3

Entropy Loss / Training Steps

- low 0.2, high 0.4
- low 0.2, high 0.2
- low 0.2, high 0.3
- low 0.1, high 0.2
- low 0.1, high 0.3

Pass@1 Score / Training Steps

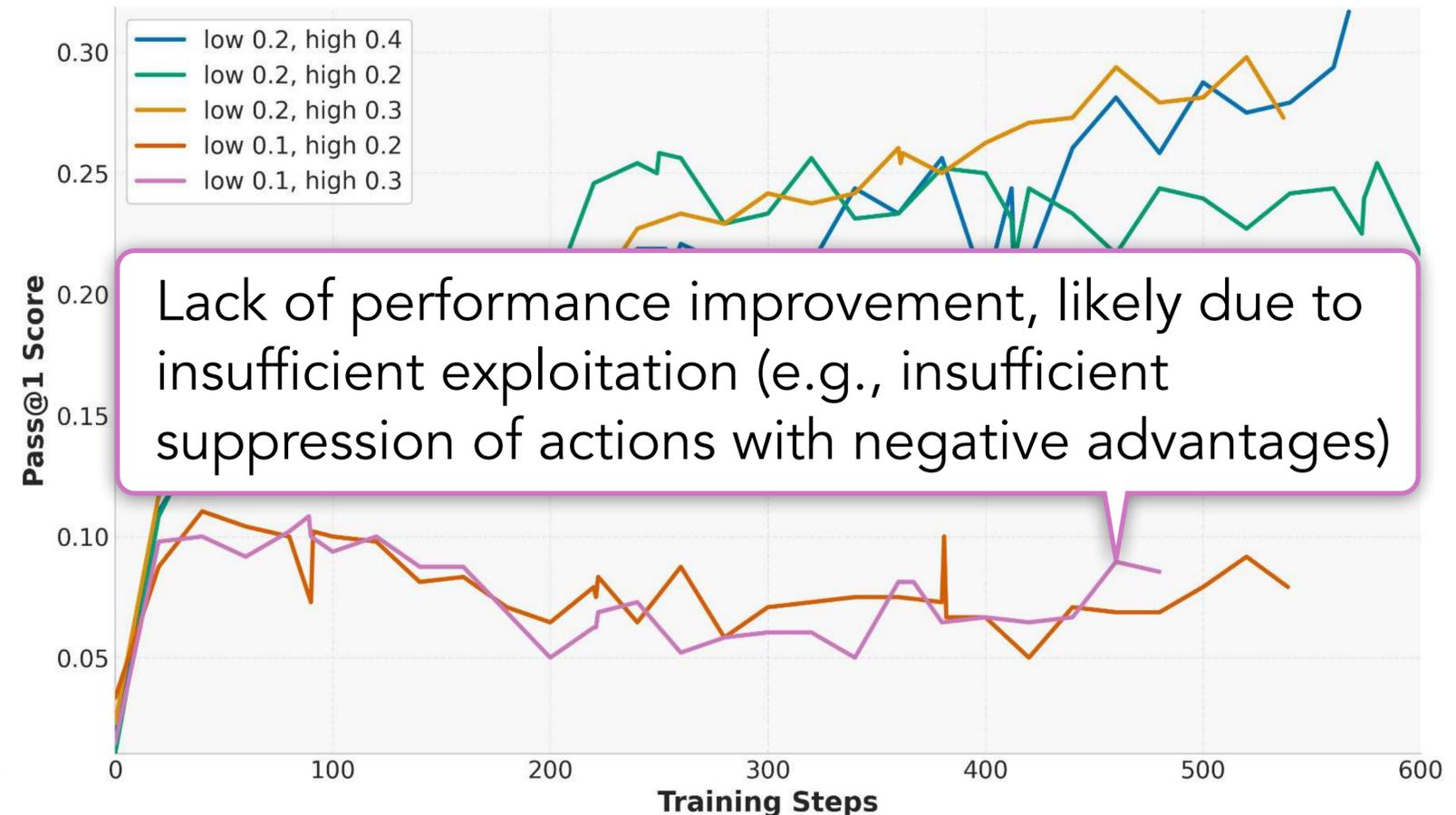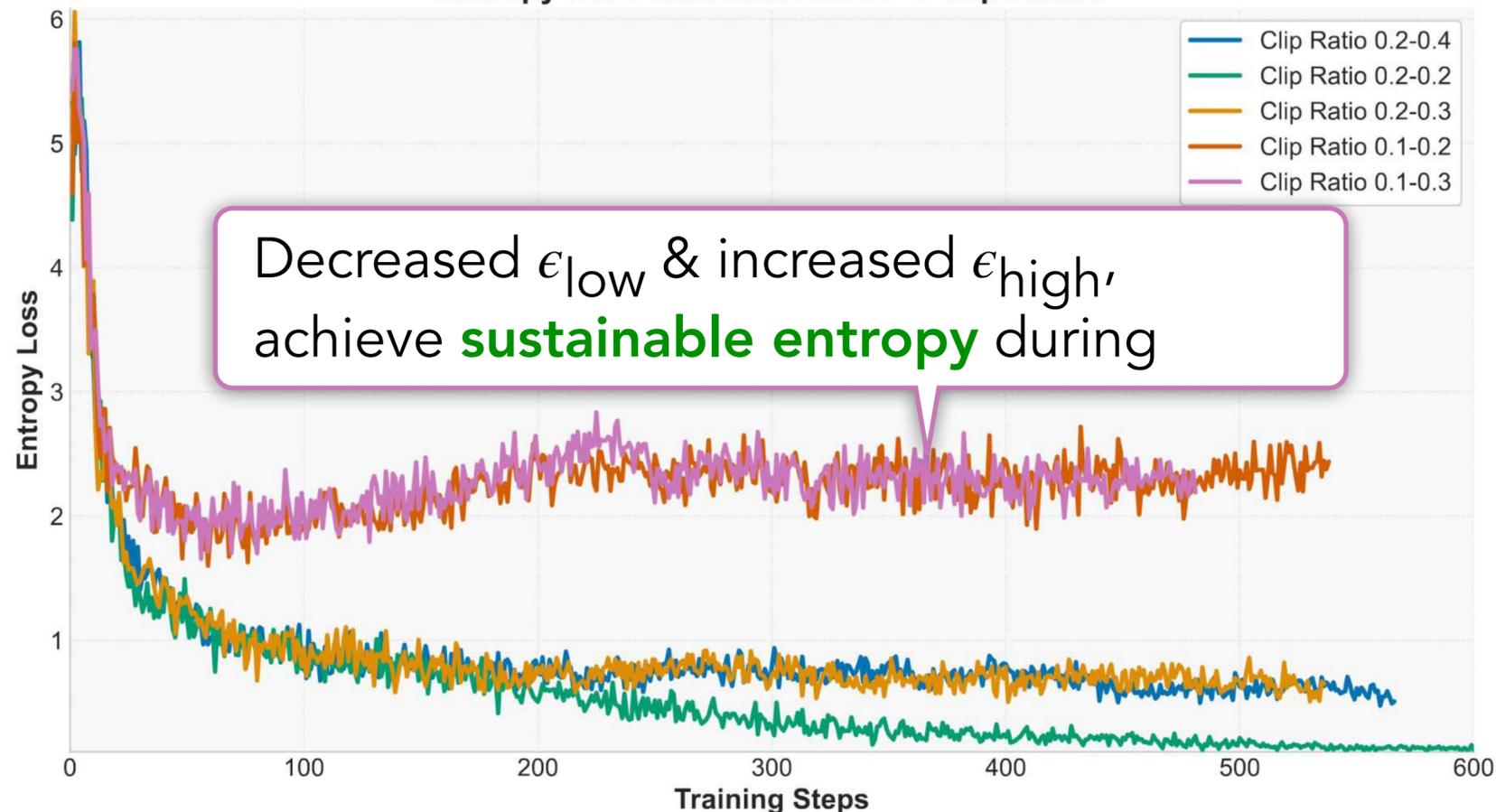# Prolonged Reinforcement Learning

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).
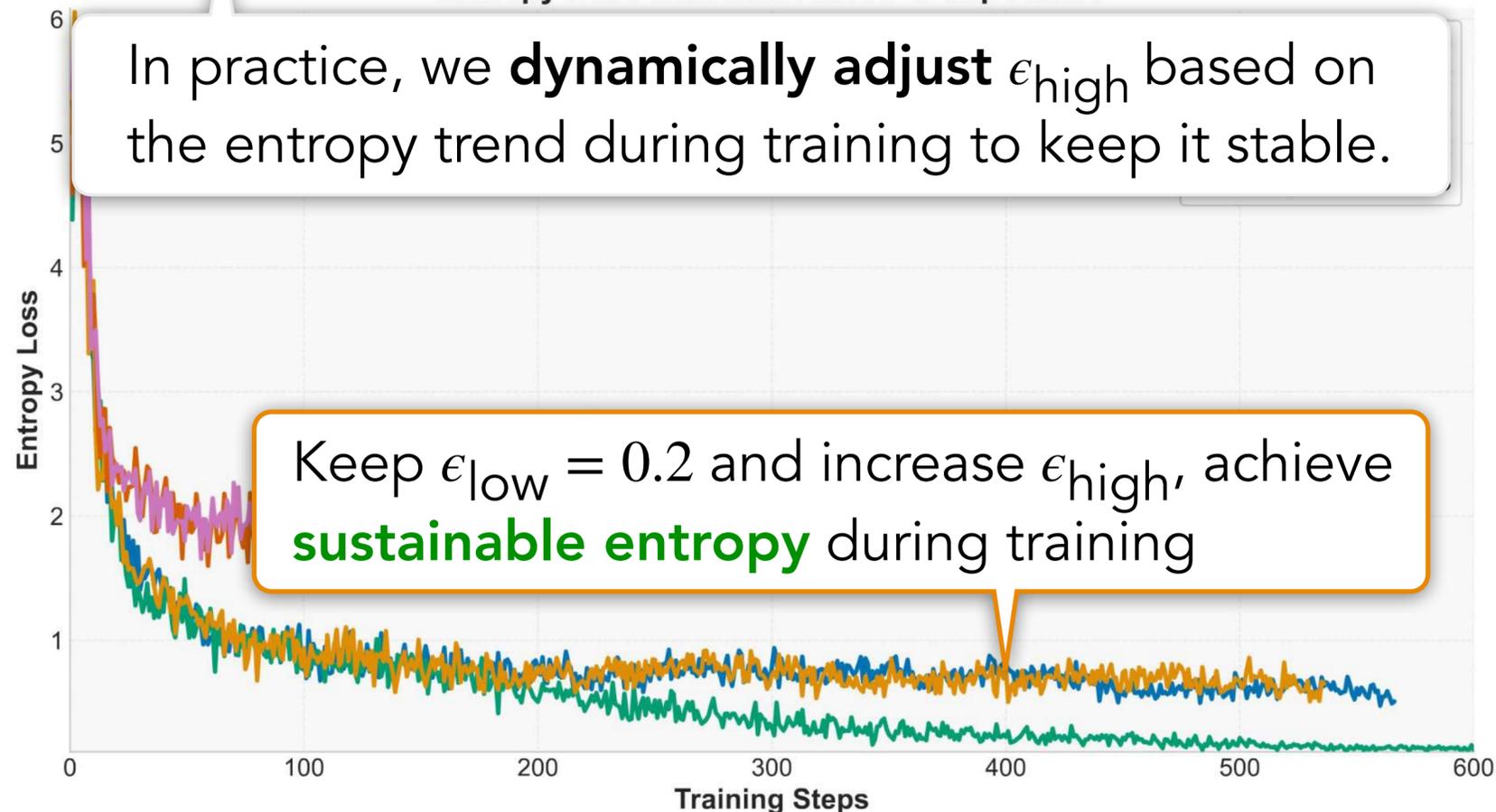
Lesson I: Maintain **sustainable entropy** via clip ratio, balance exploration and exploitation.

$$\text{clip}(r_\theta(\tau),\ 1 - \epsilon_{\text{low}},\ 1 + \epsilon_{\text{high}})$$



Decreased $\epsilon_{\text{low}}$ & increased $\epsilon_{\text{high}}$, achieve **sustainable entropy** during

Lack of performance improvement, likely due to insufficient exploitation (e.g., insufficient suppression of actions with negative advantages)

# Prolonged Reinforcement Learning

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).

