# Large Language Model Post-Training Formulation and Algorithms

Ziniu Li

The Chinese University of Hong Kong, Shenzhen

2025-03-26

# Overview of This Talk
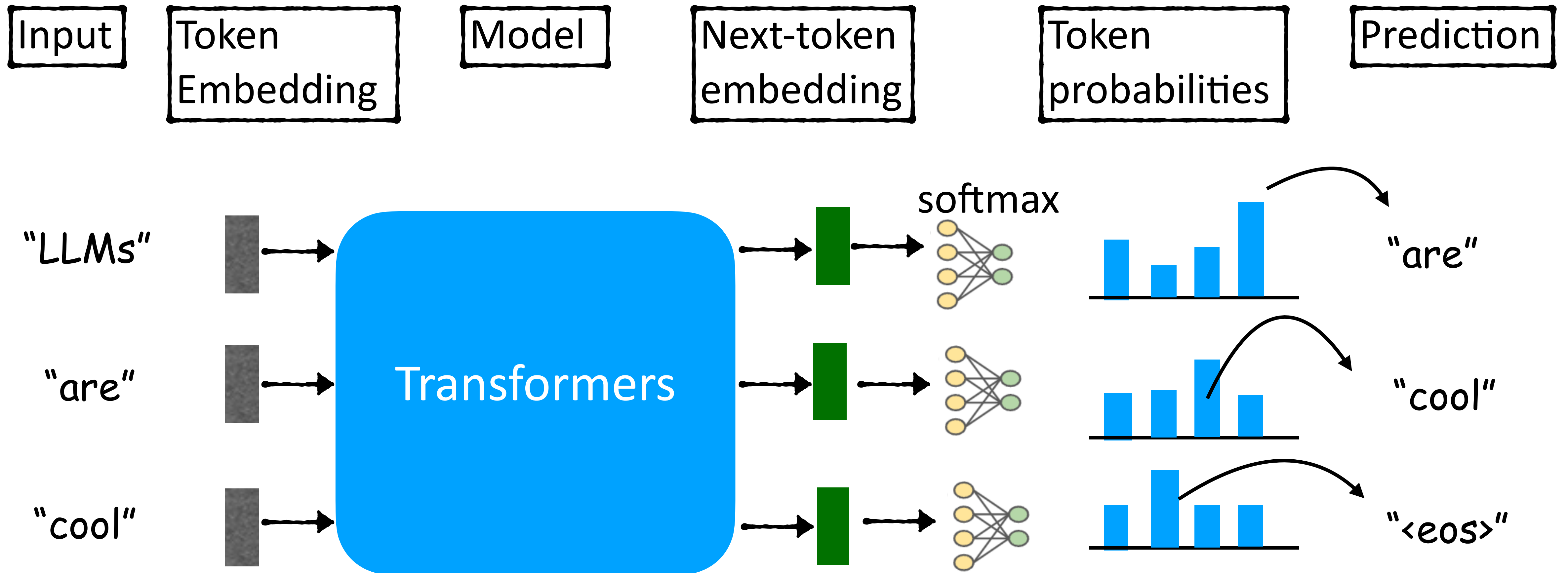
Evolution of Large Language Models

Formulation and Key Properties of LLM Training
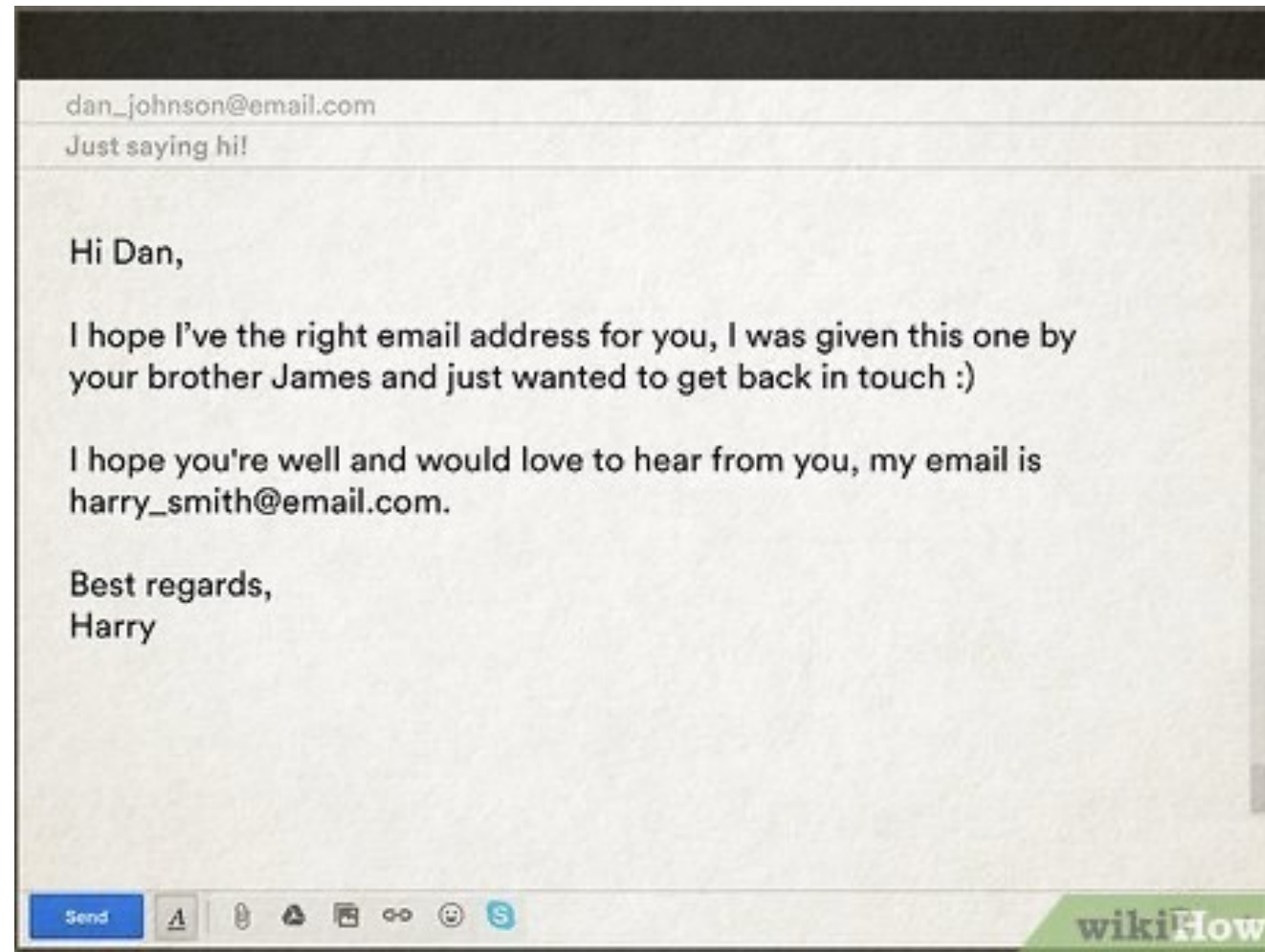
Our Research Contributions

Key Scientific Insights

# Part I: Overview of LLMs

# LLMs and Transformers

Input   Token Embedding   Model   Next-token embedding   Token probabilities   Prediction

"LLMs" → Transformers → softmax → "are"

"are" → Transformers → "cool"

"cool" → Transformers → "<eos>"

Transformers perform **next-token-prediction** and **token generation**

4

# Tasks that LLM can Solve



**Email Writing**



**Travel Plan**



**Summarization**



**Code Generation**

**Now, a single LLM can conduct all these functions**

# A Single Model for All Tasks.
# How can do this?

# LLM Training Framework

One can search "LLM Training Pipeline" and get the following figure:



But Why?

- ▸ What specific purpose does each training stage serve?
- ▸ Why do LLMs have to follow such training pipelines?

This talk provides some understanding and insights of LLM training

# LLM Pre-training

LLM Pre-training  = Transformers + Next-token-Prediction + **Textbook Data**

"Textbooks" can cover:

linguistics
world knowledge

common sense

math   coding

Ilya Sutskever
(Godfather of ChatGPT)

Next-token Prediction is enough for AGI

[https://www.youtube.com/watch?v=YEUclZdj_Sc]

"Textbook" teaches everything
(multi-task learning)

8

# Scaling Law

[Kaplan, Jared, et al. "Scaling laws for neural language models." *arXiv:2001.08361*.]



$$L = \frac{A}{D^\alpha} + \frac{B}{N^\beta} + L_0$$

L: Loss
D: dataset size
N: number of parameters
A, B: constants; $L_0$: irreducible loss

dataset size

model size

# Pre-Training Hero: GPT-3



GPT-3 is trained unsupervisedly but can beat zero-shot SOTA

GPT-3 can conduct few-shot (in-context) learning and beat fine-tuned SOTA

[Brown, Tom, et al. "Language models are few-shot learners." *Advances in neural information processing systems* 33 (2020): 1877-1901.]

# From Classification to Token Generation

- GPT-3 was primarily evaluated on **classification** tasks in 2020

- **Token generation**:

Cannot understand user's question

Give me a simple recipe for chili

I don't have a simple recipe for chili. I have a complicated one. It's not that I don't like simple chili recipes. It's just that I don't know how to make them. I've never been able to make a simple chili recipe. I've tried, but I always end up with something that's too spicy or too bland or too watery or too thick or too thin or too something.

So, I've given up on simple chili recipes. I've given up on trying to make a simple chili recipe. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated.

I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying
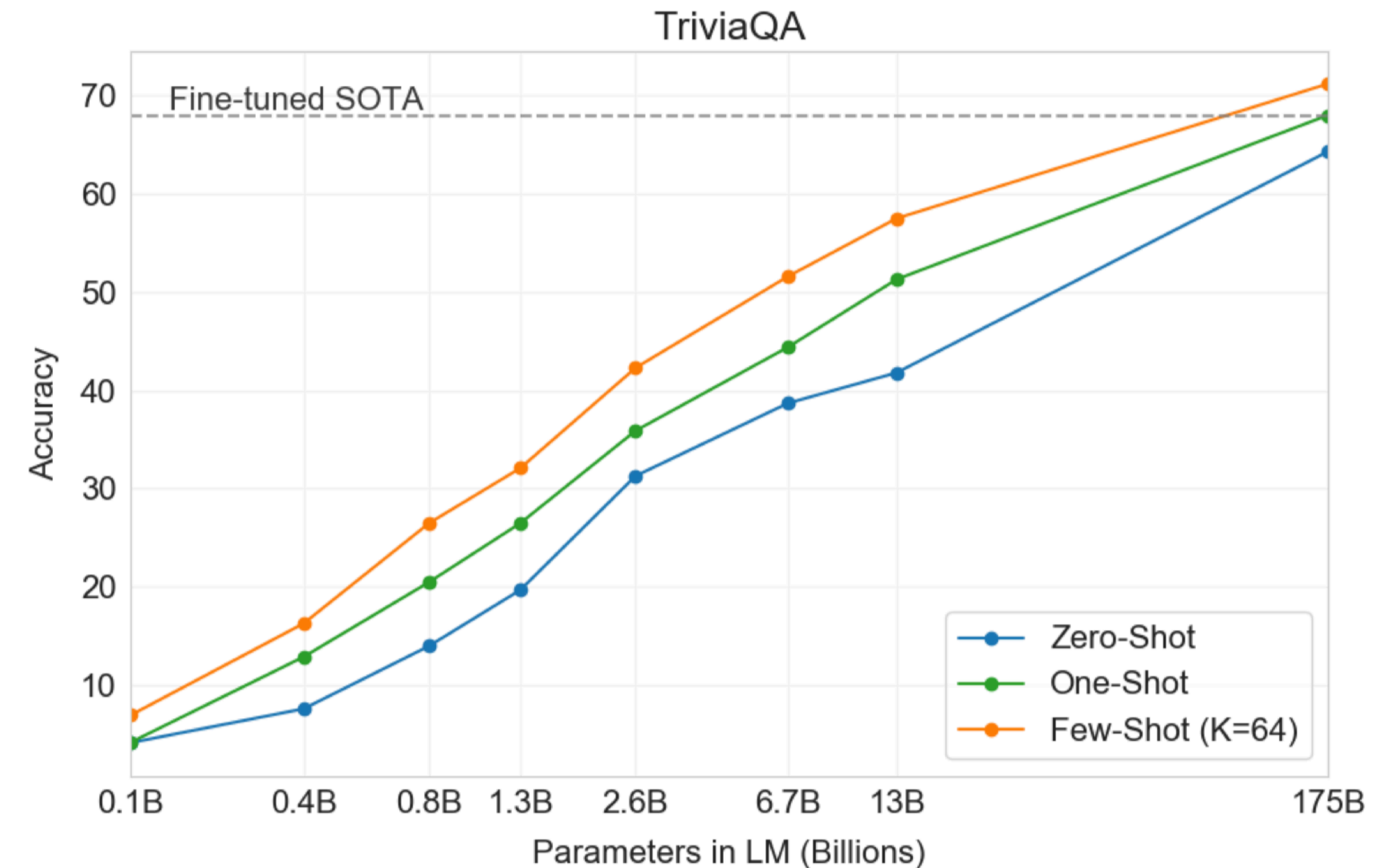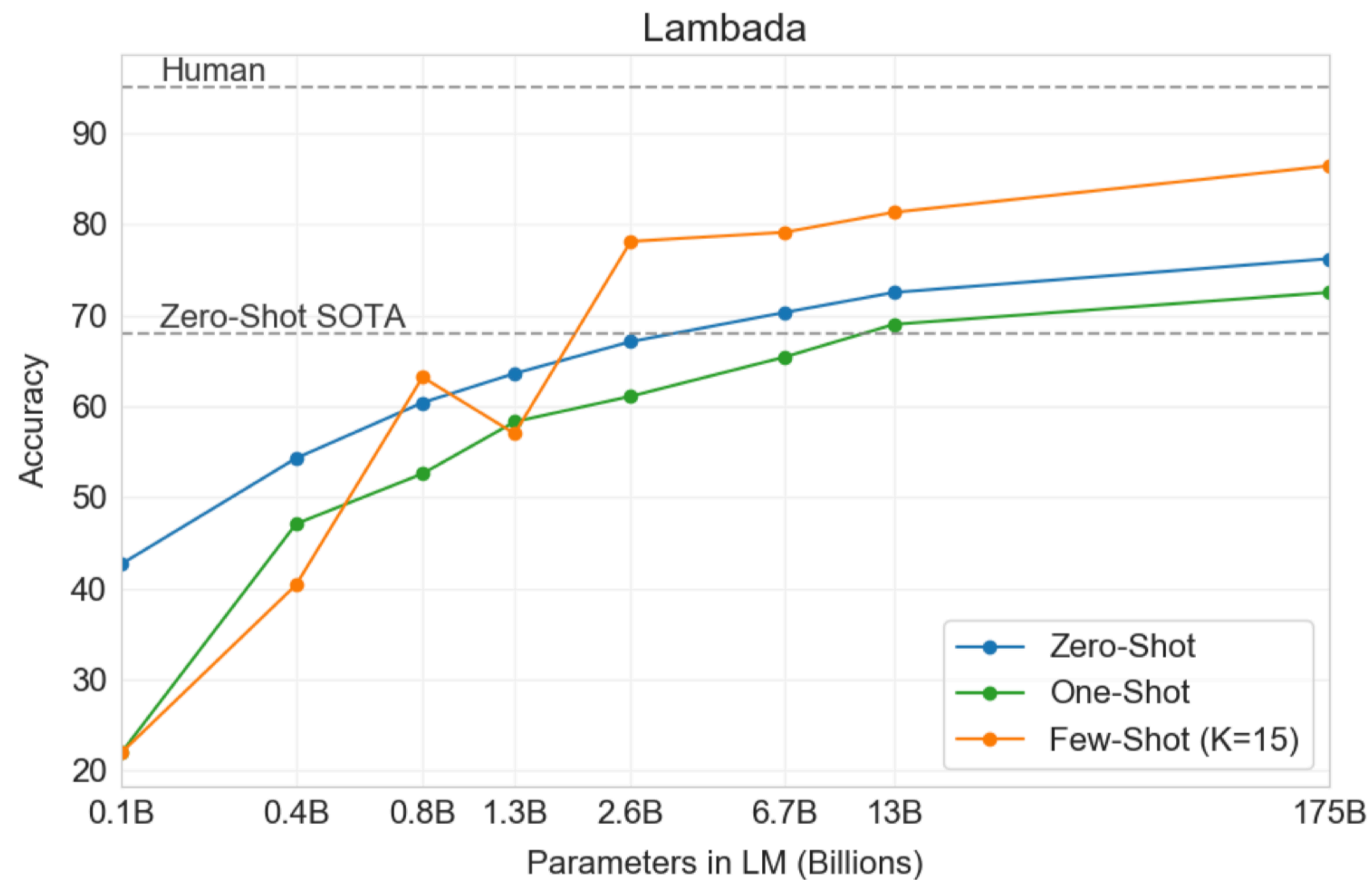
[ChatGPT and The Art of Post-Training. Barret Zoph & John Schulman. https://docs.google.com/presentation/d/11KWCKUORnPpVMSY6vXgBeFSWo7fJcuGQ9yuR6vC1pzE/edit#slide=id.p]

Repetitive responses

Pre-trained LLMs (including GPT-3):
- Possess extensive knowledge and generate coherent text
- Fails to understanding user intent and solving tasks effectively

11

make a simple chili recipe. I've tried, but I always end up with something that's too spicy or too bland or too watery or too thick or too thir

So, I've given up on simple chili recipes. I've given up on trying to make a simple chili recipe. I've given up on trying to make a chili recipe t trying to make a chili recipe that's not complicated.

I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I' recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili re up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given u not complicated. I've given up on trying to make a chili recipe that's not complicated. I've given up on trying to make a chili recipe that's n

Goal:

► Post-trained LLM:

User: give me a simple recipe for chili
Assistant: Here is a simple recipe for chili:

Ingredients:
1 pound ground beef
1 large onion, chopped
1 green bell pepper, chopped
2 (8 ounce) cans tomato sauce
1 (14.5 ounce) can diced tomatoes
1 (15 ounce) can kidney beans, drained and rinsed
2 tablespoons chili powder
1 tablespoon ground cumin
1 teaspoon garlic powder
1/2 teaspoon salt
1/2 teaspoon black pepper

Instructions:
1. In a large pot, brown ground beef over medium heat until cooked through, about 10 minutes. Drain excess fat.

[ChatGPT and The Art of Post-Training. Barret Zoph & John Schulman. https://docs.google.com/presentation/d/11KWCKUORnPpVMSY6vXgBeFSWo7fJcuGQ9yuR6vC1pzE/edit#slide=id.p]

# Summary of LLM Training Pipeline

**Pre-training**

**Post-training**

Knowledge Acquisition

Ability Reinforcement

# Post-Training Techniques



Figure is from "Weak-to-strong generalization: Eliciting strong capabilities with weak supervision."

14

# Supervised Fine-tuning

Label

Teacher    LLM

Objective $\max_{\theta} \mathbb{E}_{y \sim p(\cdot|x)}[\log f_{\theta}(y\,|\,x)]$

$x$: prompt     $y$: response/completion (label)

$p$: data distribution (from teacher)

$f_{\theta}$: distribution of LLM

SFT Data Example

Prompt
> Q: Can Geoffrey Hinton have a conversation with George Washington?

Label
> A: The answer is No because Geoffrey Hinton was born in 1947, while [...]

LLMs learn to **understand** the **question** (task) and **provide** relevant **answers**

15

# Reinforcement Learning



Reward

Teacher    LLM

Response

Objective

$$\max_{\theta} \mathbb{E}_{y \sim f_{\theta}(\cdot | x)}[r(x, y)]$$

Framework

Generation    Verification

RL Data Example

Prompt

Q: How many 'r' in strawberry?

LLM Response

A: There is one 'r' in 'stra' and another 'r' in 'berry', so the answer is 2

Teacher Feedback

Reward = -1

LLMs learn to **correct mistakes** and **enhance confidence** in answering questions

# Discussion

🤔    Why is pre-training necessary?  Why not proceed directly to post-training?

💡
- Knowledge density is **sparse** in post-training data (but rich in pre-training)
- LLMs with post-training solely **cannot generalize** well

🤔    Why implement SFT before reinforcement learning?

💡
- Pre-trained LLM outputs **lack good format** for reliable RL evaluation
- SFT establishes essential **response formatting** that enables RL optimization

# Part II: Preserving Output Diversity in Supervised Fine-Turning

# Revisiting SFT

SFT aims to align pre-trained model outputs to RL/human-preferred **format**
(outputs that are easy to **read**, **interpret**, and **verify**)

Response Space

Pre-trained LLM ▸— Poorly Formatted

Practice of SFT ▸— Well formatted ← 😞 Lose of diversity

😊 Goal of SFT ▸— Well formatted

(No diversity reduction)

# Output Diversity

**Question:** Marissa is hiking a 12-mile trail. She took 1 hour to walk the first 4 miles, then another hour to walk the next two miles. If she wants her average speed to be 4 miles per hour, what speed (in miles per hour) does she need to walk the remaining distance?