Lesson I: Maintain **sustainable entropy** via clip ratio, balance exploration and exploitation.
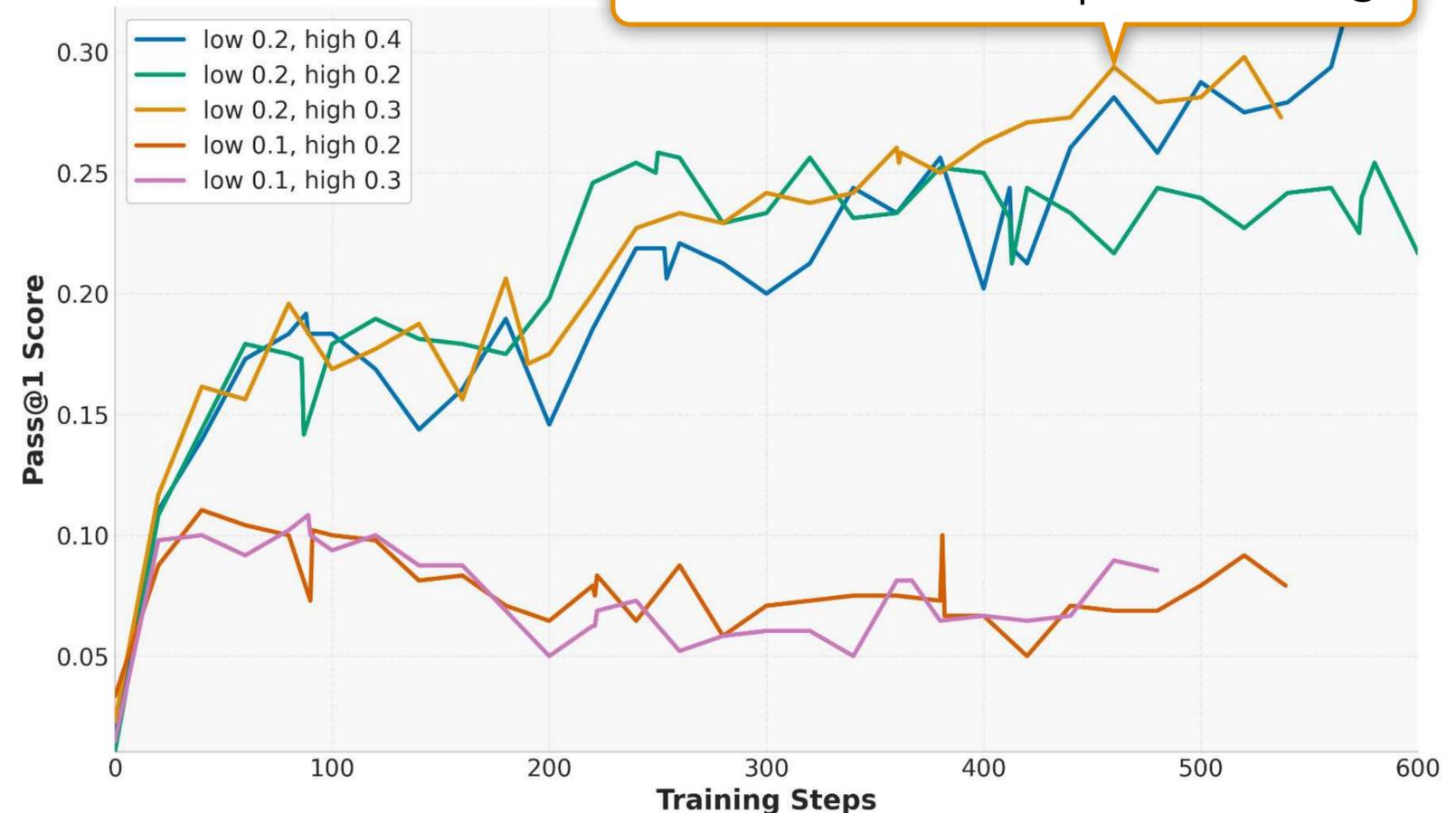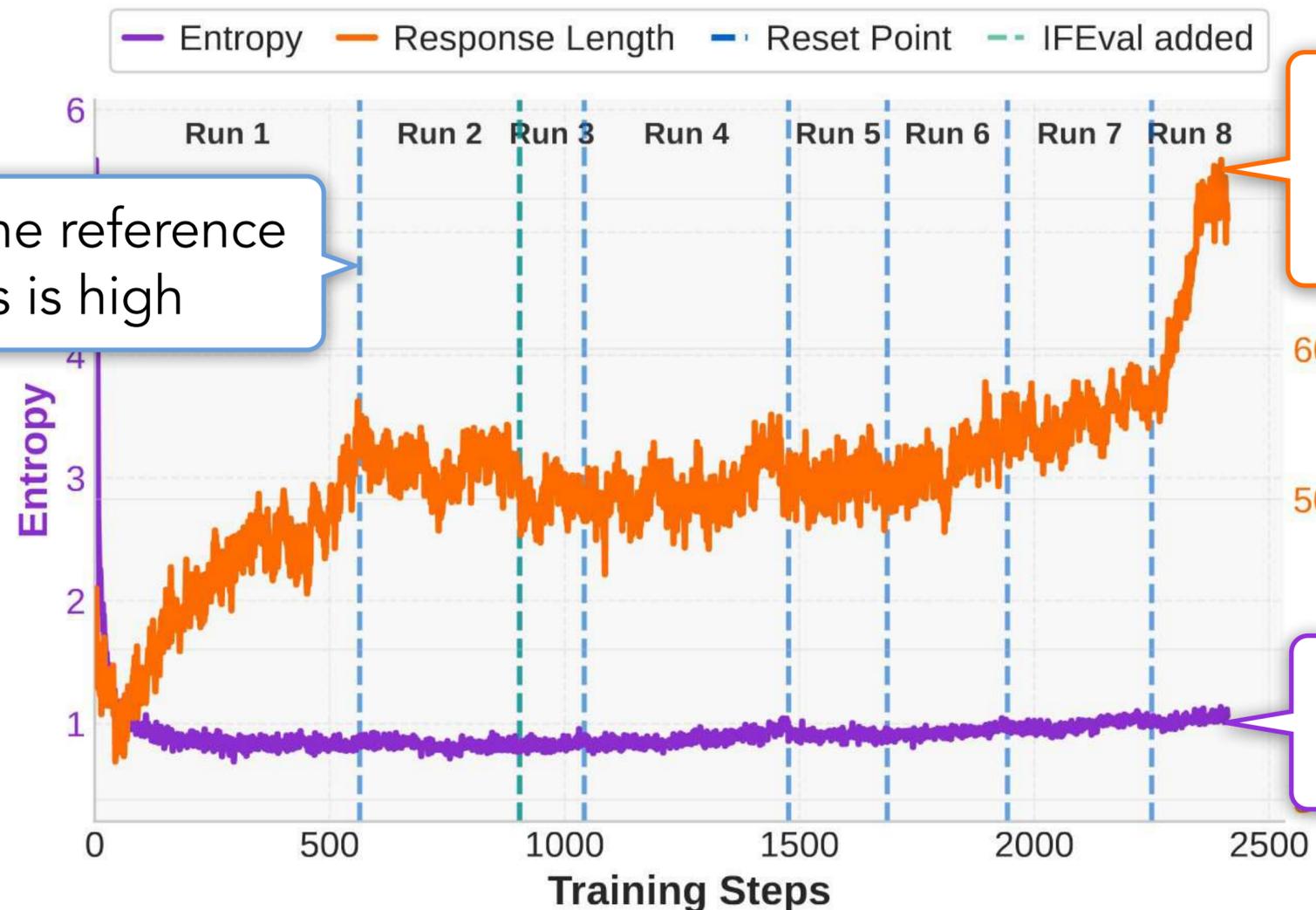
$$\text{clip}(r_\theta(\tau),\ 1 - \epsilon_{\text{low}},\ 1 + \epsilon_{\text{high}})$$

Performance keep increasing

In practice, we **dynamically adjust** $\epsilon_{\text{high}}$ based on the entropy trend during training to keep it stable.

Keep $\epsilon_{\text{low}} = 0.2$ and increase $\epsilon_{\text{high}}$, achieve **sustainable entropy** during training



Entropy Loss with Different PPO Clip Ratios



Legend:
- low 0.2, high 0.4
- low 0.2, high 0.2
- low 0.2, high 0.3
- low 0.1, high 0.2
- low 0.1, high 0.3

# Prolonged Reinforcement Learning

ProRL is built on top of one of the best practices in RL—Decoupled Clip and Dynamic Sampling Policy Optimization (**DAPO**), a variation of Group Relative Policy Optimization (GRPO).
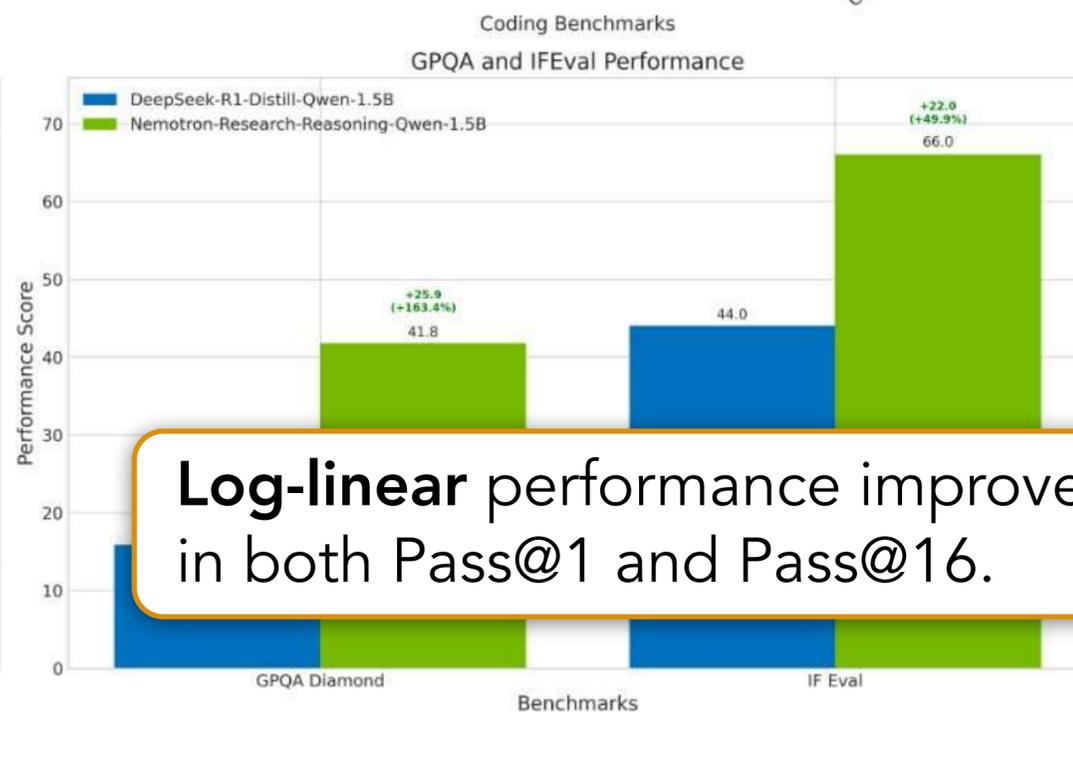
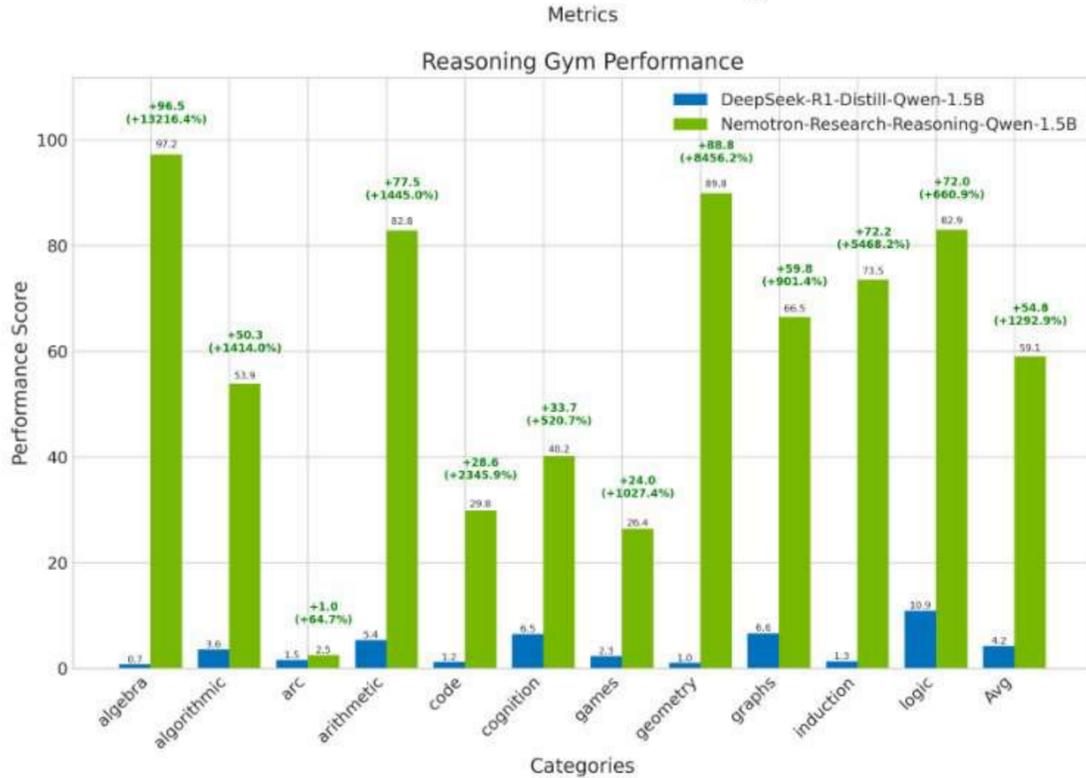Lesson III: Dynamically adjust hyper-parameters to **keep one hero run alive**
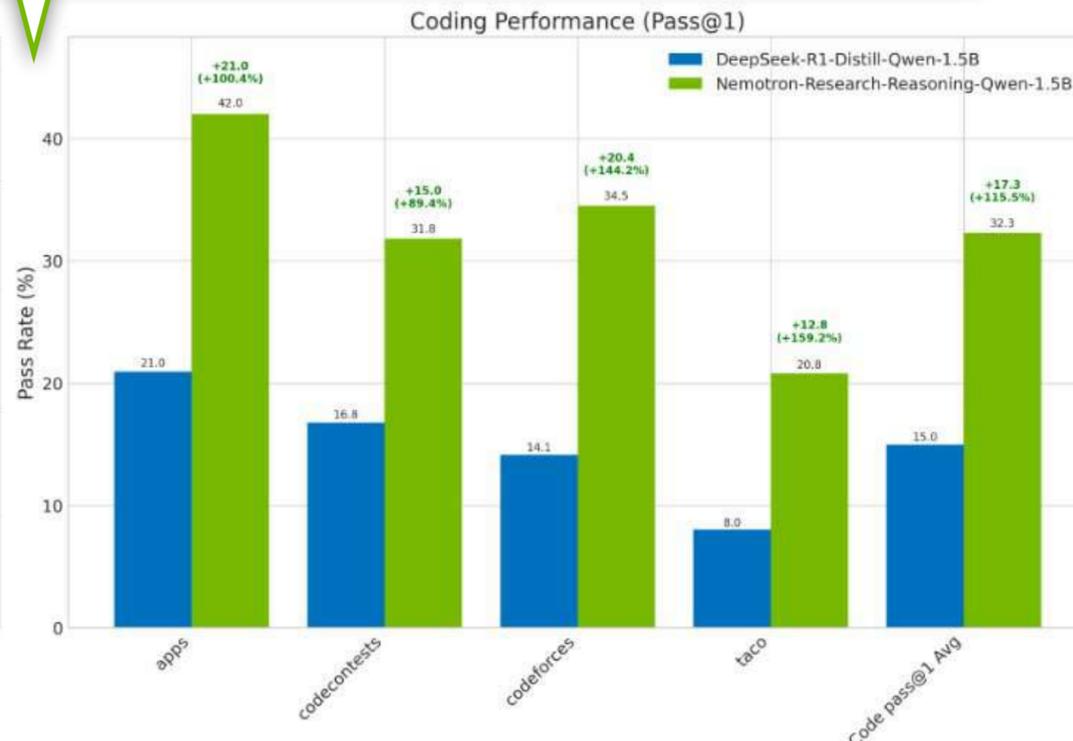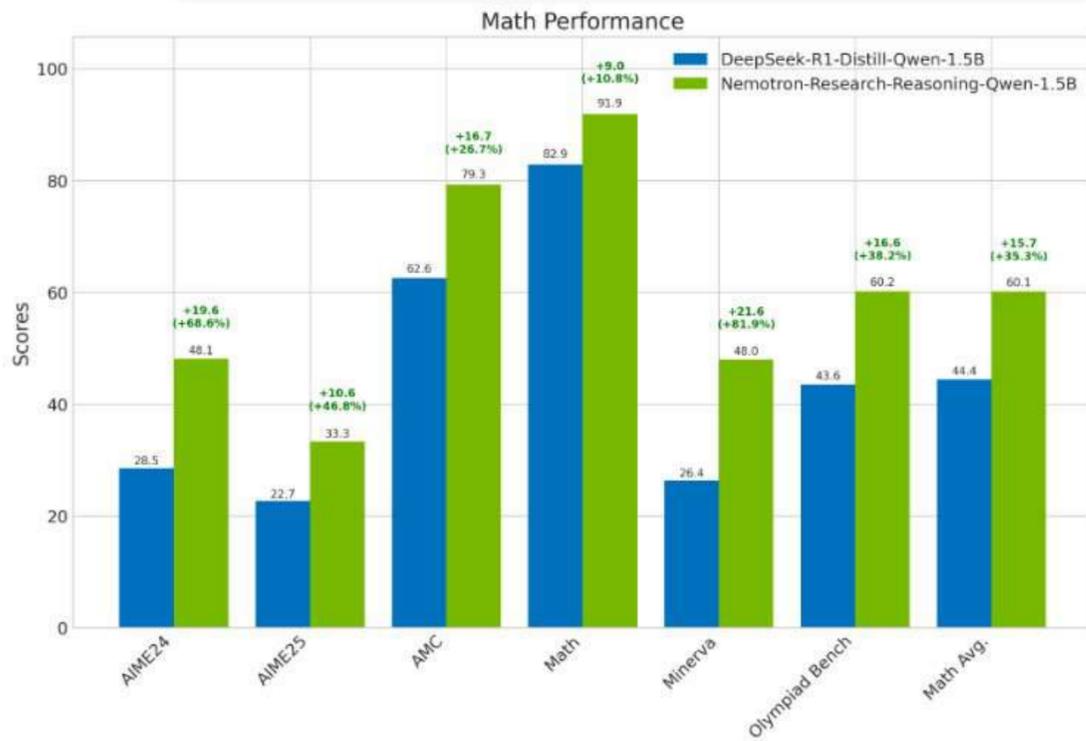