**Answer: 6**



**Monotonous Responses**

**Diverse Responses**

Greater **Diversity** Leads to Exploration of Better Solutions

# SFT Reduces Model Output Diversity

**#1**

Prompt | Give me a single-digit number

Pre-trained LLM

Response Distribution



"near uniform"

Pre-trained LLM + SFT



"biased toward 7"

[O'Mahony, Laura, et al. "Attributing mode collapse in the fine-tuning of large language models." *ICLR 2024 Workshop*. 2024.]

**#2**

Output Diversity Statistics

Output Diversity



■ Pre-training   ■ SFT

SFT reduces diversity by ~20%

# Related Issue: Model Homogenization toward GPT-4

- ► "Small" companies use GPT-4 outputs as SFT data to fine-tune their models
- ► Fine-tuned models follow GPT-4's style and behavior

## Open Problems - Preserving Diversity and Interestingness

- How to restore and preserve interestingness and diversity – all the styles and worldviews present in the base models?

[ChatGPT and The Art of Post-Training. Barret Zoph & John Schulman. https://docs.google.com/presentation/d/11KWCKUORnPpVMSY6vXgBeFSWo7fJcuGQ9yuR6vC1pzE/edit#slide=id.p]

# Let's Try to Solve the Problem

| Elements of SFT | Properties | Comment |
|---|---|---|
| **Model** | Pre-trained with rich knowledge encoded | nothing to blame |
| **Data** | Limited Size and Coverage (10B-100B in SFT v.s. 1T-10T in pre-training) | not perfect but cannot blame |
| **Algorithm** | **Minimizing cross-entropy (CE) loss** | **is this good?** 🤔 |

# CE seems Effective for …

Input | Model | Prediction | Label



Convolution Neural Network → "Dog" "Cat"

**Cross-Entropy Loss**

Back-propagation

💡 CE is Effective for **Classification**

"I like to drink" → Langauge Model → "Tea" "Coffee"

**Cross-Entropy Loss**

Back-propagation

🤔 Is CE Effective for **Generation**?

24

# Understanding Generation Tasks

| Classification | Generation |
|---|---|

**Target**

$$\mathcal{X} \mapsto \mathcal{Y}$$

(**function**: many-to-one)

$$\mathcal{X} \mapsto \Delta(\mathcal{Y})$$

(**distribution**: one-to-many)

**Illustration**



Discriminant Model



Generative Model

Remark for LLMs:
- responses are **not unique**
  (variation in formats, styles, or reasoning paths)
- (SFT) data is hard to cover all cases

25

# Theory of CE

**CE Loss (Empirical)**

$$\min_{\theta} - \sum_{(x_i, y_i) \sim D} y_i^{\top} \log f_{\theta}(y_i \mid x_i)$$

$(x_i, y_i)$: input-label pair

$f_{\theta}(y \mid x)$: the conditional prediction distribution

$\theta$: parameters of neural network

**CE Loss (Population)**

$$\max_{\theta} \mathbb{E}_{x \sim \rho} \mathbb{E}_{y \sim p(\cdot \mid x)} \log f_{\theta}(y \mid x)$$

$\rho$: prompt distribution

$p(\cdot \mid x)$: the conditional data distribution to learn

Equivalence

CE can be used to learn a distribution

If the data samples are "abundant"

✅       ✅       ❌

Classification       Pre-training       SFT

(one label sample is enough)       (huge data)       (data is limited)

**Forward KL Divergence**

$$\min_{\theta} \mathbb{E}_{x \sim \rho} \text{KL}(p(\cdot \mid x), f_{\theta}(\cdot \mid x)) + \text{constant}$$

Distribution Matching       26

# Summary

**Challenge**:
We need to protect LLM's output diversity during SFT

**Understanding**:
CE easily fits to the empirical data and loses the diversity

**Goal**:
Designing new formulation and algorithm for SFT

# Analyzing Cross-Entropy Loss

Setting: $y \sim f_\theta(\,\cdot\,|\,x)$ and $f_\theta(i\,|\,x) = \dfrac{\exp(\theta_i)}{\sum_{j=1}^{K} \exp(\theta_j)}$

Gradient of CE: assuming $i$-th token is the label

$$-\nabla_\theta \mathcal{L}_{\mathrm{CE}}(\theta) = [-f_\theta(1|x), -f_\theta(2|x), \ldots, 1 - f_\theta(i|x), \ldots, -f_\theta(K|x)].$$

Implication:

Target token (label)'s logit $\uparrow$ while other tokens' logits $\downarrow$

# Distribution Matching as Flow Transfer

**Proposition 1.** *The gradient of CE specifies a logit flow map: each source token $j$ transfers $f_\theta(j|x)$ logits to the target token $i$. Formally,*

$$
\begin{aligned}
-\nabla_\theta \mathcal{L}_{\mathrm{CE}}(\theta) &= \sum_{j=1, j \neq i}^{K} w_{i \leftarrow j} \cdot e_{i \leftarrow j} \qquad (2) \\
w_{i \leftarrow j} &= f_\theta(j|x) \\
e_{i \leftarrow j} &= [0 \; \cdots \; \underbrace{1}_{\text{i-th position}} \; \cdots \; \underbrace{-1}_{\text{j-th position}} \; \cdots \; 0]
\end{aligned}
$$

Example: $\qquad\qquad f_\theta = [0.1, 0.3, 0.6]$ $\qquad\qquad$ Label: #2

Gradient: $\qquad\qquad g = [-0.1, 0.7, -0.6]$

Flow perspective: $\quad g = 0.1 * [-1 \quad 1 \quad 0] + 0.6 * [0 \quad 1 \quad -1]$

Logits flow from **source** tokens = Logits flow to **target** token

# Limitations of CE

Procedure of CE

#1 While there exists source token $j \neq i$ with $f_{\theta_k}(j|x) > 0$, continue the following steps.

#2
- Find any $j$ with $f_{\theta_k}(j|x) > 0$
- Decrease the logit for source token $j$ by learning rate $\eta$ and weight $w_{i \leftarrow j}$:
$$\theta_{k+1}[j] = \theta_k[j] - \eta * w_{i \leftarrow j}$$
- Increase the logit for the target token $i$ in a similar manner:
$$\theta_{k+1}[i] = \theta_k[i] + \eta * w_{i \leftarrow j}$$

Limitation 1: Unbounded Transfer

Limitation 2: All-to-one Update

CE

unbounded transfer

1 2 3 4

all-to-one update

1 2 3 4

distribution collapse

# Proposed Solutions

**Procedure of Our Method**

**#1** While the target token $i \notin \mathrm{argmax}\, f_{\theta_k}(\cdot|x)$, continue the following steps.

**#2**
- Calculate the model's best prediction $j = \mathrm{argmax}\, f(\cdot|x)$
- Decrease the logit for source token $j$ by learning rate $\eta$ and weight $w_{i \leftarrow j}$:
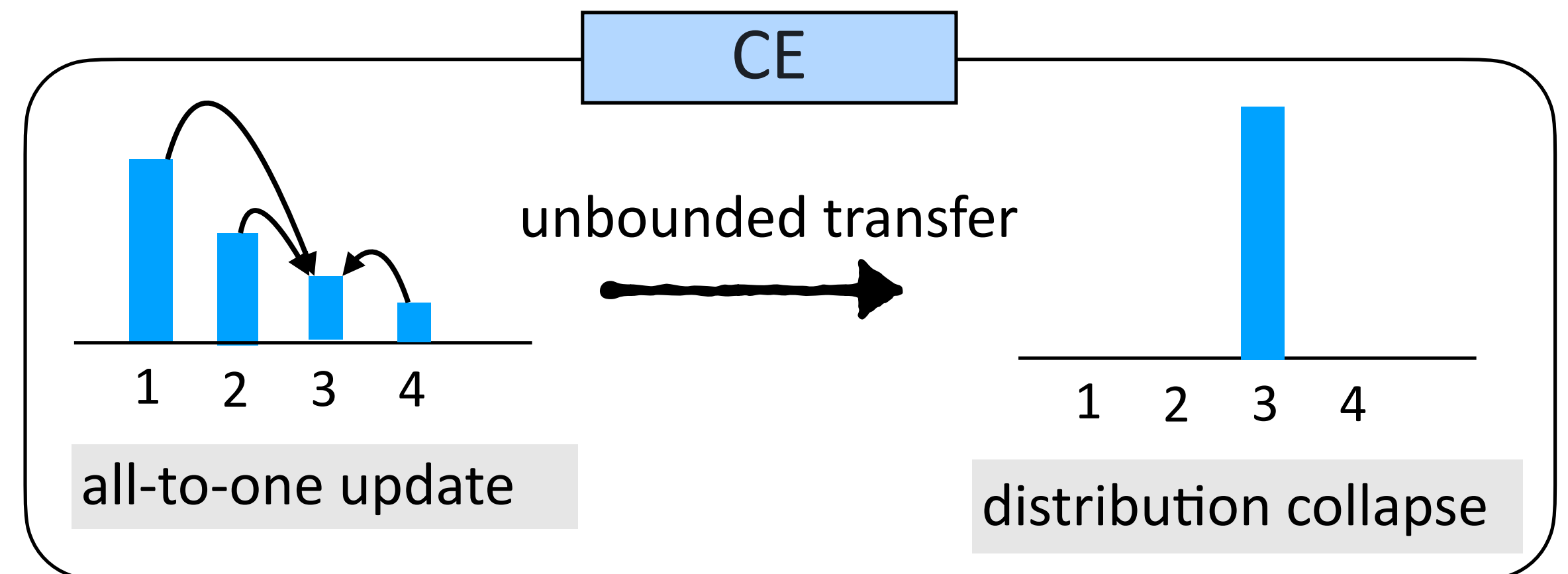$$\theta_{k+1}[j] = \theta_k[j] - \eta * w_{i \leftarrow j}$$
- Increase the logit for the target token $i$ in a similar manner:
$$\theta_{k+1}[i] = \theta_k[i] + \eta * w_{i \leftarrow j}$$

CE

**Technique 1: Adaptive Termination**

ll-to-one update

**Technique 2: Sparse Update** distribution collapse

GEM

sparse update

adaptive termination

diversity-keeping

greedy decoding can handle

# Our Insight: Dimension Increase

**Procedure of Our Method**

While the target token $i \notin \mathrm{argmax}\, f_{\theta_k}(\cdot|x)$, continue the following steps.

- Calculate the model's best prediction $j = \mathrm{argmax}\, f(\cdot|x)$

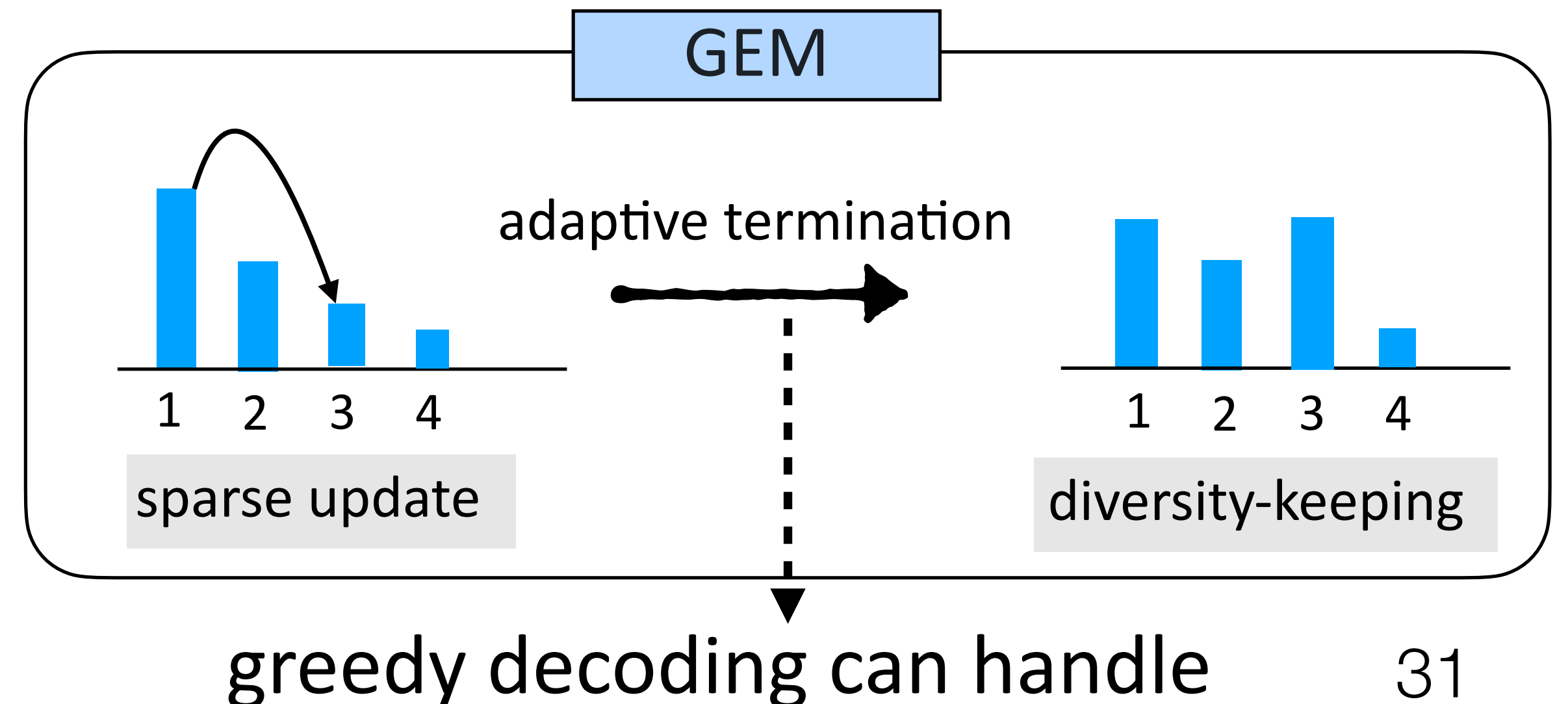- Decrease the logit for source token $j$ by learning rate $\eta$ and weight $w_{i \leftarrow j}$:

$$\theta_{k+1}[j] = \theta_k[j] - \eta * w_{i \leftarrow j}$$

- Increase the logit for the target token $i$ in a similar manner:

$$\theta_{k+1}[i] = \theta_k[i] + \eta * w_{i \leftarrow j}$$

🤔    What is the magic? Can we generalize this to neural network training?

💡    Introduce an **auxiliary variable** (dimension increase) that implements the scheme of sparse update and adaptive termination

# Towards a Game Formulation

High-level design: introduce an **another player** $q$ to the distribution matching

$$\min_f \quad \mathcal{L}(f, q) \triangleq \mathbb{E}_x \mathbb{E}_{y^{\texttt{real}} \sim p(\cdot|x)} \mathbb{E}_{y^{\texttt{gene}} \sim q(\cdot|x)} \left[ \log f(y^{\texttt{gene}}|x) - \log f(y^{\texttt{real}}|x) \right]$$

$$\max_q \quad \mathcal{Q}(f, q) \triangleq \mathbb{E}_x \mathbb{E}_{y^{\texttt{gene}} \sim q(\cdot|x)} \left[ \log f(y^{\texttt{gene}}|x) \right] + \beta \cdot \mathcal{H}(q(\cdot|x)).$$

Intuitive Understanding:
- $f$: increase the likelihood on real data and decrease likelihood on the generated data

- $q$: increase the energy induced by $\log f$ with entropy regularization

33

# Understanding the Game

**main player**

$$-\nabla_\theta \mathcal{L}(f_\theta, q) = \sum_{j=1, j\neq i}^{K} w_{i\leftarrow j} \cdot e_{i\leftarrow j},$$

$$w_{i\leftarrow j} = q(j|x).$$

→ **flow transfer**

→ **controller**

**meta-controller**

$$\operatorname*{argmax}_{q} \mathcal{Q}(q, f) = \begin{cases} \delta_j(x) \text{ with } j = \operatorname{argmax} f_i(\cdot|x) & \text{if } \beta = 0 \\ \texttt{softmax}(1/\beta * \log f(y|x)) & \text{if } \beta > 0 \end{cases}$$

$f$

$q$

1 2 3 4

$\beta = 0$    1 2 3 4

$\beta = 0.5$    1 2 3 4

$\beta = 1$    1 2 3 4

$\beta = 5$    1 2 3 4

$\beta \to 0$: sparse update

$\beta \to 1$: same as CE

$\beta \to \infty$: uniform update

34

# Connection with Probability Transfer

**Proposition 2.** *For a data distribution satisfying $p(y|x) > 0$, with $\beta > 0$, the game in Equations* (3) *and* (4) *posses a unique Nash equilibrium point:*

$$\begin{cases} f^\star = softmax(\beta * \log p) \\ q^\star = p \end{cases} \tag{7}$$

*Furthermore, $f^\star$ corresponds to the optimal solution to the distribution matching problem (with $1/\beta = (\gamma + 1)$), which minimizes the <u>reverse</u> KL divergence with entropy regularization:*

$$f^\star = \underset{f}{\arg\min} \, \mathbb{E}_x \left[ D_{\mathrm{KL}}(f(\cdot|x), p(\cdot|x)) - \gamma \mathcal{H}(f(\cdot|x)) \right]. \tag{8}$$

| **Terminology** | **Reserve** KL Minimization | **Entropy** Maximization |
|---|---|---|
| **Role** | Fit the data distribution | Protect the output diversity |

For $\beta = 0$, there are **multiple** Nash equilibrium points with non-closed-form solutions → future work

# Training Algorithm

Idea: block-wise gradient-descent and coordinate descent

$$\begin{cases} f_{\theta_{k+1}} = f_{\theta_k} - \nabla_\theta \mathcal{L}(f_\theta, q_k) \mid_{\theta=\theta_k} \\ q_{k+1} = \mathrm{argmax}_q \, \mathcal{Q}(f_{\theta_{k+1}}, q) = \mathtt{softmax}(1/\beta * \log f_{\theta_{k+1}}) \end{cases}$$

Feature 1: **Single**-model optimization

↳ There is no need of storing and explicit training of $q$

Optimization with the token space (**discrete**)

Feature 2: **Variance-reduced** gradient estimation

$$\mathcal{L}_{\mathrm{GEM}}(\theta) = \sum_i \sum_{y^{\mathrm{gene}}} {\color{red}q_k(y^{\mathrm{gene}}|x_i)} \cdot \left[ \log f_\theta(y^{\mathrm{gene}}|x_i) - \log f_\theta(y_i^{\mathrm{real}}|x_i) \right]$$

↳ We use the exact distribution (in GANs, stochastic approximation is used)