Periodically reset the reference policy when KL loss is high

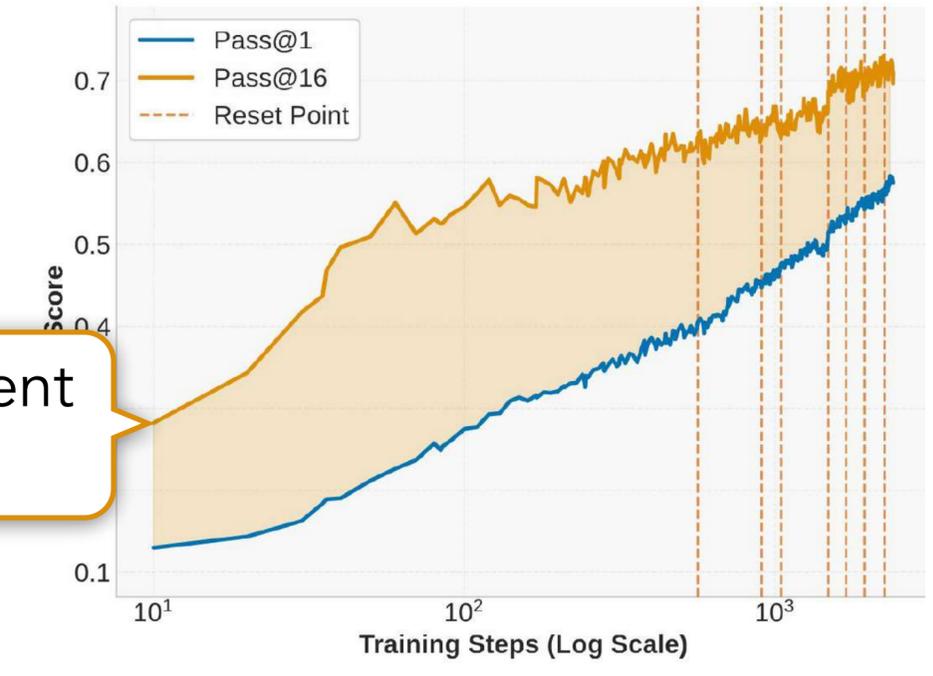Keep response length under 8k for training efficiency; extend context length in the final run.

**Dynamically adjust** $\epsilon_{high}$ to maintain **sustainable entropy**

Evaluated across diverse domains—math, code, STEM, reasoning gym, and instruction following—**ProRL results in significant gains** over the base model, **Distill-R1-1.5B**, in all cases.
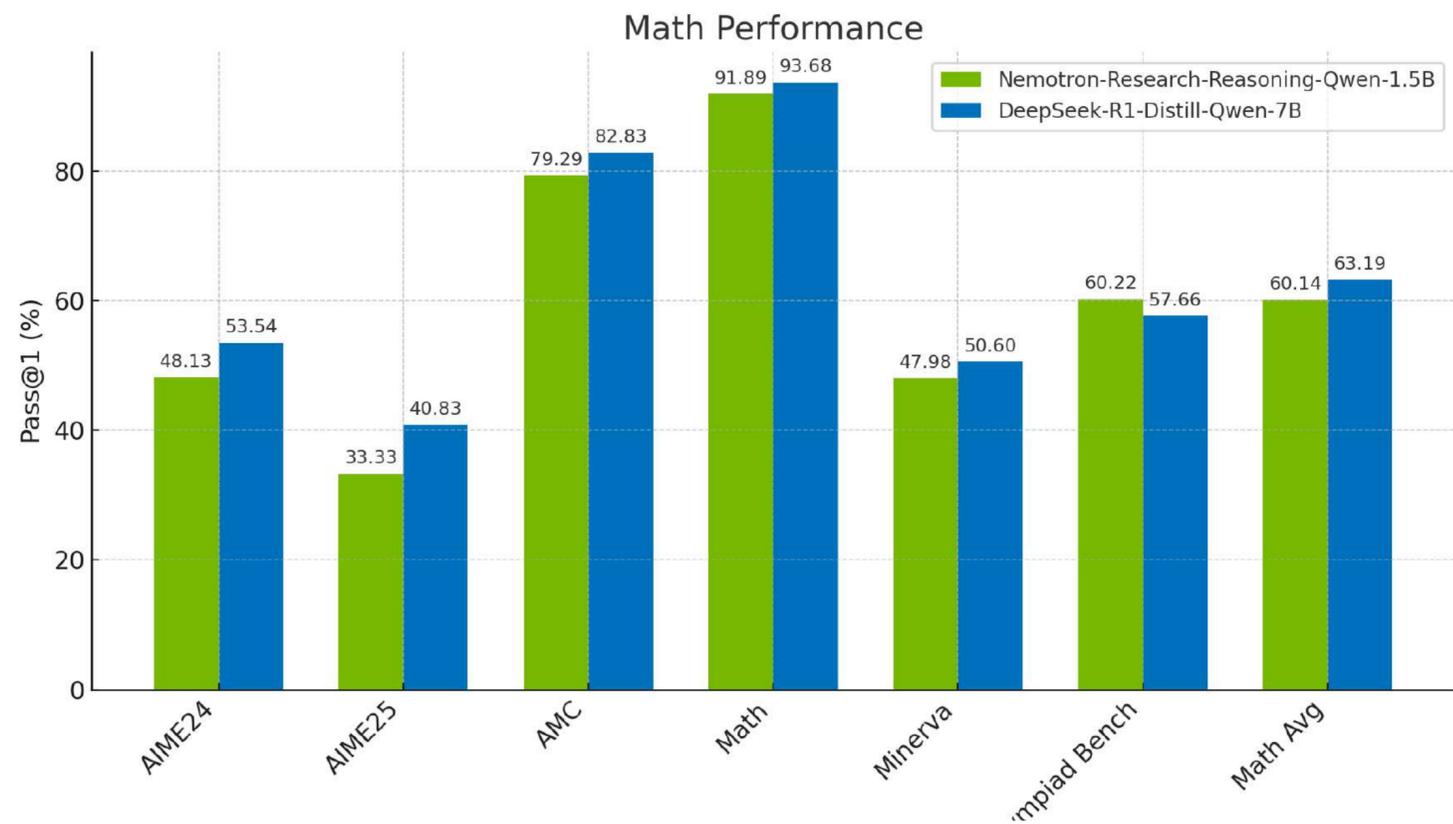
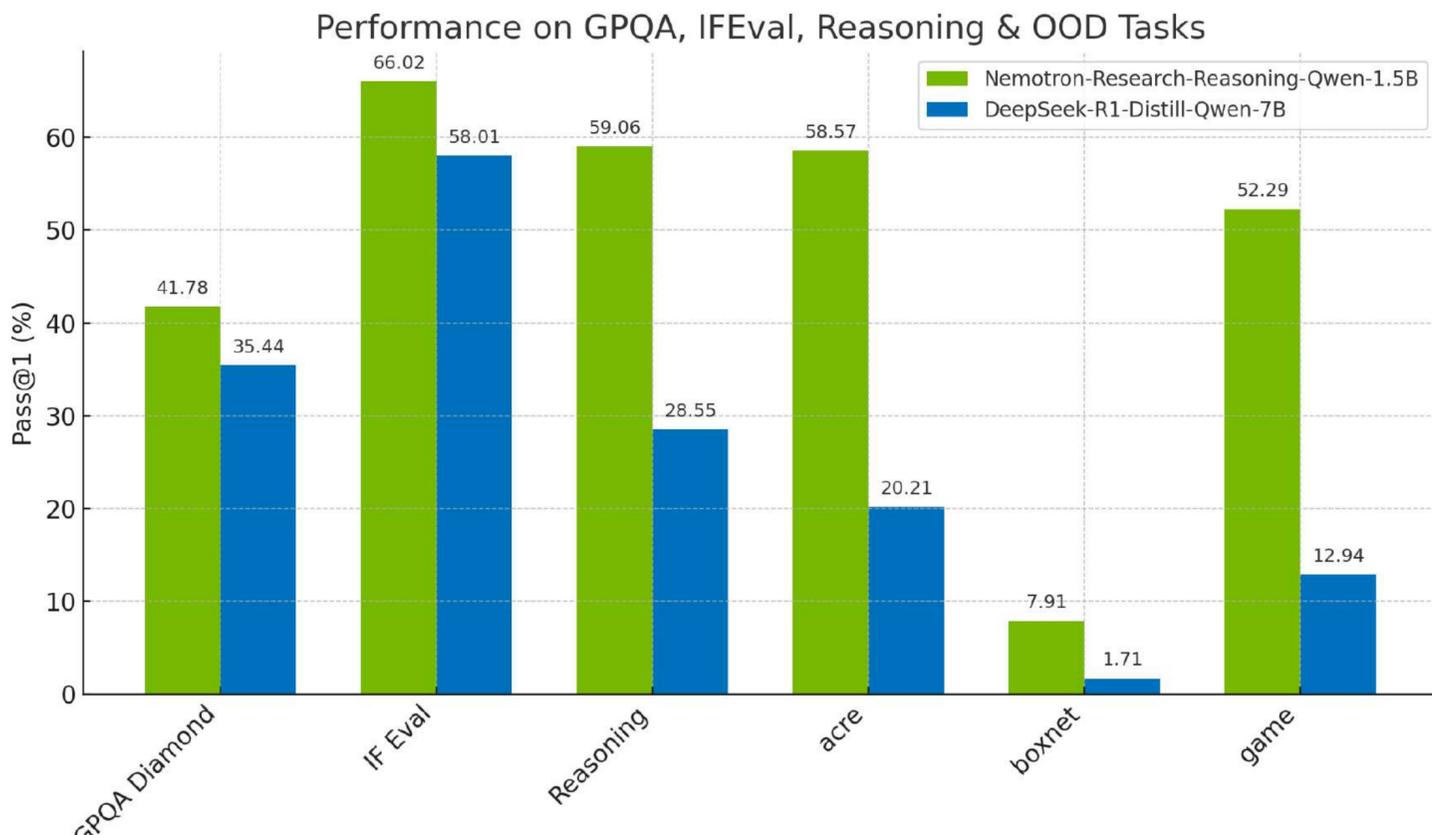The gain is particularly large in out-of-domain tasks and reasoning gym where the base model struggles.
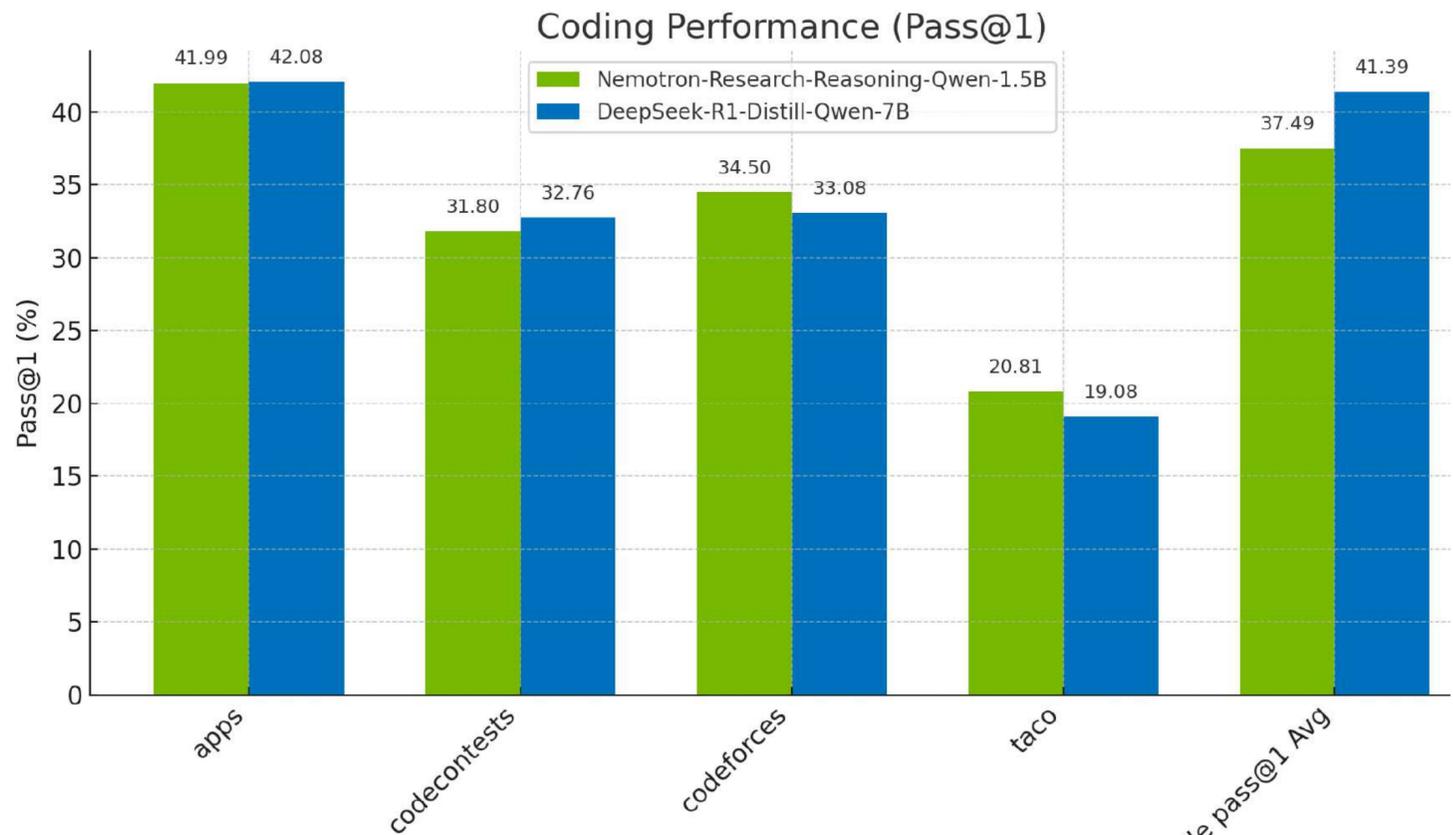
**Log-linear** performance improvement in both Pass@1 and Pass@16.

# The SOTA 1.5B Reasoning Model



**Math Performance**

Our ProRL-ed model, **Nemotron-Reasoning -1.5B**, achieves comparable or better performance than the 4.5× larger model, **DeepSeek-R1-7B**.