# Discussion: Difference with GANs

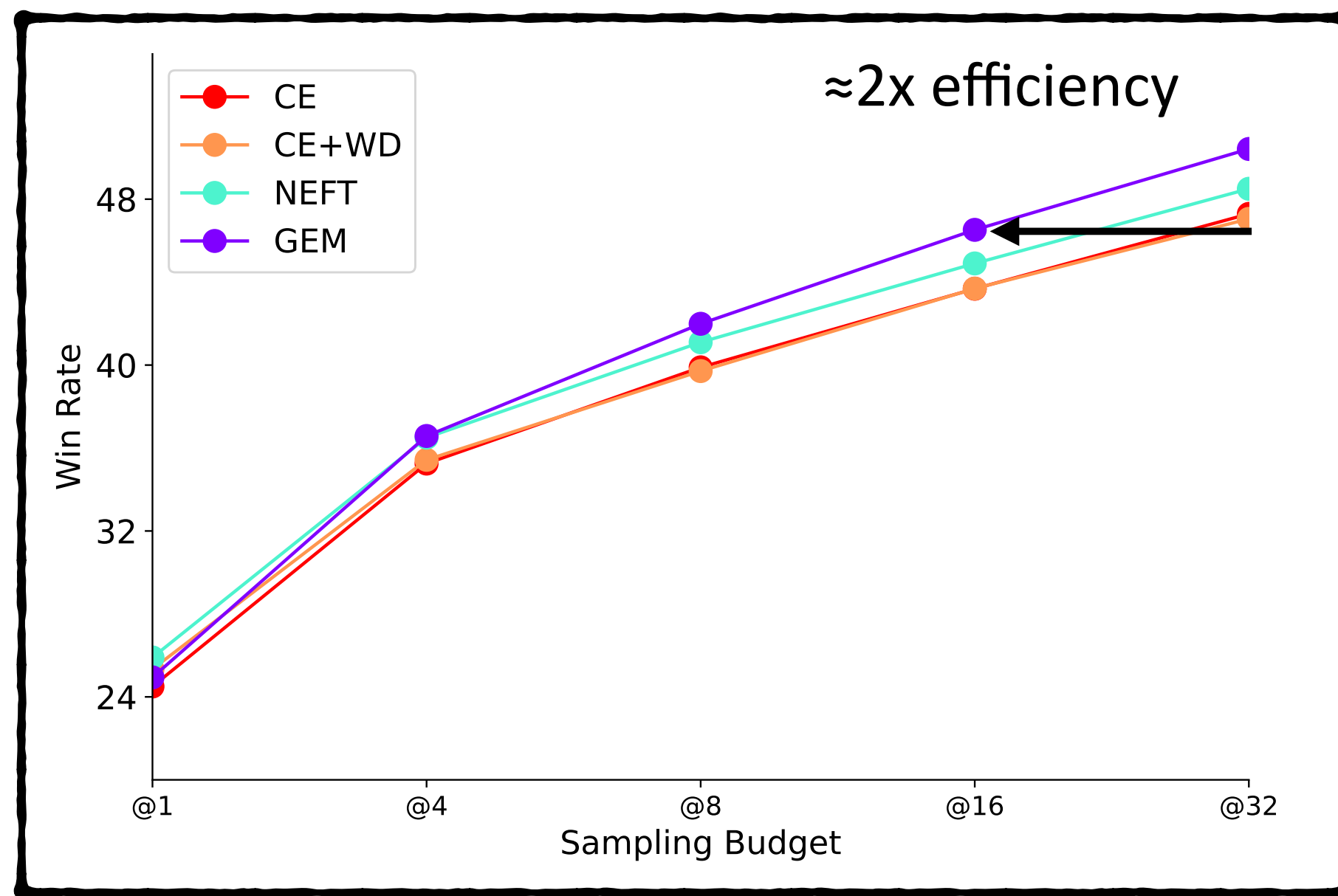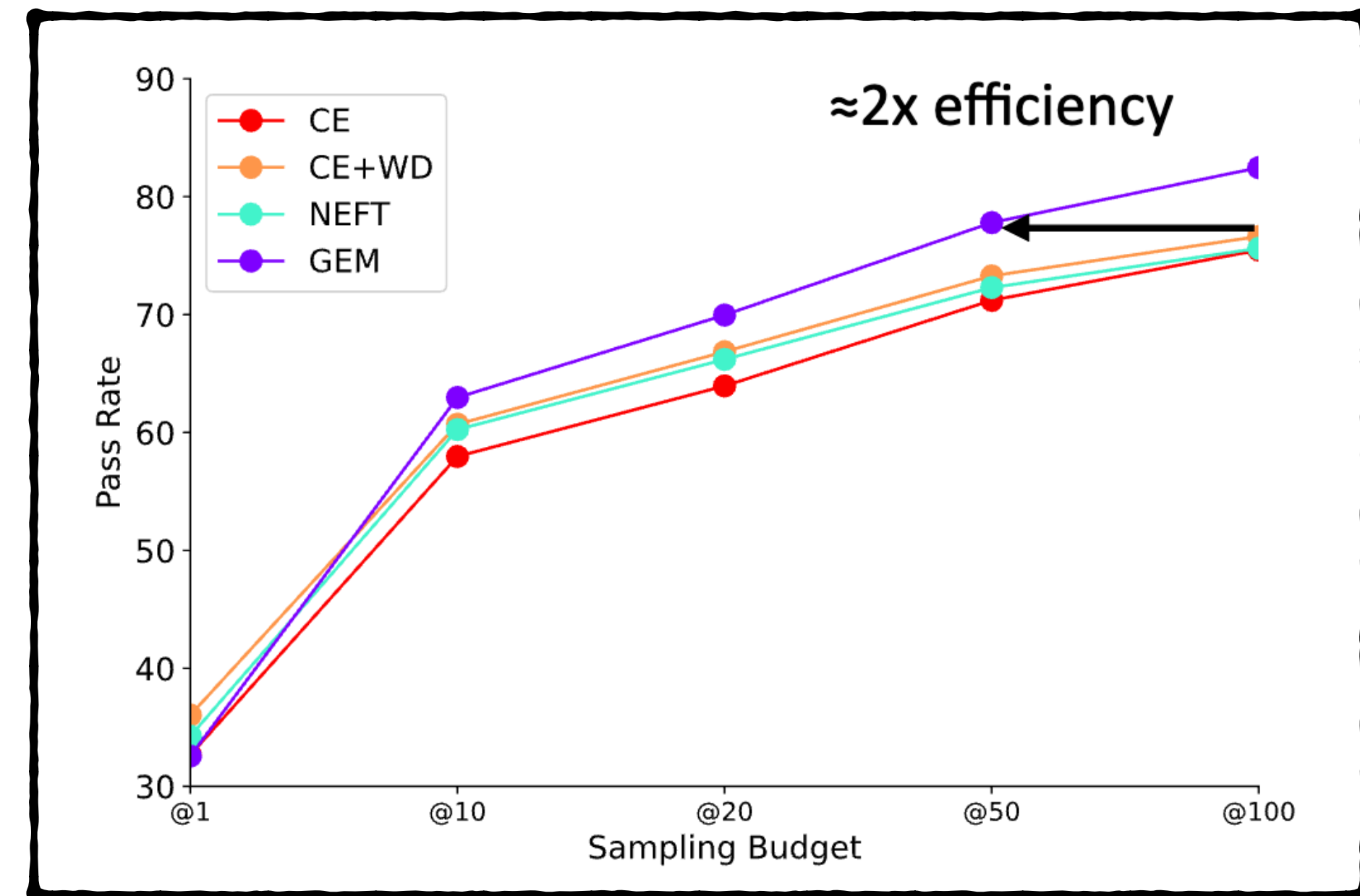| | **GAN** (generative adversarial network) | **GEM** (game-theoretic entropy maximization) |
|---|---|---|
| **Task** | Image Generation | Text Generation |
| **Challange** | Estimation the distance among two images is hard | Overfitting the data and losing output diversity |
| **Idea** | Introduction of discriminator | Introduction of flow-controller |
| **Computation Complexity** | High | Low |

# Test-Time Scaling

▸ Evaluation Method: Best-of-N Sampling
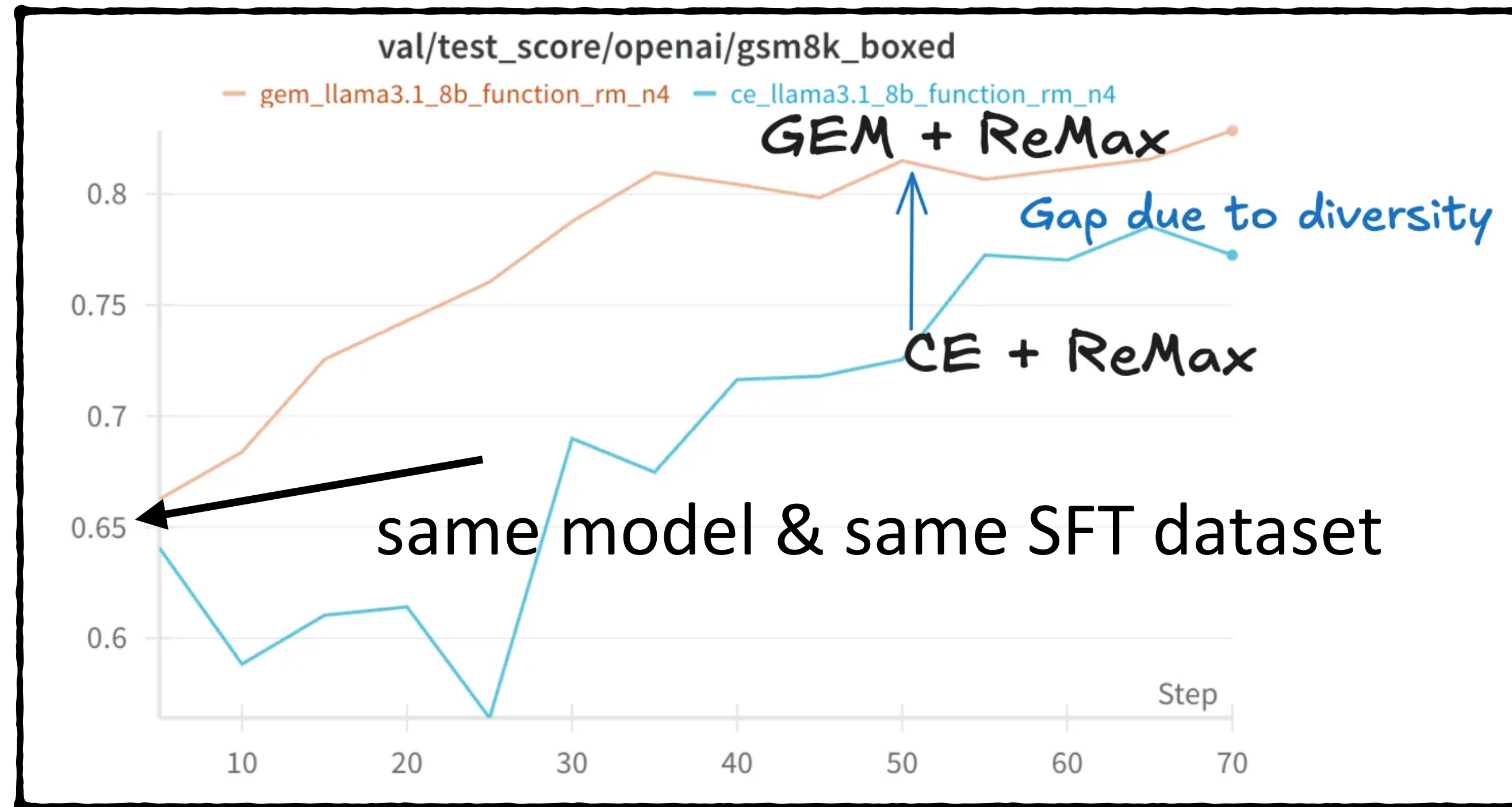
▸ Model: Llama-3.1-8B;  Dataset: Ultrafeedback



RLHF Alignment (Chat)

Code Generation

GEM requires about **2x** less sampling budget for comparable performance
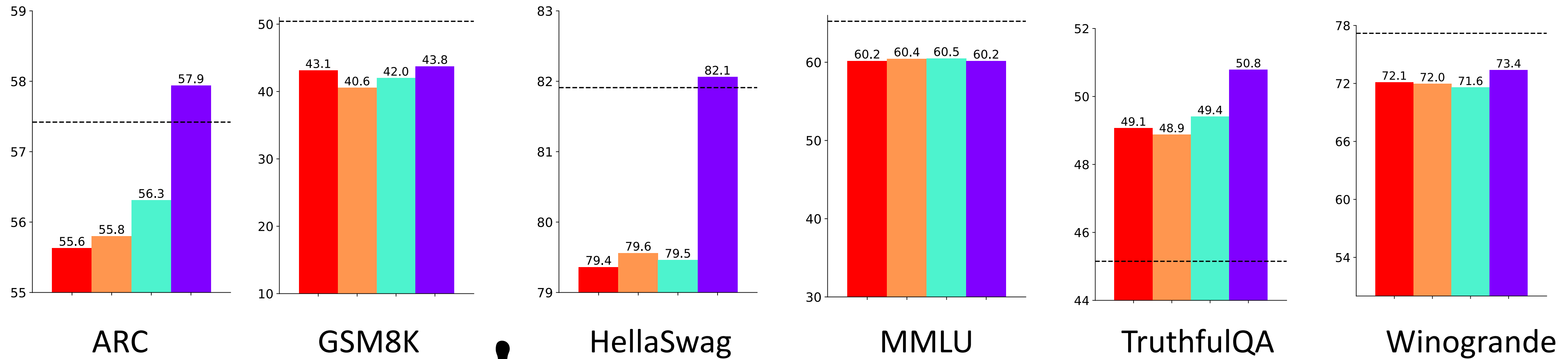
# Math Reasoning



val/test_score/openai/gsm8k_boxed

— gem_llama3.1_8b_function_rm_n4    — ce_llama3.1_8b_function_rm_n4

GEM + ReMax

Gap due to diversity

CE + ReMax

same model & same SFT dataset

[https://tangible-polo-203.notion.site/]

- ▸ Task: optimize CoT (reasoning steps) to answer math questions
- ▸ Reward: accuracy of final reward
- ▸ Model: Qwen-2.5-3B
- ▸ RL Algo: ReMax

[Li, Ziniu, et al. "Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models." ICML 2024.]
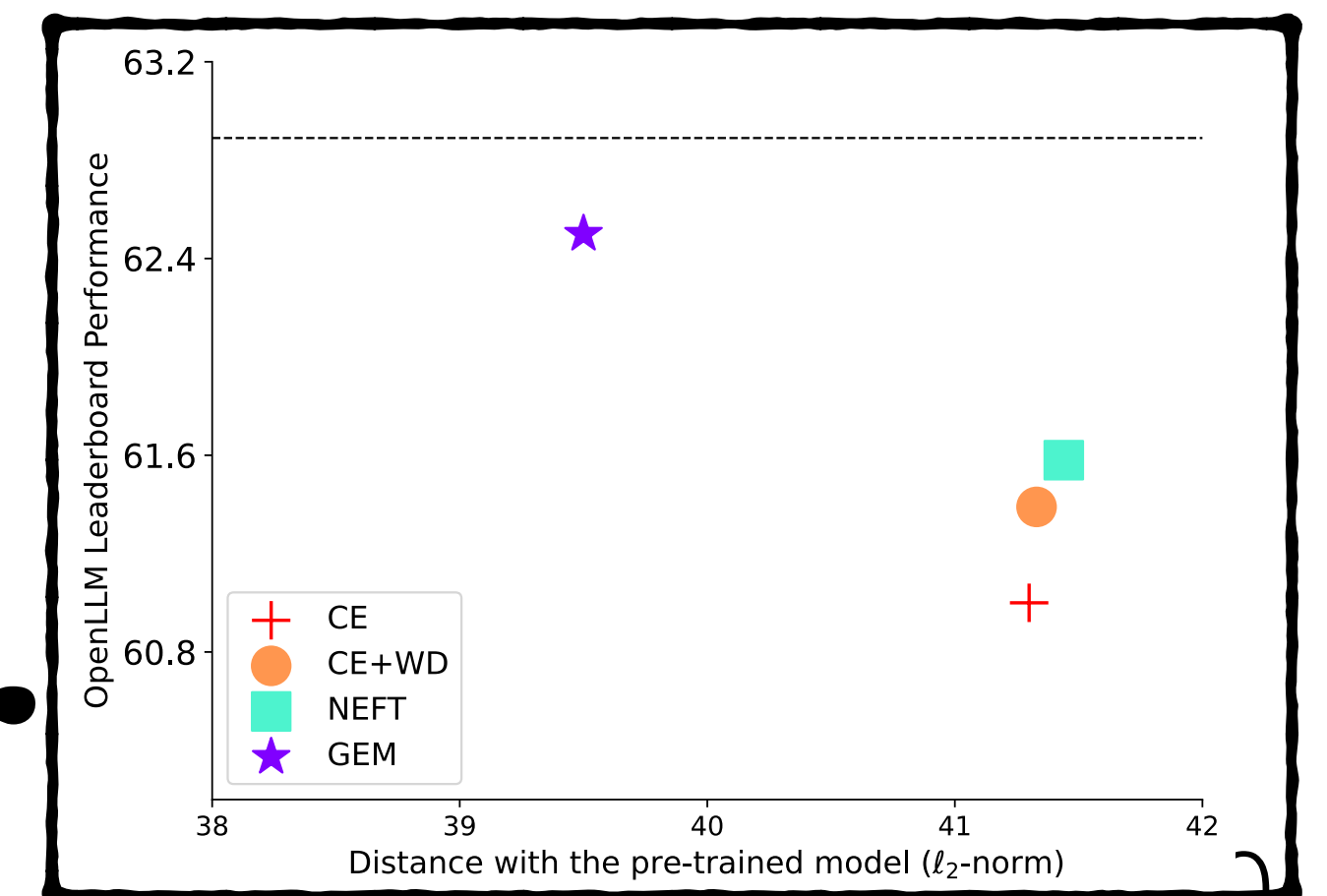
GEM improves the performance limit of RL training

# Alignment Tax



Legend: - - - Pre-trained | CE | CE+WD | NEFT | GEM

**ARC**
- CE: 55.6
- CE+WD: 55.8
- NEFT: 56.3
- GEM: 57.9

**GSM8K**
- CE: 43.1
- CE+WD: 40.6
- NEFT: 42.0
- GEM: 43.8

**HellaSwag**
- CE: 79.4
- CE+WD: 79.6
- NEFT: 79.5
- GEM: 82.1

**MMLU**
- CE: 60.2
- CE+WD: 60.4
- NEFT: 60.5
- GEM: 60.2

**TruthfulQA**
- CE: 49.1
- CE+WD: 48.9
- NEFT: 49.4
- GEM: 50.8

**Winogrande**
- CE: 72.1
- CE+WD: 72.0
- NEFT: 71.6
- GEM: 73.4

GEM fine-tunes the model with 83% less alignment tax

GEM-tuned model shows less overfitting to the data

Scatter plot: OpenLLM Leaderboard Performance vs Distance with the pre-trained model ($\ell_2$-norm)
- CE: +
- CE+WD: ●
- NEFT: ■
- GEM: ★