**Coding Performance (Pass@1)**

**Performance on GPQA, IFEval, Reasoning & OOD Tasks**

**Striking findings against RL**

Chemistry between the base LLM and RL matters
Conclusions from *effortless* RL != *effortful* RL

# Conclusion (🌹 vs 💀)

🌹 Conclusions from effortless RL != effortful RL

🌹 An effortful, proloooooonged RL on 1.5B model can go faaaaaaaaaaarrrr

🌹 ProRL shows that a dynamic control of entropy is critical

BroRL: Scaling Reinforcement Learning via Broadened Exploration

Jian Hu[1]    Mingjie Liu[1]    Ximing Lu[1]    Fang Wu[2]    Zaid Harchaoui[3]    Shizhe Diao[1]
Yejin Choi[1]    Pavlo Molchanov[1]    Jun Yang[1]    Jan Kautz[1]    Yi Dong[1]
[1]NVIDIA    [2]Stanford University    [3]University of Washington

💀 Can effortful RL succeed on say, GPT2?

# The *Entropy* Mechanism of Reinforcement Learning for Reasoning Language Models

Ganqu Cui[1]*Yuchen Zhang[1,4]*Jiacheng Chen[1]*Lifan Yuan[3], Zhi Wang[5], Yuxin Zuo[2], Haozha...
Yuchen Fan[1], Huayu Chen[2], Weize Chen[2], Zhiyuan Liu[2], Hao Peng[3], Lei Bai[1], Wanli Ouyang...
Yu Cheng[1,6]†Bowen Zhou[1,2]†Ning Ding[2,1]†
[1] Shanghai AI Laboratory [2] Tsinghua University [3] UIUC [4] Peking University [5] Nanjing University [6] CUHK

# The Unreasonable Effectiveness of Entropy Minimization in LLM Reasoning

Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han, Hao Peng
University of Illinois Urbana-Champaign

# Beyond the 80/20 Rule: High-Entropy Minority Tokens Drive Effective Reinforcement Learning for LLM Reasoning

Shenzhi Wang[1,2], Le Yu[1], Chang Gao[1], Chujie Zheng[1], Shixuan Liu[1], Rui Lu[2],
Dang[1], Xiong-Hui Chen[1], Jianxin Yang[1], Zhenru Zhang[1], Yuqiong Liu[1], An Yang[1],
Andrew Zhao[2], Yang Yue[2], Shiji Song[2], Bowen Yu[1,✉,†], Gao Huang[2,✉], Junyang Lin[1]
[1] Qwen Team, Alibaba Inc.      [2] LeapLab, Tsinghua University

**EPO**😅  **EPO**😅

# Evolutionary Policy Optimization

Jianren Wang[*1]   Yifan Su[*1]   Abhinav Gupta[1]   Deepak Pathak[1]

# EPO: ENTROPY-REGULARIZED POLICY OPTIMIZATIO FOR LLM AGENTS REINFORCEMENT LEARNING

Wujiang Xu[1]*Wentian Zhao[2], Zhenting Wang[1], Yu-Jhe Li[2],
Can Jin[1], Mingyu Jin[1], Kai Mei[1], Kun Wan[2], Dimitris N. Metaxas[1]
[1] Rutgers University     [2] Adobe Inc.

**Striking findings against RL**

🤯 We cracked RLVR with... Random Rewards?!

Training Qwen2.5-Math-7B with our Spurious Rewards improved MATH-500 by:
- Random rewards: +21%
- Incorrect rewards: +25%
- (FYI) Ground-truth rewards: + 28.8%
How could this even work⁉️ Here's why: 🧵
Blogpost: tinyurl.com/spurious-rewar...

Chemistry between the base LLM and RL matters
Conclusions from effortless RL != effortful RL

Conclusions from effortless SFT != effortful SFT

Stella Li ➡️ CogSci2025
@StellaLisy

{rulins,stelli,rx31,sgeng}@cs.washington.edu

[1] LeapLab, Tsinghua University   [2] Shanghai Jiao Tong University

**Rosie Zhao***
Harvard University
Kempner Institute

**Alexandru Meterez***
Harvard University
Kempner Institute

**Sham Kakade**
Harvard University
Kempner Institute

**Cengiz Pehlevan**
Harvard University
Kempner Institute

**Samy Jelassi**[†]
Harvard University

**Eran Malach**[†]
Harvard University
Kempner Institute

Figure 1: MATH-500 accuracy after 150 steps of RLVR on various training signals. We show that

**Striking findings against RL**

Conclusions from effortless SFT != effortful SFT

Why is Qwen's chemistry so good with RL?
SFT or SFT-style "mid-training" during pre-training

**Spurious Rewards: Rethinking Training Signals in RLVR**

Rulin Shao[1*] Shuyue Stella Li[1*] Rui Xin[1] Scott Geng[1] Yiping Wang[1]
Sewoong Oh[1] Simon Shaolei Du[1] Nathan Lambert[2] Sewon Min[3] Ranjay Krishna[1,2]
Yulia Tsvetkov[1] Hannaneh Hajishirzi[1,2] Pang Wei Koh[1,2] Luke Zettlemoyer[1]

[1]University of Washington  [2]Allen Institute for Artificial Intelligence
[3]University of California, Berkeley
{rulins,stelli,rx31,sgeng}@cs.washington.edu

**Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?**

*May 19, 2025*

Yang Yue[1*†], Zhiqi Chen[1*], Rui Lu[1], Andrew Zhao[1], Zhaokai Wang[2], Yang Yue[1],
Shiji Song[1], and Gao Huang[1✉]

[1]LeapLab, Tsinghua University  [2]Shanghai Jiao Tong University
* Equal Contribution  † Project Lead  ✉ Corresponding Author

Figure 1: MATH-500 accuracy after 150 steps of RLVR on various training signals. We show that

# 2025: The Rise of L**R**Ms (as opposed to L**L**Ms)



### Does Math Reasoning Improve General LLM Capabilities? Understanding Transferability of LLM Reasoning

Maggie Huan[1,2,*], Yuetai Li[3,*], Tuney Zheng[4,*], Xiaoyu Xu[5], Seungone Kim[1], Minxin Du[5], Radha Poovendran[3], Graham Neubig[1], Xiang Yue[1,†]

[1]Carnegie Mellon University  [2]University of Pennsylvania  [3]University of Washington
[4]M-A-P  [5]The Hong Kong Polytechnic University

ziyuh@seas.upenn.edu   yuetaili@uw.edu   xyue2@andrew.cmu.edu

**Abstract:** Math reasoning has become the poster child of progress in large language models (LLMs), with new models rapidly surpassing human-level performance on benchmarks like MATH and AIME. But as math leaderboards improve week by week, it is worth asking: *do these gains reflect broader problem-solving ability or just narrow overfitting?* To answer this question, we evaluate over 20 open-weight reasoning-tuned models across a broad suite of tasks, including math, scientific QA, agent planning, coding, and standard instruction-following. We surprisingly find that most models that succeed in math fail to transfer their gains to other domains. To rigorously study this phenomenon, we conduct controlled experiments on Qwen3-14B models using math-only data but different tuning methods. We find that reinforcement learning (RL)-tuned models generalize well across domains, while supervised fine-tuning (SFT)-tuned models often forget general capabilities. Latent-space representation and token-space distribution shift analyses reveal that SFT induces substantial representation and output drift, while RL preserves general-domain structure. Our results suggest a need to rethink standard post-training recipes, particularly the reliance on SFT-distilled data for advancing reasoning models.

github.com/ReasoningTransfer/Transferability-of-LLM-Reasoning
huggingface.co/ReasoningTransferability

Then, SFT cannot generalize?
It depends...

### SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training

Tianzhe Chu ♠* Yuexiang Zhai ♥♣* Jihan Yang ♦ Shengbang Tong ♦
Saining Xie ♠♦ Dale Schuurmans ♠♣ Quoc V. Le ♣ Sergey Levine ♥ Yi Ma ♠♥

Conclusions from effortless SFT != effortful SFT

In this talk:

===

ProRL: Prolonged RL
Prismatic Synthesis

RL as Pretraining

===

"Smaller but Better"
"Algorithms for the Win"

*David vs. Goliath:*

*the Art of Scaling Intelligence*

*the Era of Extreme-Scale Neural Models*

eary

r)

Three components to innovate:

Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏

# The Era of Brute-Force Scaling is Over
## The Era of Smart Scaling Begins



-- Ilya Sutskever from his test time award talk at NeurIPS 2024

# Synthetic Data to the Rescue!

Already used in both mid / post-training:



Sometimes, for the entirety of LLM training…

Common practice & assumption

(1) rely on the largest / strongest teacher models

# Obvious Concerns: Mode Collapse



Conclusions from effortless synthetic data
!= effortful synthetic data

nature

Explore content ✓    About the journal ✓    Publish with us ✓

nature › articles › article

PRICE $8.99    THE NEW YORKER    NOV. 20, 2023

**Abstract**

Stable diffusion re
(ref. 2) and GPT-4 (

THE CURSE OF RECURSION:
TRAINING ON GENERATED DATA MAKES MODELS FORGET

Ilia Shumailov*          Zakhar Shumaylov*          Yiren Zhao          Yarin Gal
University of Oxford     University of Cambridge    Imperial College London    University of Oxford

Nicolas Papernot                    Ross Anderson
University of Toronto & Vector Institute    University of Cambridge & University of Edinburgh

# Gradient as Data Representation

**Input $x$**

Bill walks 0.5 mile south, then 0.75 mile east, and finally 0.5 mile south. How many miles is he, in a direct line, from his starting point?

**Output $y$**

Bill walked 0.5+0.5 = 1.0 miles south total and walked 0.75 miles east in total. Therefore, using Pythagorean… the answer is 1.25 miles.

Use a small reference model to compute $\nabla_\theta P(y \mid x)$

**Rademacher Projection** *for dimensionality reduction*

Gradient encodes the mapping from $x$ to $y$;
which naturally represents
"reasoning"
between input and output!

Sample A

**1024-dim representation of data point**

# G-Vendi Score

# Prismatic Synthesis

Sample data

R1-32B

Overgenerated data

Filters

Filtered data

R1-7B

Prismatic Synthesis

Sample data

R1-32B

Overgenerated data

Filters

Filtered data

R1-7B

# Prismatic Synthesis improves long-CoT reasoning



Outperforms baselines, while
(1) using **20x smaller R1-32B** as data generator,
(2) relying on **zero human-labeled answers**, i.e., entirely model generated!

Legend: OpenThinker-7B, OpenThinker2-7B, R1-distill-7B, PrismMath-7B

Categories: MATH-500, AIME24, AIME25, AMC23, MATH^2, OlympiadBench, GSM8k-Platinum, OOD Avg

# Concluding Thoughts

- Reasoning requires data that transcends the internet data

- Synthetic data to the rescue

- RL can be viewed as a form of (implicit) synthetic data generation via exploration

- RL or not, existing methods often lack the bird eye view on the overall diversity

- Systematic diversification (+ quality check) can make a big difference, even to overcome 20x size difference of the R1 teacher models

In this talk:

===

ProRL: Prolonged RL
Prismatic Synthesis

RL as Pretraining

===

"Smaller but Better"
"Algorithms for the Win"

h: Underdogs,
Art of Battling

*David vs. Goliath:*

*the Art of Scaling Intelligence*

*the Era of Extreme-Scale Neural Models*

Three components to innovate:

Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌍

# RLP
## Reinforcement as a Pretraining Objective
— *ICLR 2026*—

**Ali Hatamizadeh**[†1]**, Syeda Nahida Akter**[†2*]**,   Shrimai Prabhumoye**[†1,3]**,  Jan Kautz**[1]**,
Mostofa Patwary**[1]**,  Mohammad Shoeybi**[1]**,   Bryan Catanzaro**[1]**,   Yejin Choi**[1,4]

NVIDIA[1], Carnegie Mellon University[2], Boston University[3], Stanford University[4]
ahatamizadeh@nvidia.com, sprabhumoye@nvidia.com

# The End of Ever More Pretraining

- Scaling compute is no longer the bottleneck

- High-quality data is finite and increasingly exhausted

- The era shifts from scaling data to extracting more value per token

-- Ilya Sutskever from his test time award talk at NeurIPS 2024



Pre-training as we know it will end

Compute is growing:

- Better hardware
- Better algorithms
- Larger clusters

Data is not growing:

- We have but one internet
- **The fossil fuel of AI**

IA

# The Problem with Standard LLM Training

Reasoning is an afterthought — we can do better

**Pretraining**
(Gather World Knowledge)

**SFT**
=
Supervised Finetuning
(Mimics reasoning format)

**RLHF/RLVR**
=
Reinforcement Learning
(Reasoning as an add-on)

**Imitation Learning**

**Exploration Learning**

# Vanilla Pretraining vs RLP Pretraining

Same context — but RLP induces reasoning.

Photosynthesis is the process plants, algae and some bacteria use to make their own food using _____

**Vanilla Pretraining (Next Token Prediction)**

**RLP Pretraining**

<think>Photosynthesis relies on solar energy. Hence the next token must be sunlight. </think>

**P(next token | context)**

(Pattern Completion)

**P(next token | context, thought)**

(Reasoning driven prediction)

**sunlight**

**Key difference:** RLP produces an explicit reasoning trace before predicting the token — making the "why" visible and trainable, not just the final answer.

# From NTP to RLP: Teaching Models to Think First

Treat CoT as an action; reward by information gain; update only thought tokens.

## How RLP works

- Sample a thought before predicting the next token.

- Compare likelihoods with vs. without the thought

$$r(c_t) = \log p_\theta\big(x_t \mid x_{<t}, c_t\big) \ - \ \log \bar{p}_\phi\big(x_t \mid x_{<t}\big)$$

- Verifier-free, position-wise credit at every step.

# RLP Framework

# Training Recipe

- **Single network** acts as policy + reasoned predictor.

- **EMA "no-think" teacher** with slow lag $\phi \leftarrow \tau\,\phi + (1 - \tau)\,\theta, \qquad \tau = 0.999$

- **Group-relative advantages** with inclusive-mean correction

$$\bar{r} = \frac{1}{G} \sum_{j=1}^{G} r\big(c_t^{(j)}\big), \qquad A^{(i)} = \frac{G}{G-1}\Big(r\big(c_t^{(i)}\big) - \bar{r}\Big)$$

- **Clipped, per-token surrogate** (GRPO-style) on **thought tokens** only

$$\mathscr{L}_{\text{clip}} = -\,\mathbb{E}\left[\frac{1}{|c_t^{(i)}|} \sum_u \min\Big(\rho_u^{(i)} A^{(i)}, \text{clip}\big(\rho_u^{(i)}; 1 - \epsilon_\ell, 1 + \epsilon_h\big) A^{(i)}\Big)\right]$$

with per-token importance ratios

$$\rho_u^{(i)} = \exp\Big(\log \pi_\theta\big(\ell_u^{(i)} \mid x_{<t}, \ell_{1:u-1}^{(i)}\big) - \log \pi_{\theta_{\text{old}}}\big(\ell_u^{(i)} \mid x_{<t}, \ell_{1:u-1}^{(i)}\big)\Big)$$

# Leverage Early Reasoning with "RLP"

**Q1** — Can RLP improve the reasoning ability of a base model without any prior task-specific tuning? (And better than Base + CPT?)
Does the gain of RLP survive even after SFT+RLVR?

**Q2** — Can RLP, when applied to an earlier pre-training checkpoint with 200B less tokens (to match the pretraining FLOPs), match the performance of a fully trained base model?

**Training Ingredients**

**Model**: Qwen3-1.7B-Base
**Checkpoint:** Final model
**Token Budget**: 1B Tokens
**CPT Data ($\mathscr{D}$):**
- General Pretraining Corpora (Nemotron-Pretraining-Dataset)

**Comparison with:**
- Base model
  ‣ **BASE**
  ‣ **BASE+POST**
- Token Matched: NTP-based CPT on base model
  ‣ **BASE+CPT**
  ‣ **BASE+CPT+POST**

**Training Ingredients**

**Model**: NVIDIA-Nemotron-Nano-12B-v2-Base
**Checkpoint:** Model trained till 19.8T tokens
**Token Budget**: 250M Tokens
**CPT Data ($\mathscr{D}$):**
- General Pretraining Corpora (Nemotron-Pretraining-Dataset)

**Comparison with:**
- FLOP Matched: Base model (20T)
  ‣ **BASE**
  ‣ **BASE+POST**

RLP with *200B* less tokens!

# RLP shows a significant improvement on Qwen3-1.7B-Base

**RLP** outperforms BASE by *+19%* and **CPT** by *+17%* on average across math and science benchmarks.

After identical SFT + RLVR, **RLP compounds its advantage:** *+8%* relative over BASE+Post; *+7%* relative over CPT+Post

Unlike NPT or continuous pretraining, RLP instills reasoning during pretraining itself
The early advantages compound through post training, giving models stronger and more robust reasoning capabilities.



Left chart — Average Accuracy by category (Math, Science, Science Pass@1[4], Overall)
Legend: Base, Base + CPT, Base + RLP

Right chart — Average Accuracy by category (Math, Science, Science Pass@1[4], Overall)
Legend: Base + Post, Base + CPT + Post, Base + RLP + Post

# RLP shows a significant improvement on Qwen3-1.7B-Base

Does the improvement sustain under compute equivalent baselines?

Even with flop matching training, **RLP outperforms CPT** (exposed to **35× more data**) by ***+14%*** on average!

$$T_{flop} = (n \times l_{seq} \times bs \times iters) + T_{inp}$$

$n$ = rollouts

$bs$ = batch size

$l_{seq}$ = sequence length

## Training Ingredients

**Model**: Qwen3-1.7B-Base

**Checkpoint:** Final model

**Token Matching ($T_{inp}$):** 170M Tokens

**Flop Matching ($T_{flop}$):** 6B Tokens

**CPT Data ($\mathscr{D}$):**
- Nemotron-CrossThink



Legend: Base | Base + CPT [TM] | Base + CPT [FM] | Base + RLP

# RLP scales with LLM size (12B) and frameworks

**RLP**, trained on **~200B fewer** tokens, achieves a **35% average gain**, with the largest boost in science (**+23% absolute**), showing robust cross-domain benefits at

After identical post-training, **RLP** outperforms BASE by **+2% absolute margin,** which has seen **~200B** more tokens during pretraining**.**

The benefits of RLP persist and even amplify when scaling to larger models and across architectures.

Efficiency meets intelligence—proof that reasoning, not sheer scale, drives progress.



Left chart — Average Accuracies:
- 37 (Science, Base+RLP)
- 35 (Science, Base)
- 33 (Science Pass@1[4], Base)
- 47 (Overall, Base)
- 20T (Base)
- 19.8T (Base+RLP)
- Categories: Math, Science, Science Pass@1[4], Overall
- Legend: Base, Base+RLP

Right chart — Average Accuracies:
- 59 (Science, Base + Post)
- Categories: Math, Science, Science Pass@1[4], Overall
- Legend: Base + Post, Base + RLP + Post

# Quiet-STaR: Language Models Can Teach Themselves to Think Before Speaking

**Eric Zelikman**
Stanford University

**Georges Harik**
Notbad AI Inc

**Yijia Shao**
Stanford University

**Varuna Jayasiri**
Notbad AI Inc

**Nick Haber**
Stanford University

**Noah D. Goodman**
Stanford University

# Reinforcement Pre-Training

**Qingxiu Dong**[*†‡]    **Li Dong**[*†]

**Yao Tang**[†]    **Tianzhu Ye**[†§]    **Yutao Sun**[†§]    **Zhifang Sui**[‡]    **Furu Wei**[†◇]
[†] Microsoft Research
[‡] Peking University
[§] Tsinghua University
https://aka.ms/GeneralAI

# Reinforcement Learning on Pre-Training Data

Siheng Li[1,3,*,†], Kejiao Li[1,†], Zenan Xu[1,†], Guanhua Huang[1], Evander Yang[1], Kun Li[1,3,*],
Haoyuan Wu[1], Jiajia Wu[1], Zihao Zheng[1], Chenchen Zhang[1], Kun Shi[1], Kyrierl Deng[1], Qi Yi[1],
Ruibin Xiong[1], Tingqiang Xu[1,*], Yuhao Jiang[1], Jianfeng Yan[1], Yuyuan Zeng[1], Guanghui Xu[1],
Jinbao Xue[2], Zhijiang Xu[2], Zheng Fang[2], Shuai Li[2], Qibin Liu[2], Xiaoxue Li[2], Zhuoyu Li[2],
Yangyu Tao[2], Fei Gao[2], Cheng Jiang[2], Bo Chao Wang[2], Kai Liu[2], Jianchen Zhu[2],
Wai Lam[3], Bo Zhou[1,‡], Di Wang[1]
[1]LLM Department, Tencent    [2]HunYuan Infra Team
[3]The Chinese University of Hong Kong
✉ chaysezhou@tencent.com

# RLP: The Core Idea

Chain-of-thought as an exploratory action during pretraining

```
Reward(thought) = log P(next token | context + thought) − log P(next
                             token | context alone)
```

## 💡 Verifier-Free

No external judge needed. The reward signal comes directly from information gain on the next token — applicable to any ordinary pretraining text.

## ⚡ Dense Reward

Position-wise credit assigned at every token where reasoning helps. Unlike sparse post-training rewards, RLP provides continuous feedback throughout training.

## 📈 Scalable

Works on the same massive pretraining corpora used for standard training. No specialized datasets required — reasoning emerges from plain text.

RLP bridges next-token prediction and the emergence of useful chain-of-thought reasoning — all during pretraining.

# Front Loading Reasoning

FLR systematically injects reasoning-style data at different phases of training—pretraining, SFT, RL—while varying its diversity, quantity, and quality.

**FLR**

RQ1: Is the inclusion of reasoning data in pretraining beneficial for the base model?
Lesson I: Yes, **highly beneficial**, and the benefit scales with **diversity and quantity**!

Using reasoning data in pretraining yields an absolute **+16% average improvement** over the no-reason baseline.

Pretraining gains improve with scale and diversity, driving math, science and code improvements!



Left chart — Average Accuracies:
- Average: No-Reason Base 53, Reason Base 61
- MATH: No-Reason Base 47, Reason Base 67
- SCIENCE: No-Reason Base 47, Reason Base 52
- CODE: No-Reason Base 41, Reason Base 49
- GPR

Legend: No-Reason Base, Reason Base

Right chart — Average Accuracies:
- Math: Reason Base [SHQ] 52.6, Reason Base [LDQ] 75.6, Reason Base [LMQ] 72.4
- Science: 46.9, 54.4, 54.5
- Code: 44.3, 49.9, 52.6
- GPR
- Overall: 55.0, 64.1, 64.1

Legend: Reason Base [SHQ], Reason Base [LDQ], Reason Base [LMQ]

# Front Loading Reasoning

FLR systematically injects reasoning-style data at different phases of training—pretraining, SFT, RL—while varying its diversity, quantity, and quality.

**FLR**

RQ3: Can a model with no-reason pretraining "catch up" by getting more SFT compute?
Lesson IV: **No**! Early reasoning builds an **irreplaceable** foundation that SFT can't replicate



Doubling SFT data for the baseline improves performance by **+4.09%**, but it still falls short of even the weakest reasoning-pretrained model (**+3.32%**).

Legend:
- No-Reason Base + SFT
- No-Reason Base + SFT (2x)
- Reason Base [SHQ] + SFT

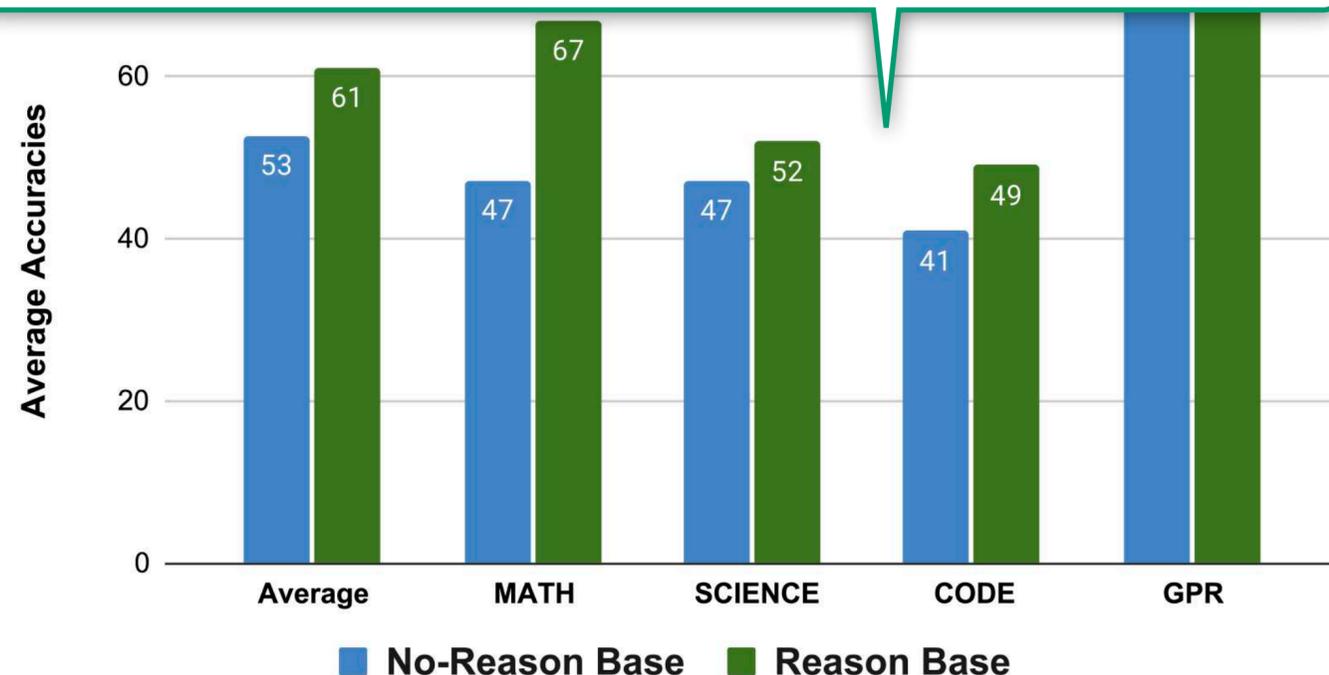| Category | No-Reason Base + SFT | No-Reason Base + SFT (2x) | Reason Base [SHQ] + SFT |
|---|---|---|---|
| Math | 42.8 | 48.05 | 50.52 |
| Science | 35.8 | 40.69 | 40.00 |
| Code | 10.5 | 14.60 | 24.76 |
| Instruction | 30.6 | 32.70 | 34.06 |
| Overall | 29.9 | 34.01 | 37.33 |