40

# PRESERVING DIVERSITY IN SUPERVISED FINE-TUNING OF LARGE LANGUAGE MODELS

Ziniu Li[1,2], Congliang Chen[1,2], Tian Xu[3], Zeyu Qin[4], Jiancong Xiao[5], Zhi-Quan Luo[1,2], and Ruoyu Sun[1,2,†]
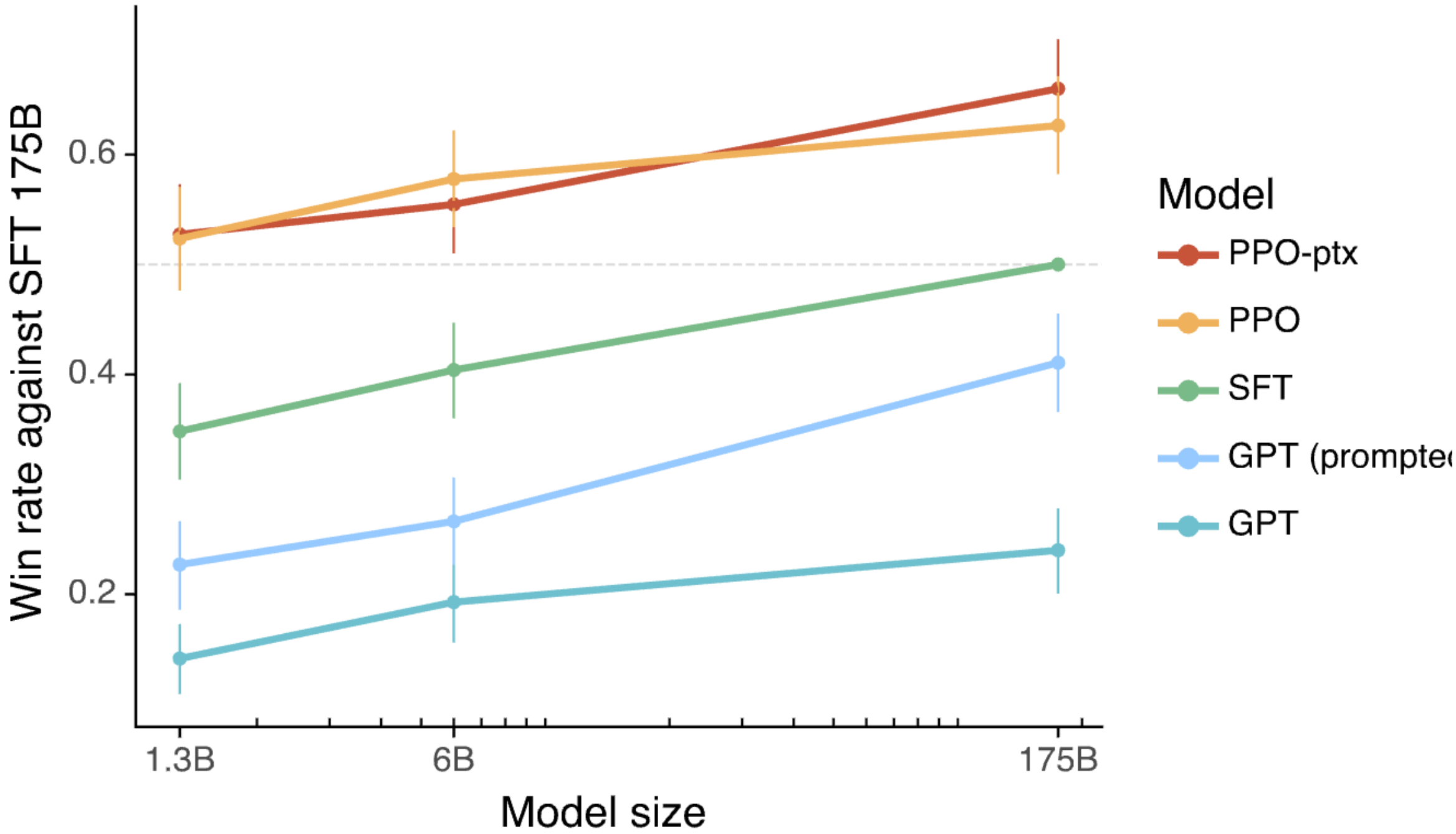
ICLR 2025
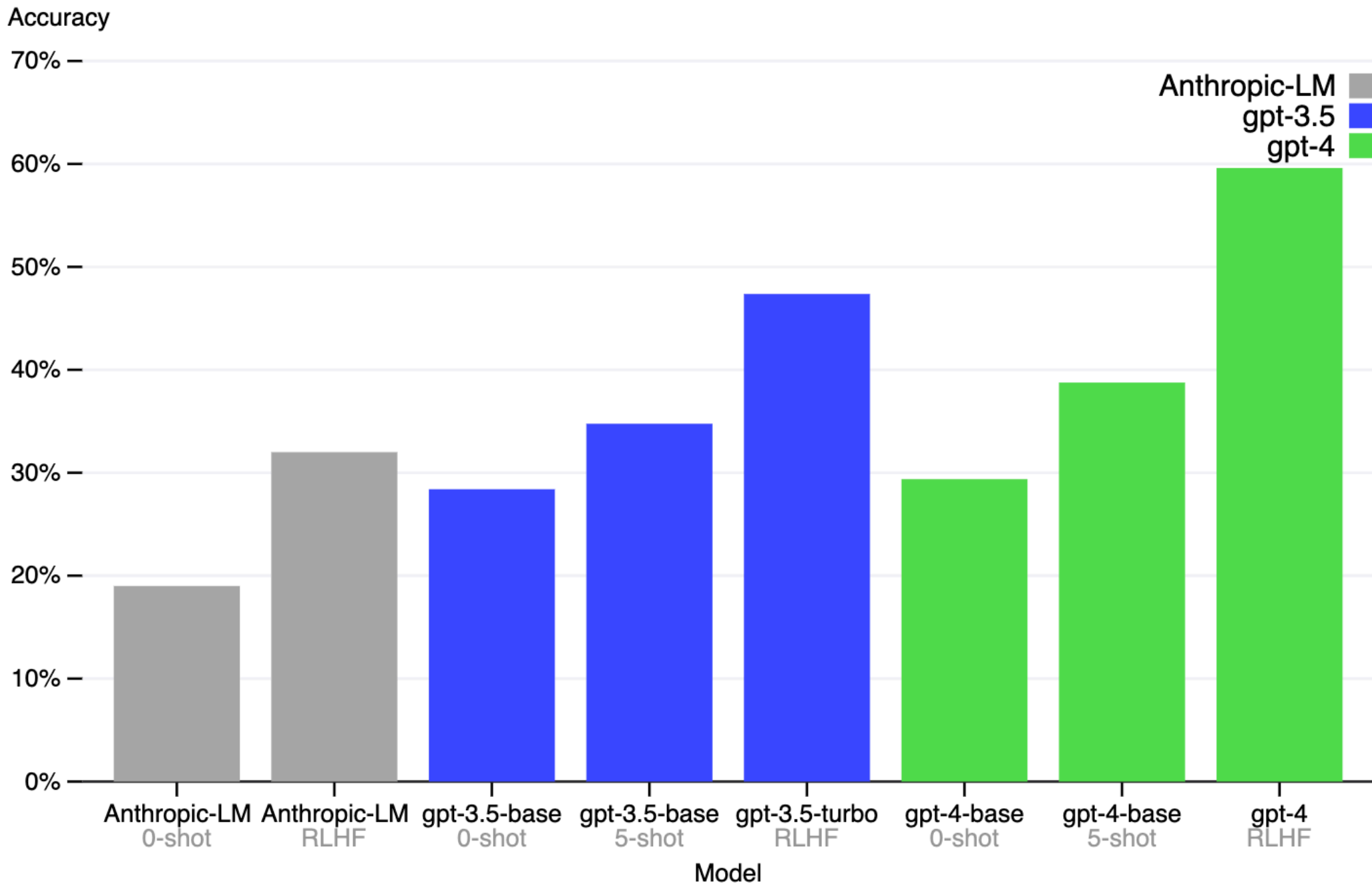
NeurIPS 2024 FITML Workshop Best Paper Runner-up



Paper

Code

# Part IV: Efficient and Scalable Reinforcement Learning in LLMs

# RL Task: Alignment



Only PPO Achieves a Win Rate Above 50%
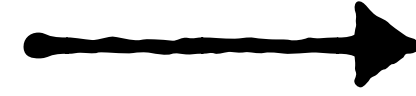
[Ouyang, Long, et al. "Training language models to follow instructions with human feedback." *NeurIPS 2022.*]