# Front Loading Reasoning

FLR systematically injects reasoning-style data at different phases of training—pretraining, SFT, RL —while varying its diversity, quantity, and quality.
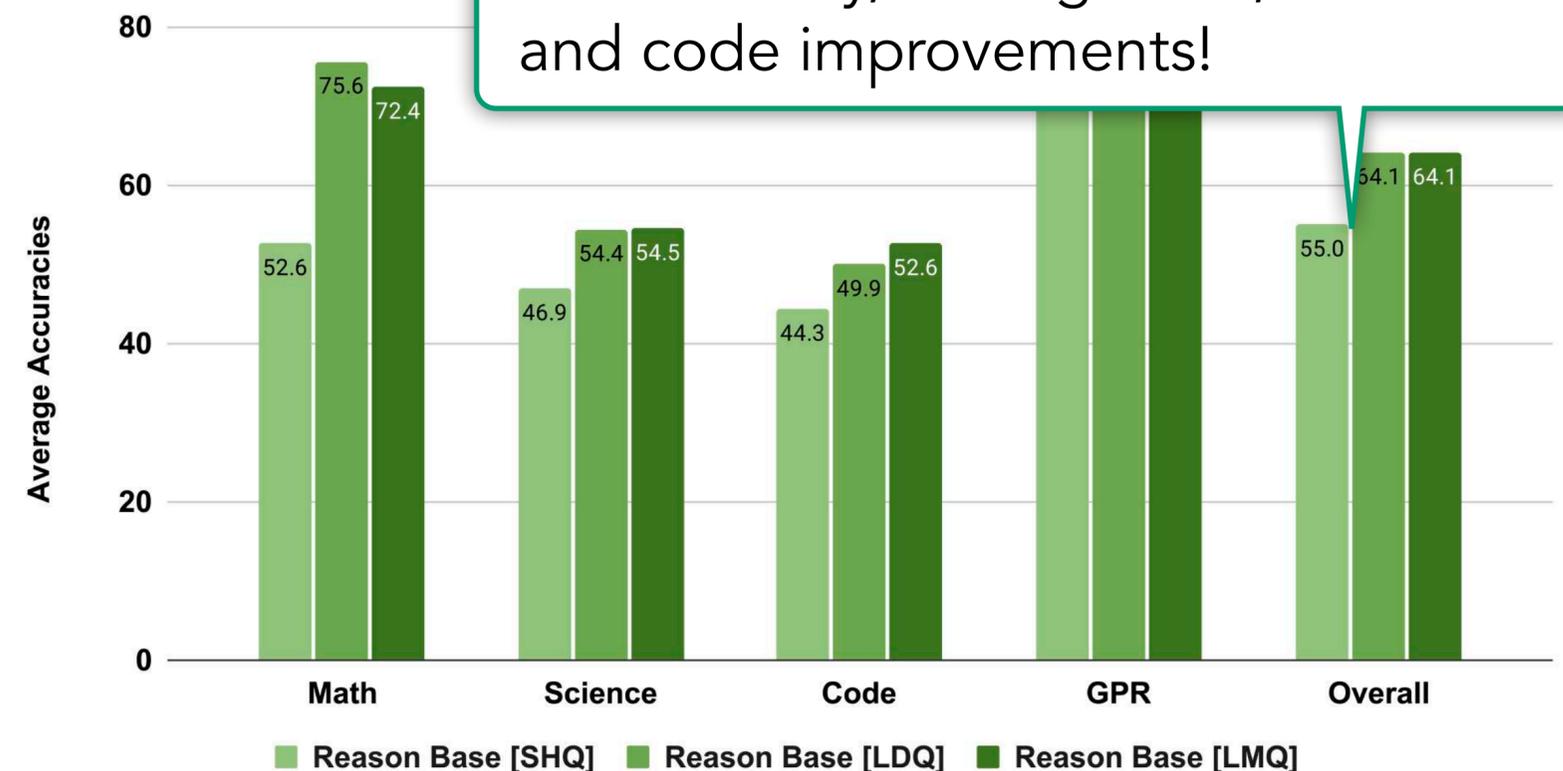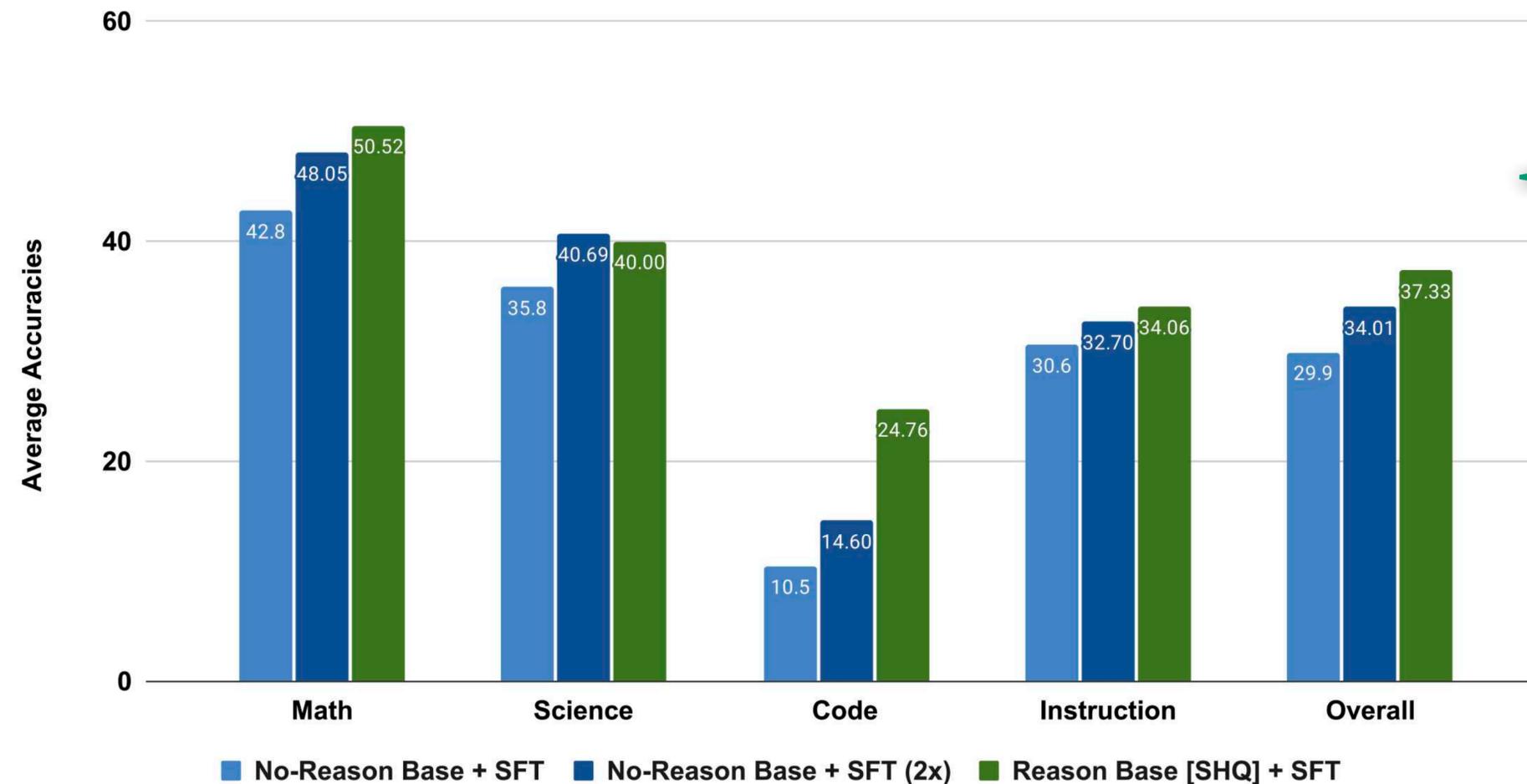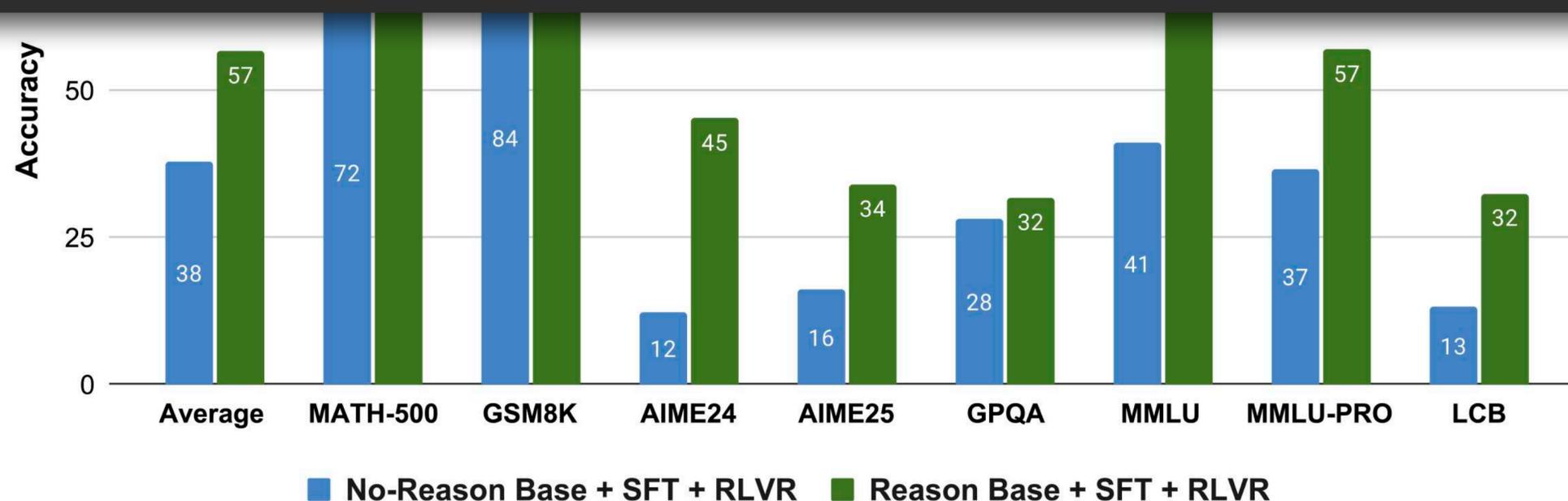
**FLR**

RQ5: Does the gain matter after heavy RLVR?

Lesson VI: **Yes**! "Front-loading" reasoning data creates a **durable, compounding gain**

💡 Front-loading reasoning data into pretraining creates a durable, compounding advantage.
💡 High-quality pretraining data can have a latent effect unlocked by SFT
💡 The optimal data strategy is asymmetric: prioritize diversity in pretraining and quality in SFT.

Reason-Base model finished with a stunning **+19%** average lead. On the most difficult competition math problems (AIME), this advantage ballooned to a **+39.3%** gain.



Accuracy

| | No-Reason Base + SFT + RLVR | Reason Base + SFT + RLVR |
|---|---|---|
| Average | 38 | 57 |
| MATH-500 | 72 | |
| GSM8K | 84 | |
| AIME24 | 12 | 45 |
| AIME25 | 16 | 34 |
| GPQA | 28 | 32 |
| MMLU | 41 | |
| MMLU-PRO | 37 | 57 |
| LCB | 13 | 32 |

■ No-Reason Base + SFT + RLVR   ■ Reason Base + SFT + RLVR

# Concluding Remarks

*David vs. Goliath: the Art of Scaling Intelligence
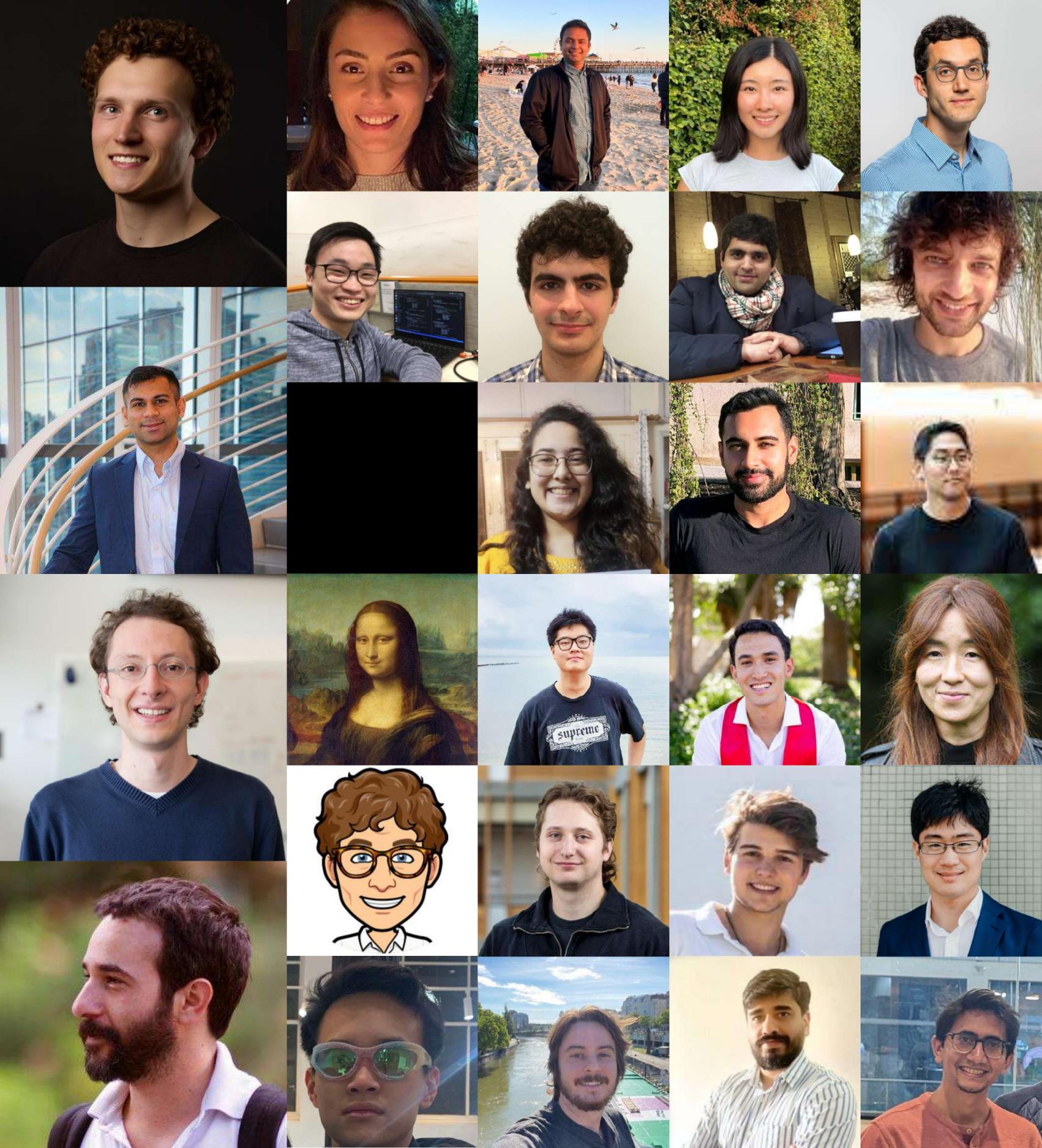in the Era of Extreme-Scale Neural Models*

Three components to innovate:
Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏
~    implicit: open-science & open-source community
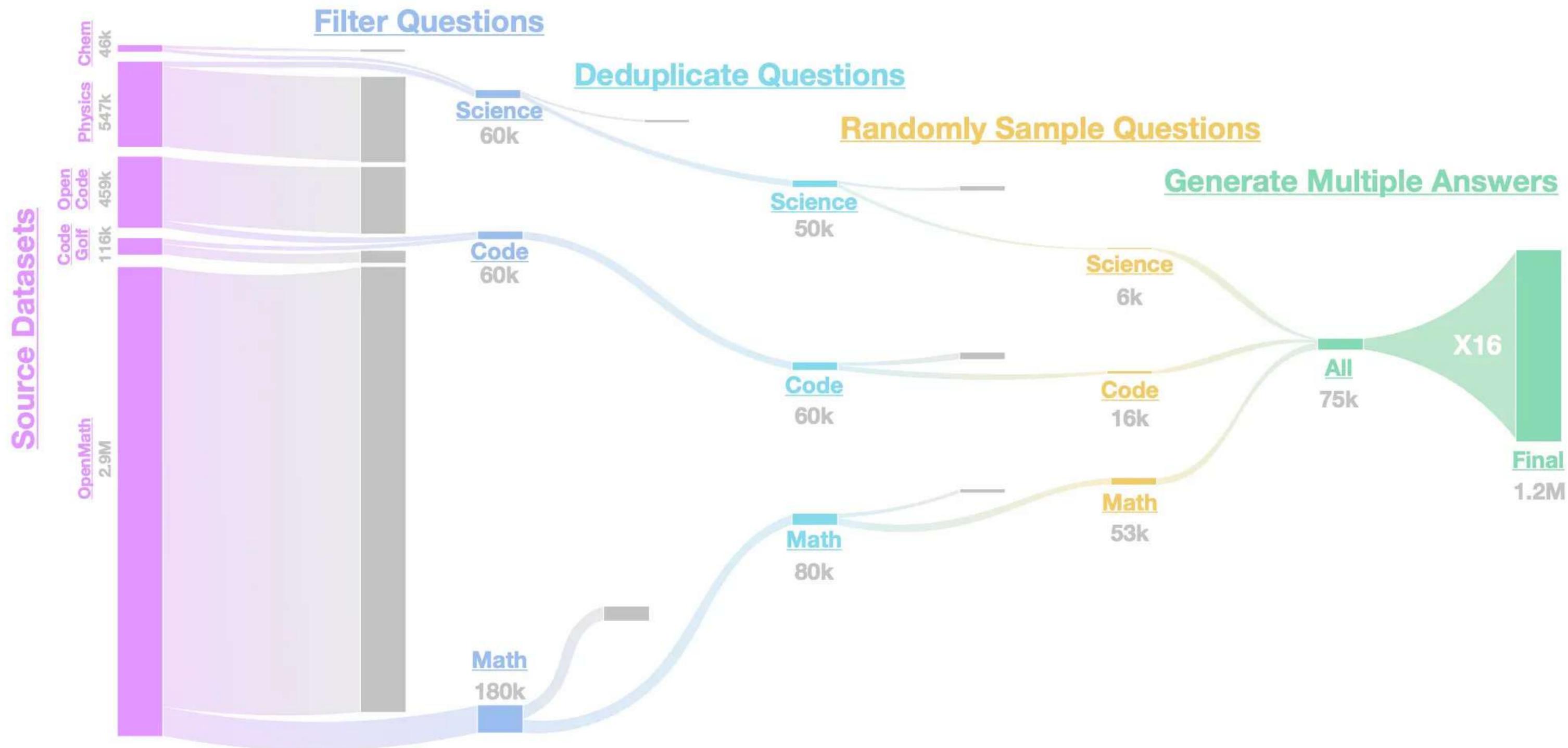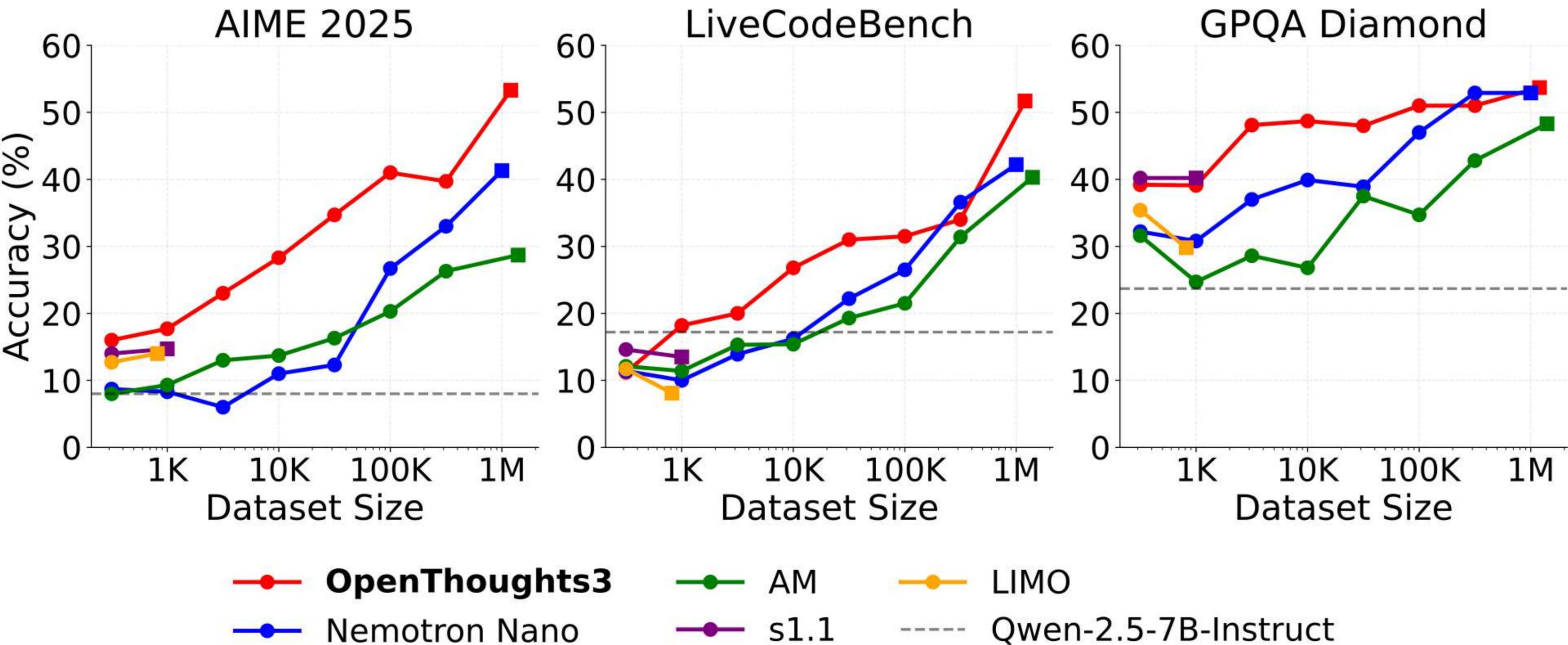~    explicit: cross-institutional, cross-boarder collaborations

# The OpenThoughts Team

Etash Guha*, Ryan Marten*, Sedrick Keh*, Negin Raoof*, Georgios Smyrnis*, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su, Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Stutee Acharya, Jeffrey Li, Achal Dave, Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy, Alexandros G. Dimakis, Ludwig Schmidt*

And several others with ongoing contributions

# OpenThoughts3-1M

# OpenThoughts3 is the SOTA reasoning dataset recipe



Figure: Three line charts comparing accuracy (%) versus dataset size for AIME 2025, LiveCodeBench, and GPQA Diamond benchmarks across OpenThoughts3, Nemotron Nano, AM, s1.1, LIMO, and Qwen-2.5-7B-Instruct.

| Benchmark | OpenThoughts3-7B | DS-R1-Qwen-7B | NemoNano-1M | AM-1.4M | OpenR1-Distill-7B | *SFT vs. RL* Nemotron-Nano-8B | AceReason-7B | Skywork-7B | *Base Model* Qwen2.5-7B-Instruct |
|---|---|---|---|---|---|---|---|---|---|
| Base Model | 🔷 | 🔷ᴹ | 🔷 | 🔷 | 🔷ᴹ | ∞ | 🐳 | 🐳 | N/A |
| Train Size | 1.2M | 800K | 1M | 1.4M | 350K | 3.9M | 57K | 119K | N/A |
| Method | SFT | SFT | SFT | SFT | SFT | SFT/RL | RL | RL | N/A |
| Trained by us | Yes | No | Yes | Yes | No | No | No | No | N/A |
| Open Data | 🟢 | 🔴 | 🟢 | 🟢 | 🟢 | 🟢 | 🟡 | 🟡 | N/A |
| Average | **55.3** | 42.9 | 47.3 | 42.1 | 47.2 | 53.2 | 52.9 | 51.6 | 24.0 |
| *Math* AIME24 | **69.0** | 51.3 | 55.0 | 28.3 | 57.7 | 62.0 | **71.0** | 68.3 | 15.0 |
| *Math* AMC23 | **93.5** | 92.0 | 87.0 | 82.2 | 87.0 | **94.0** | 93.8 | 91.0 | 53.0 |
| *Math* MATH500 | 90.0 | 88.0 | 86.8 | 87.4 | 88.0 | 89.4 | 89.8 | **90.2** | 70.8 |
| *Code* CodeElo | 31.0 | 19.9 | 28.6 | 21.0 | 30.1 | 30.9 | 32.9 | **37.0** | 5.5 |
| *Code* LCB 05/23-05/24 | 64.5 | 48.7 | 58.0 | 54.5 | 37.9 | **68.0** | 60.5 | 60.4 | 36.2 |
| *Code* CodeForces | **32.2** | 21.1 | 28.3 | 24.8 | 29.3 | **32.9** | 30.9 | **32.5** | 10.2 |
| *Sci* GPQA-D | 53.7 | 33.2 | 52.9 | 48.3 | **58.9** | 52.9 | 52.9 | 50.2 | 24.6 |
| *Sci* JEEBench | **72.4** | 50.4 | 61.0 | 61.1 | 68.7 | 70.7 | 64.3 | 55.3 | 33.9 |
| *Held Out* HMMT 02/25 | **42.7** | 25.0 | 24.7 | 19.0 | 25.7 | 26.7 | 33.3 | 32.7 | 2.0 |
| *Held Out* HLE MCQ | 10.2 | 12.4 | 2.1 | 9.5 | 12.4 | 12.0 | 10.9 | 10.7 | **12.7** |
| *Held Out* AIME25 | **53.3** | 38.0 | 41.3 | 28.7 | 39.7 | 48.0 | 50.7 | 47.3 | 8.0 |
| *Held Out* LCB 06/24-01/25 | **51.7** | 34.5 | 42.2 | 40.3 | 30.7 | **50.9** | 44.3 | 43.8 | 16.3 |

Super effortful SFT can win over effortful RL

# Open Thoughts
## DATA RECIPES FOR REASONING MODELS

**Etash Guha**[*,1,2], **Ryan Marten**[*,3], **Sedrick Keh**[*,4], **Negin Raoof**[*,5], **Georgios Smyrnis**[*,6],
Hritik Bansal[ζ,7], Marianna Nezhurina[ζ,8,9,16], Jean Mercat[ζ,4], Trung Vu[ζ,3], Zayne Sprague[ζ,6],
Ashima Suvarna[7], Benjamin Feuer[10], Liangyu Chen[1], Zaid Khan[11], Eric Frankel[2],
Sachin Grover[12], Caroline Choi[1], Niklas Muennighoff[1], Shiye Su[1], Wanjia Zhao[1], John Yang[1],
Shreyas Pimpalgaonkar[3], Kartik Sharma[3], Charlie Cheng-Jie Ji[3], Yichuan Deng[2],
Sarah Pratt[2], Vivek Ramanujan[2], Jon Saad-Falcon[1], Jeffrey Li[2], Achal Dave, Alon Albalak[13],
Kushal Arora[4], Blake Wulfe[4], Chinmay Hegde[10], Greg Durrett[6], Sewoong Oh[2],
Mohit Bansal[11], Saadia Gabriel[7], Aditya Grover[7], Kai-Wei Chang[7], Vaishaal Shankar,
Aaron Gokaslan[14], Mike A. Merrill[1], Tatsunori Hashimoto[1], Yejin Choi[1],
Jenia Jitsev[8,9,16], Reinhard Heckel[15], Maheswaran Sathiamoorthy[3],
Alexandros G. Dimakis[†,3,5], Ludwig Schmidt[†,1]

[1]Stanford University, [2]University of Washington, [3]BespokeLabs.ai, [4]Toyota Research Institute
[5]UC Berkeley, [6]UT Austin, [7]UCLA, [8]JSC, [9]LAION, [10]NYU, [11]UNC Chapel Hill
[12]ASU, [13]Lila Sciences, [14]Cornell Tech [15]TUM [16]Open-Ψ (Open-Sci) Collective

| | AceReason-7B | Skywork-7B | Base Model Qwen2.5-7B-Instruct |
|---|---|---|---|
| | 🐋 | 🐋 | N/A |
| | 57K | 119K | N/A |
| | RL | RL | N/A |
| | No | No | N/A |
| | 🟡 | 🟡 | N/A |
| | 52.9 | 51.6 | 24.0 |
| | **71.0** | **68.3** | 15.0 |
| | **93.8** | 91.0 | 53.0 |
| | 89.8 | **90.2** | 70.8 |
| | 32.9 | **37.0** | 5.5 |
| | 60.5 | 60.4 | 36.2 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| *Held Out* | HMMT 02/25 | **42.7** | 25.0 | 24.7 | 19.0 | 25.7 | | | |
| | HLE MCQ | 10.2 | 12.4 | 2.1 | 9.5 | 12.4 | | | |
| | AIME25 | **53.3** | 38.0 | 41.3 | 28.7 | 39.7 | 48.0 | 47.3 | 8.0 |
| | LCB 06/24-01/25 | **51.7** | 34.5 | 42.2 | 40.3 | 30.7 | **50.9** | 44.3 | 43.8 | 16.3 |

Unconventional collaboration wins!

Super effortful SFT can win over effortful RL

# Concluding Thoughts

- Nothing is easy in life

- No pain no gain

- But everything is doable with efforts

- The power of collaboration

- Ego is the enemy

Three components to innovate:
Unconventional data 🔥
Unconventional algorithms 🚀
Unconventional collaboration 🌏

Chemistry between the base LLM and RL matters
Conclusions from effortless RL != effortful RL
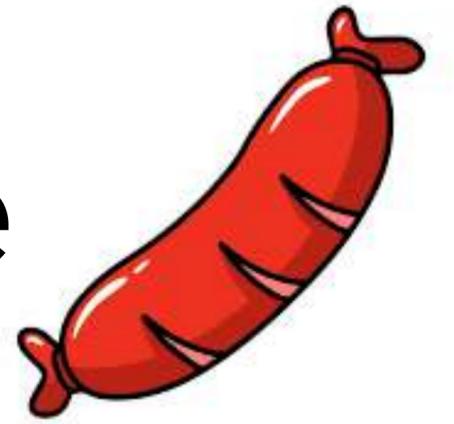
Conclusions from effortless SFT != effortful SFT

# NVIDIA LEADS OPEN SOURCE AI MOMENTUM AS CHINESE LABS CLOSE IN



October 23, 2024

*Considering repos above 10 space likes and 500 downloads*

# NVIDIA LEADS OPEN SOURCE AI MOMENTUM AS CHINESE LABS CLOSE IN

1. Small models are in high demand, with daily downloads hitting tens of millions (!!! 🔥) for top repositories
2. China's contribution to the open source ecosystem has skyrocketed, with Alibaba now surpassing Hugging Face 😮
3. NVIDIA has emerged as a rising star ⭐ (or a risen star? 😎), for open-source contributions 💚

# LLM 101: How sausages are made

"Fine-tuning"

**SFT**
=
sequential
fine-tuning

**RL**
=
Reinforcem
ent learning

Pre-
training

(On the internet data)

(On curated exam-style data)
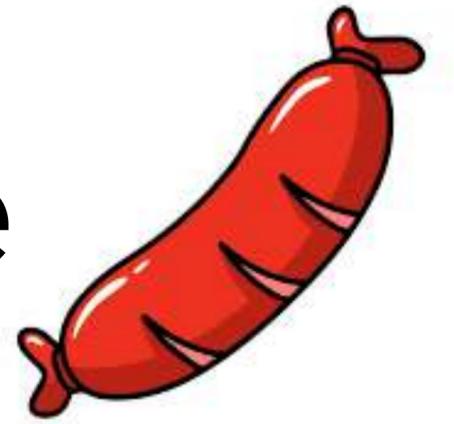
(On curated exam-style data)

Imitation learning

Exploration learning

# LLM 101: How sausages are made

"Fine-tuning"

It's all about data
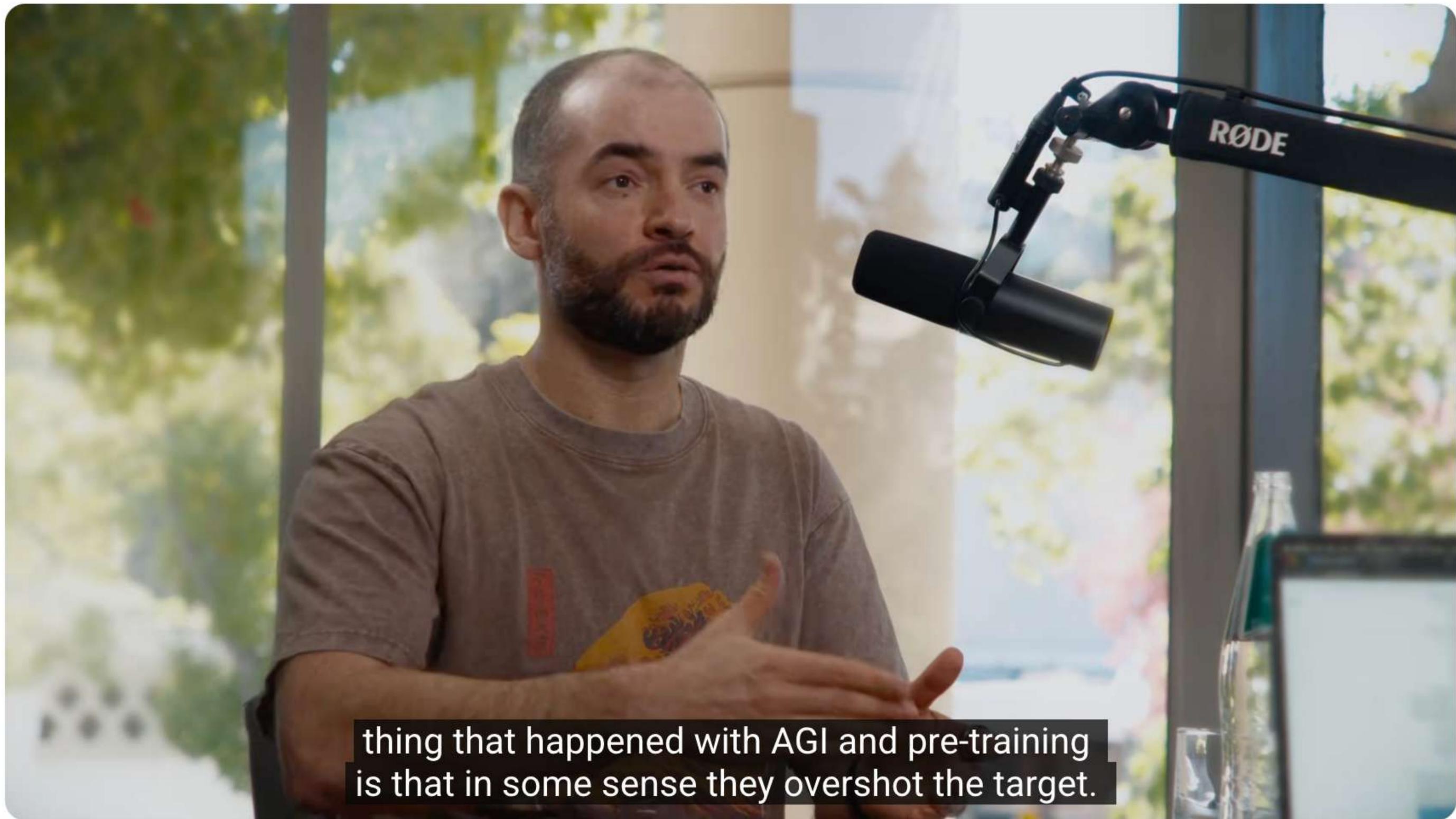
Use entirely of the internet data

Internet data not enough, so ask humans to write exam data

When even that's not enough, ask AI to synthesize data

=> Make OOD data as in-dist data via brute-force data synth!

Imitation learning          Exploration learning

thing that happened with AGI and pre-training is that in some sense they overshot the target.

**Ilya Sutskever – We're moving from the age of scaling to the age of research**

Supervised learning is not something that happens in nature.

**Richard Sutton – Father of RL thinks LLMs are a dead end**

Reinforcement Learning is terrible.

**Andrej Karpathy — "We're summoning ghosts, not building animals"**

The Universe of Knowledge
= The Universe of Synthetic Data

Knowledge in the **Internet** Data

Knowledge in the (Conventional) **Distillation**

Knowledge in the **Human/Experts'** Annotation

Knowledge in the Extreme Scaling of Reasoning

Knowledge in the Unconventional Simulation

# Universe of Know...

= The Universe of Synthetic Data

Discovery of new knowledge is the act of a scientist, even if the target domain is not for scientific fields per say

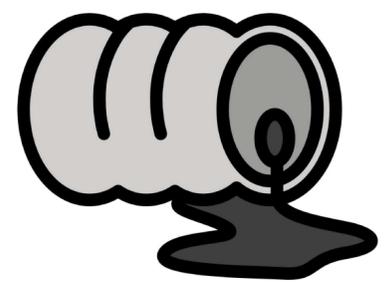Many real-world use cases require the capability to fill in this data gap

Those with a stronger synthetic data flywheel will win. It will be increasingly reasoning & compute-heavy, and often not be open.

Knowledge in the **Internet** Data

Knowledge in the (Conventional) **Distillation**

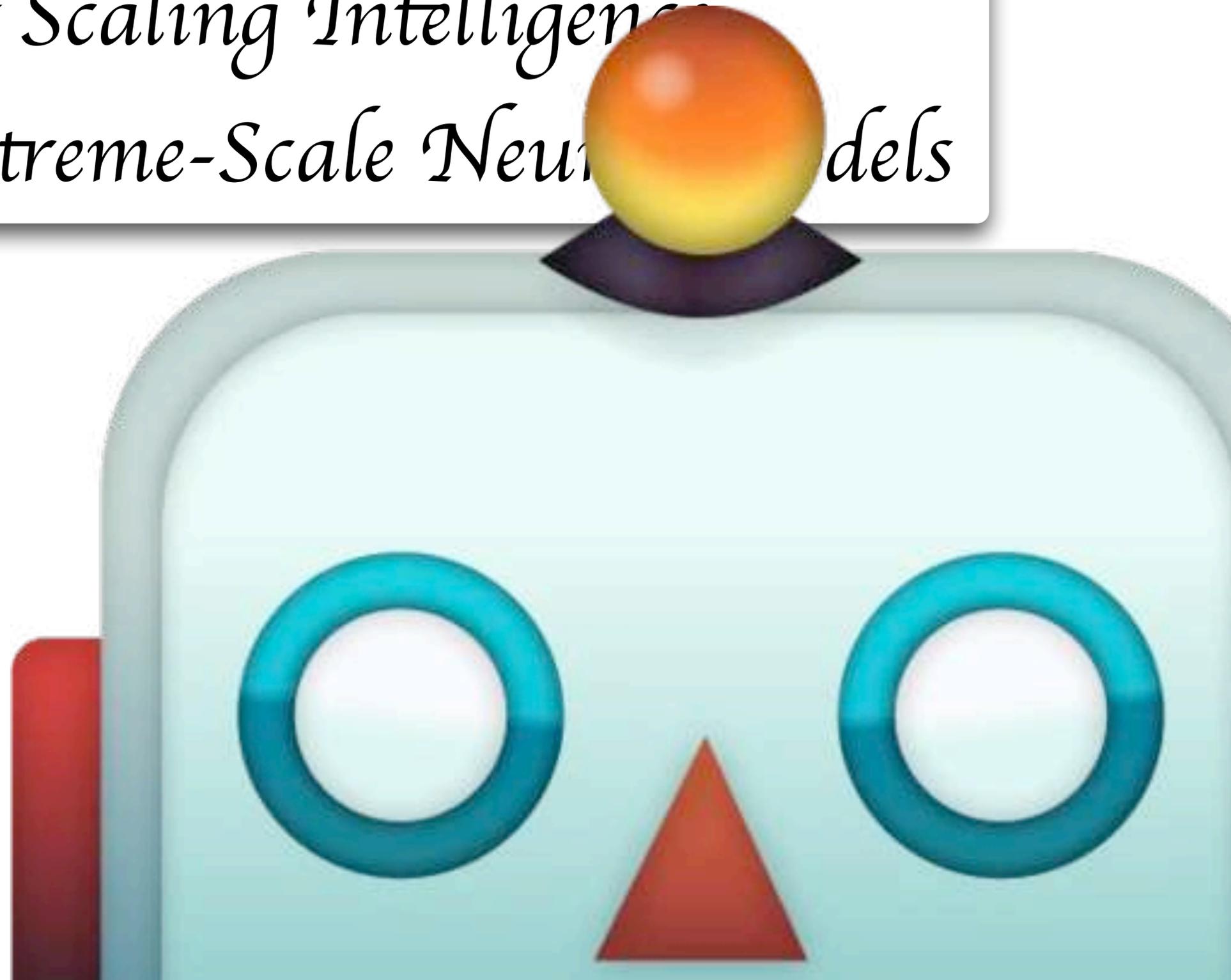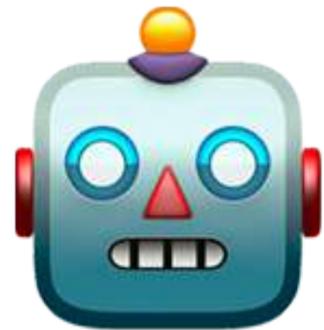Knowledge in the **Human/Experts'** Annotation

Knowledge in the Extreme Scaling of Reasoning

Knowledge in the Unconventional Simulation

# David vs. Goliath:
# the Art of Scaling Intelligence in the Era of Extreme-Scale Neural Models

# David vs. Goliath:
## the Art of Scaling Intelligence in the Era of Extreme-Scale Neural Models

AI mirrors human intelligence:

- Reasoning is often memorized knowledge
- Exploitation vs exploration trade offs
- Diverse examples/experiences enhance learning
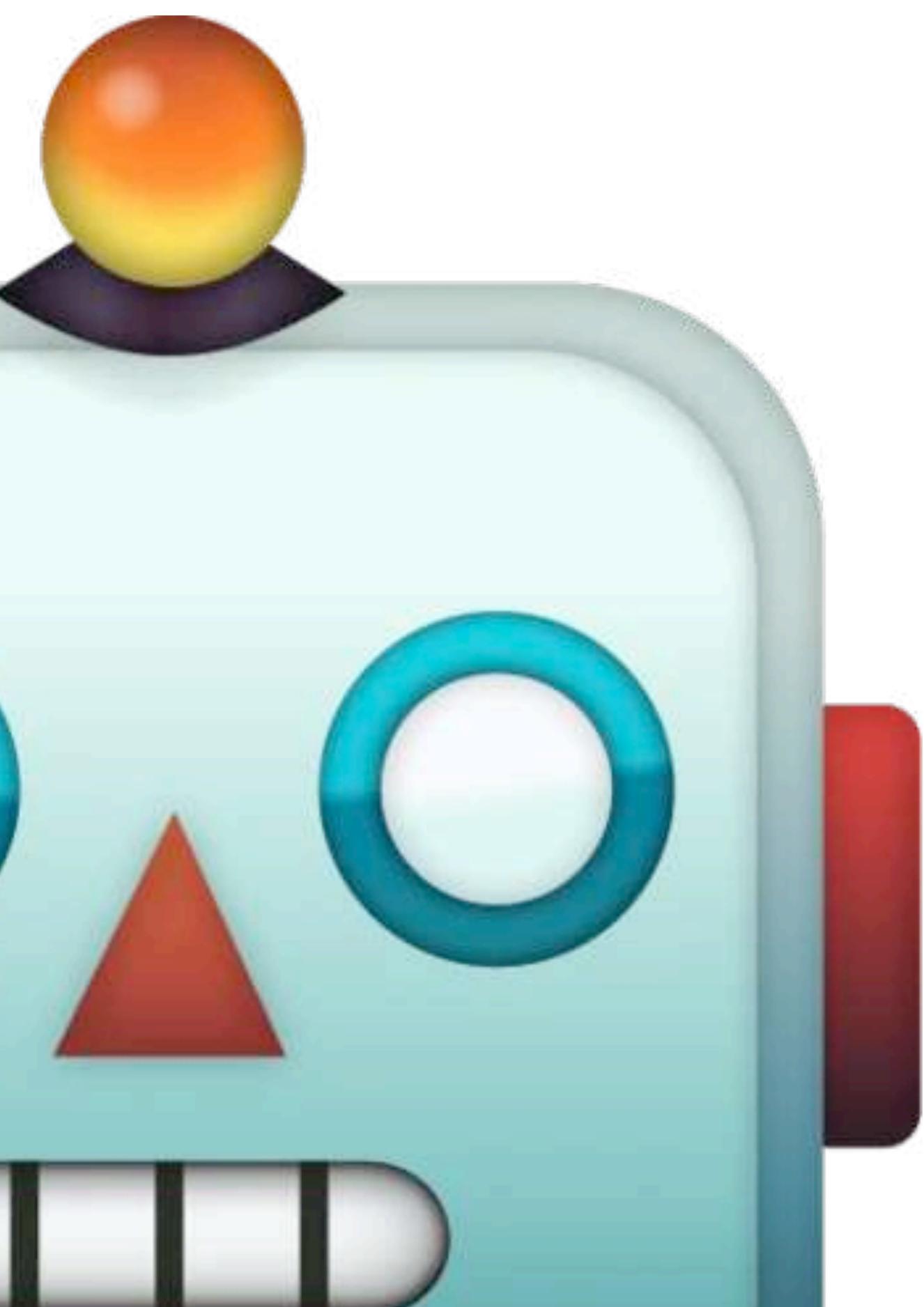
AI diverges from human intelligence:

- More data, compute, and memory
- Less abstraction and conceptualization
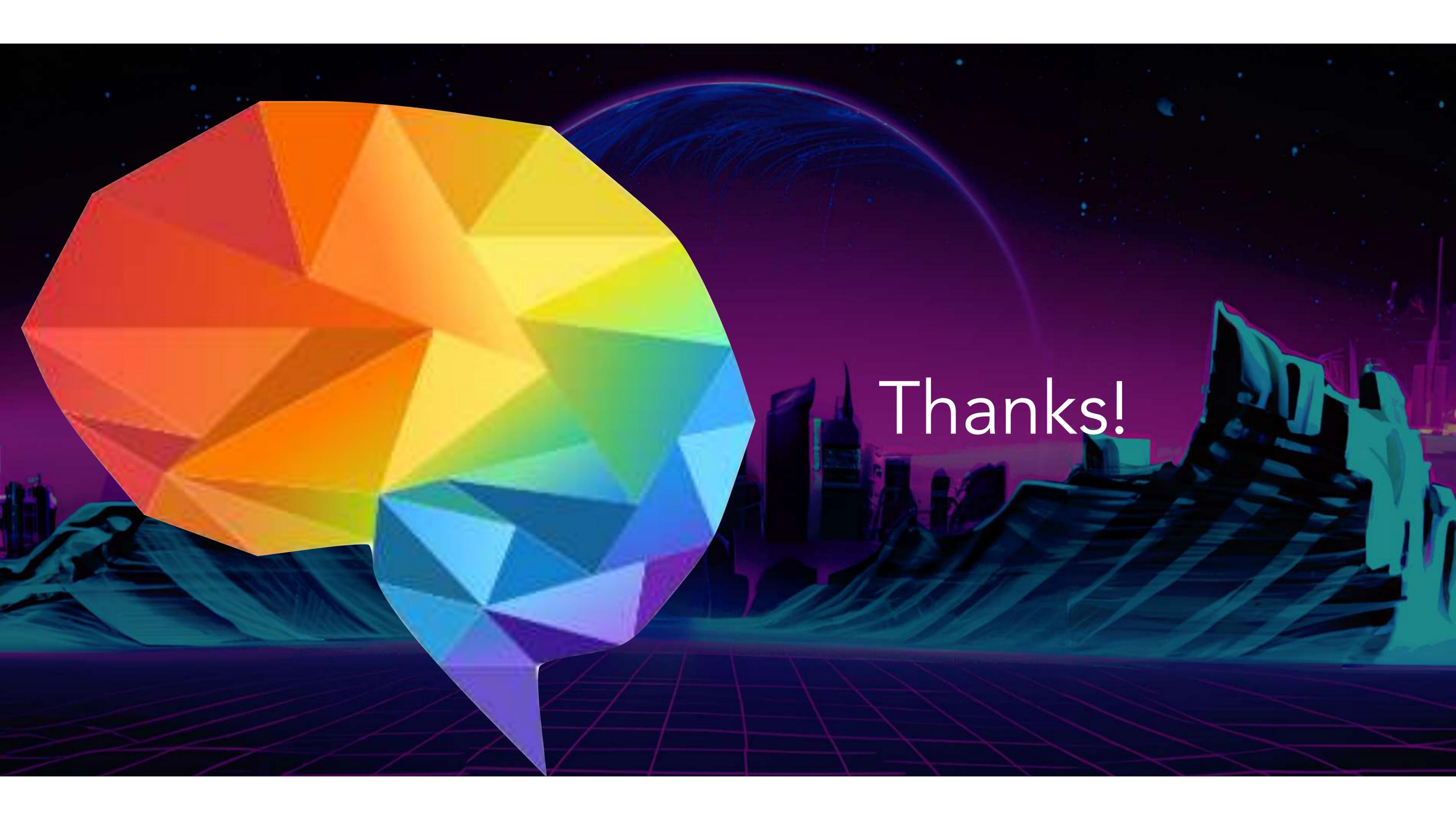- Clear separation between training and testing

How sausages are made:

👩‍🚀 Human's internet data, the artifact of human intelligence

👩‍🚀 Human's audacity to throw in $$$$$$$

👩‍🚀 Human's annotations at scale for alignment

👩‍🚀 Human's intuitions and insights on a lot of engineering details

Open Research Questions

- New theories of intelligences?

- New theories of knowledge and reasoning?

- That humans do not have a context window size of 1M tokens — is this a limitation or blessing?

Thanks!