RLHF Enhances Acc. by More Than 10%
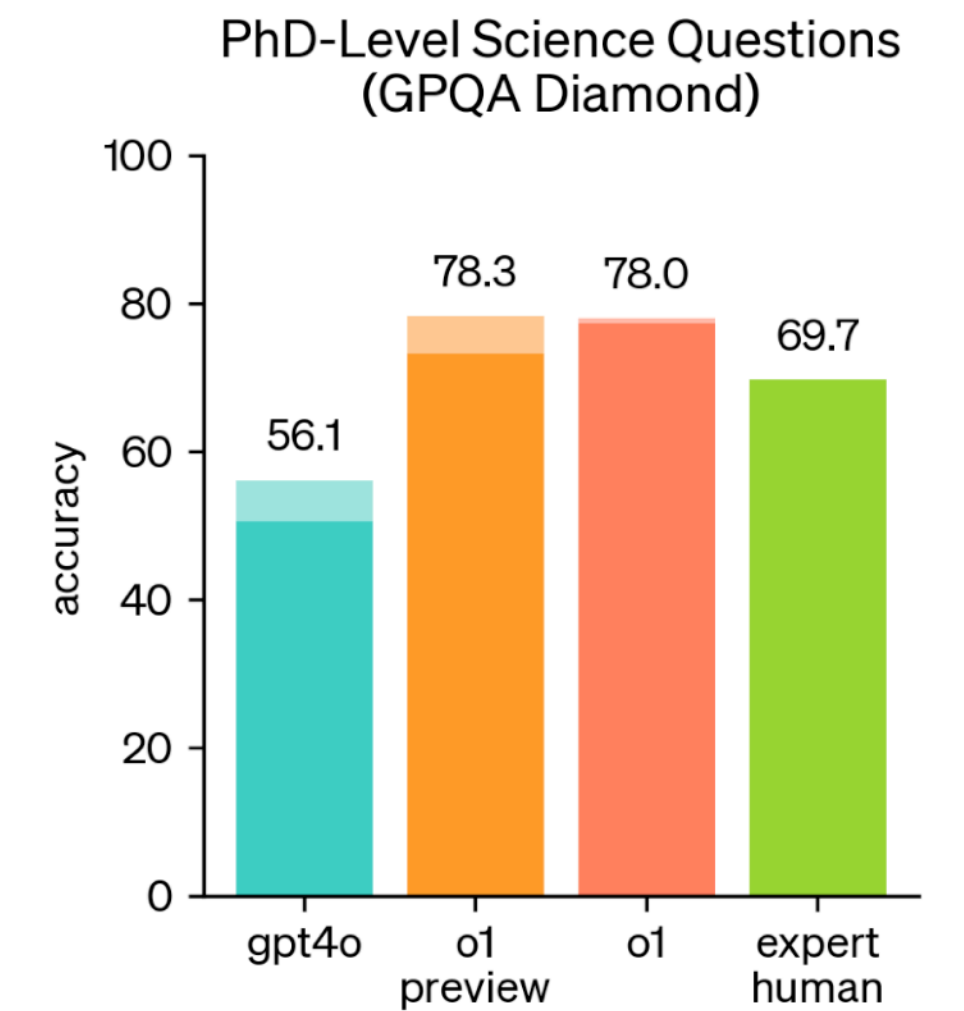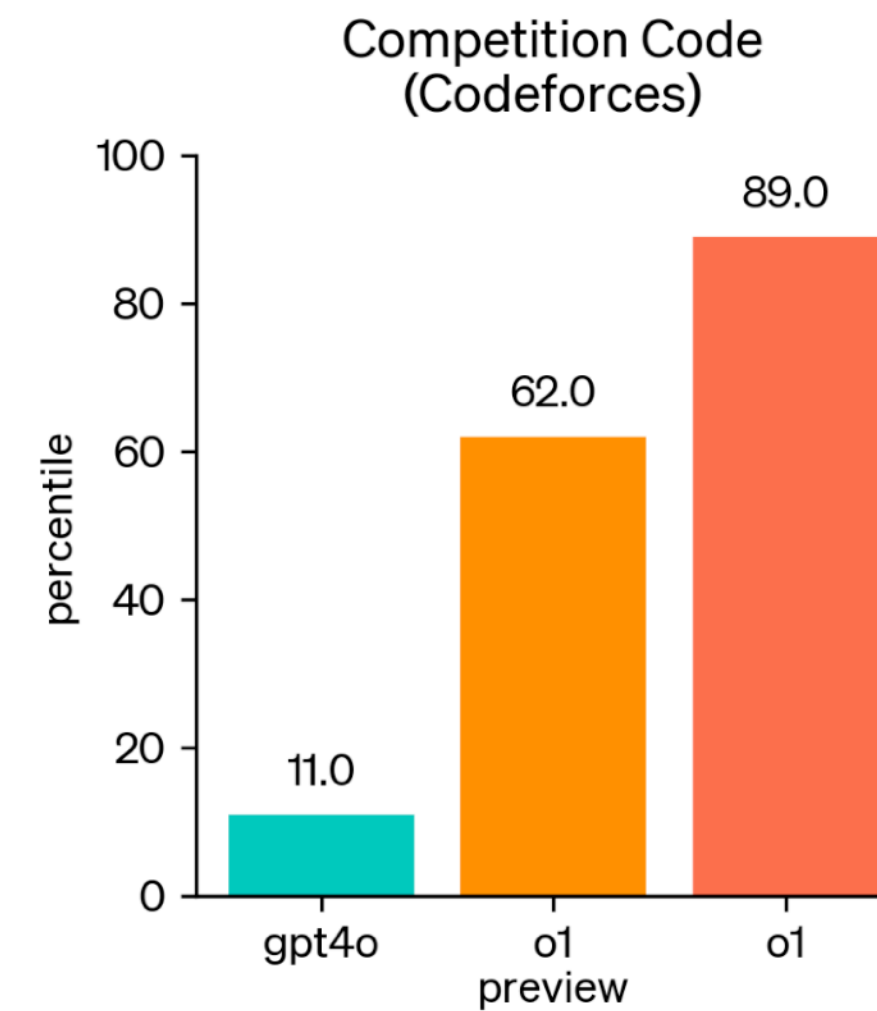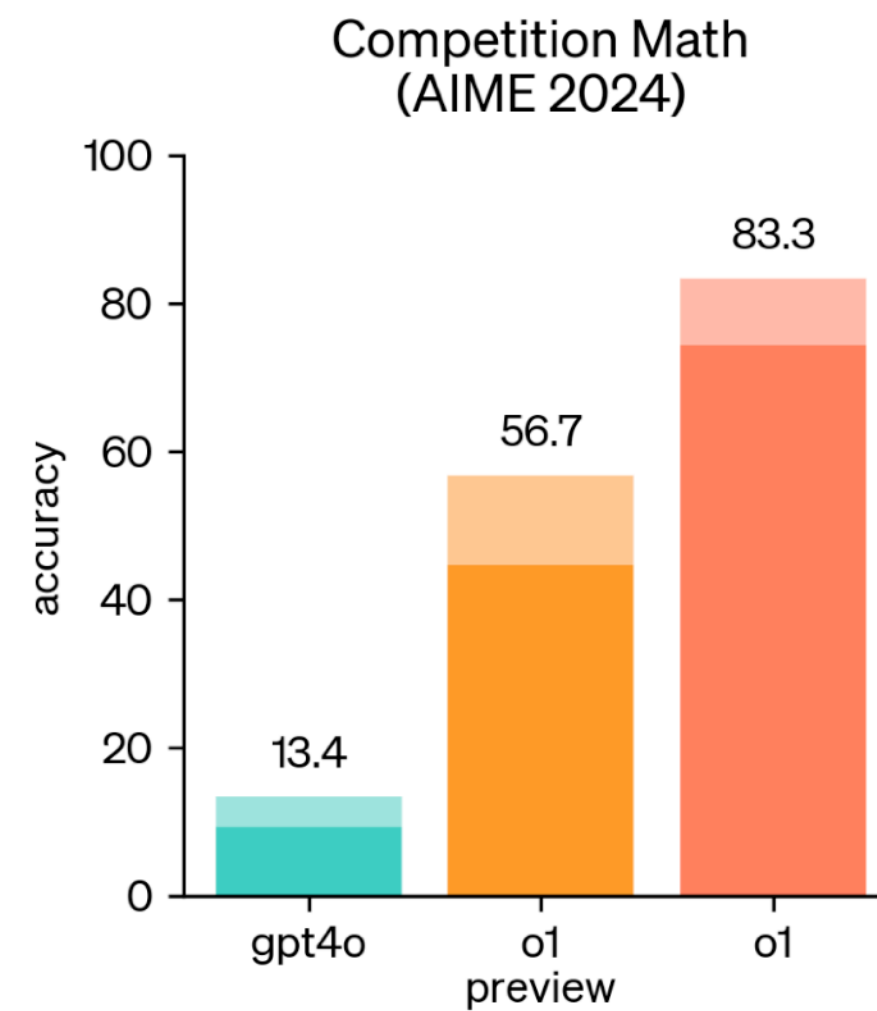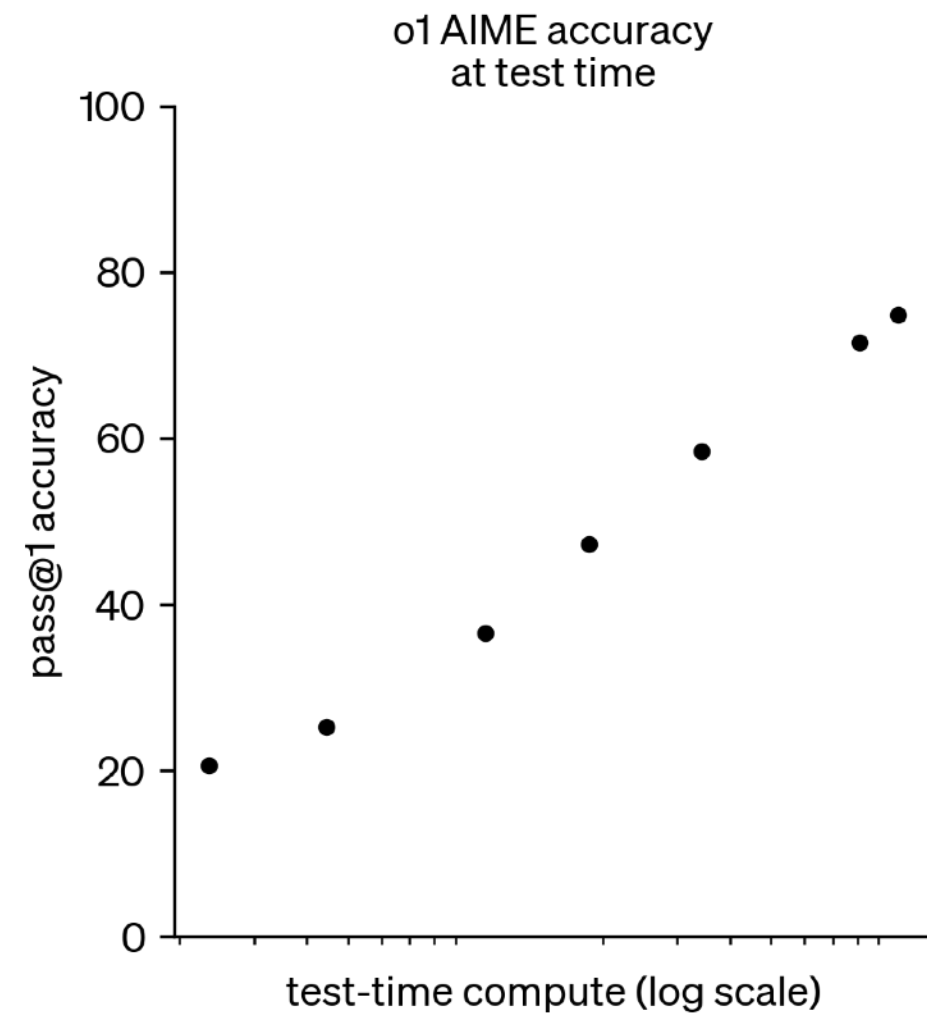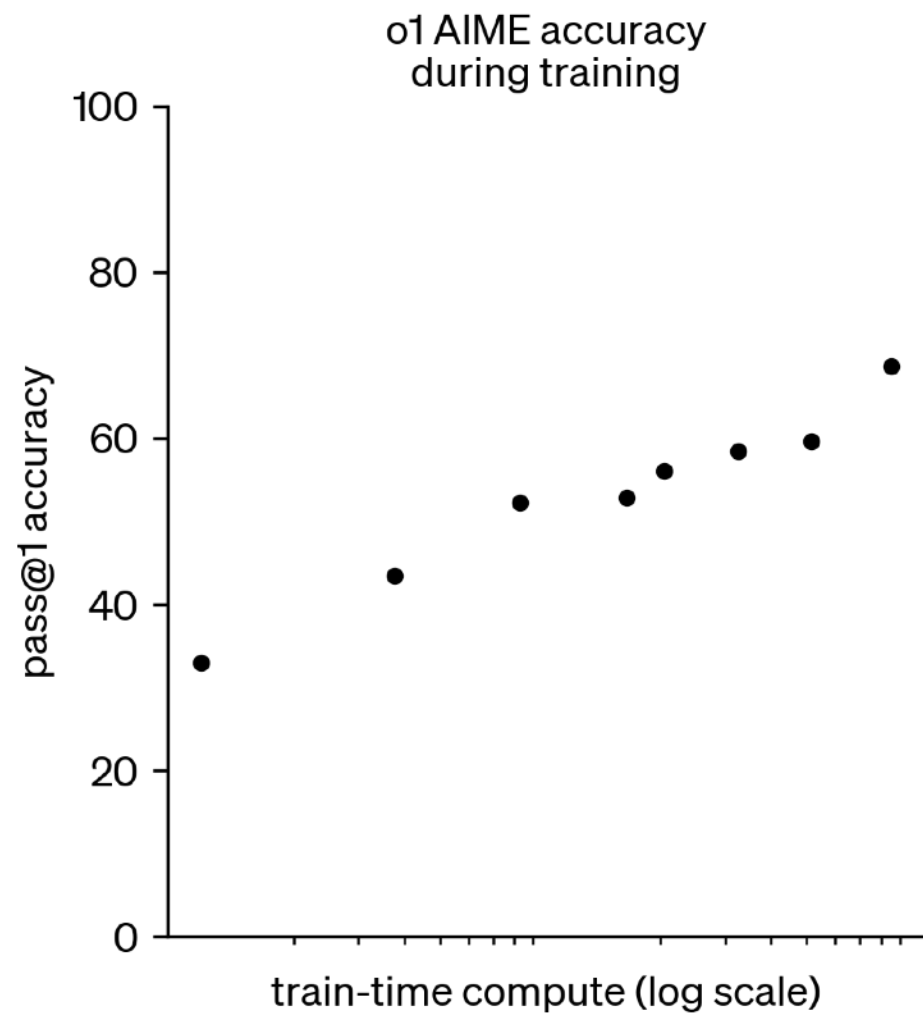
[Achiam, Josh, et al. "Gpt-4 technical report." *arXiv preprint arXiv:2303.08774* (2023).]

# RL Task: Eliciting Reasoning



Test-time Scaling → Huge Improvement in Challenging Tasks

[https://openai.com/index/learning-to-reason-with-llms/]

RL training enables models to think deep

o1 can exceeds GPT-4o by 40+ points on MATH, code, and PhD-Level QA

# How does RL work in LLMs?

**Alignment**
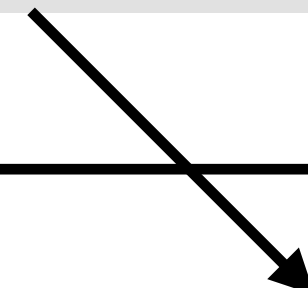
Collect preference data

↓

Train a reward model

↓

Optimize LLM against reward model
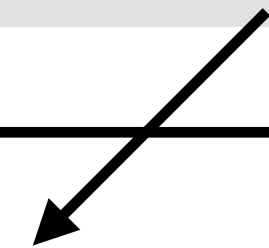
**Reasoning**

Collect data with answers

↓

Design a rule-based reward function
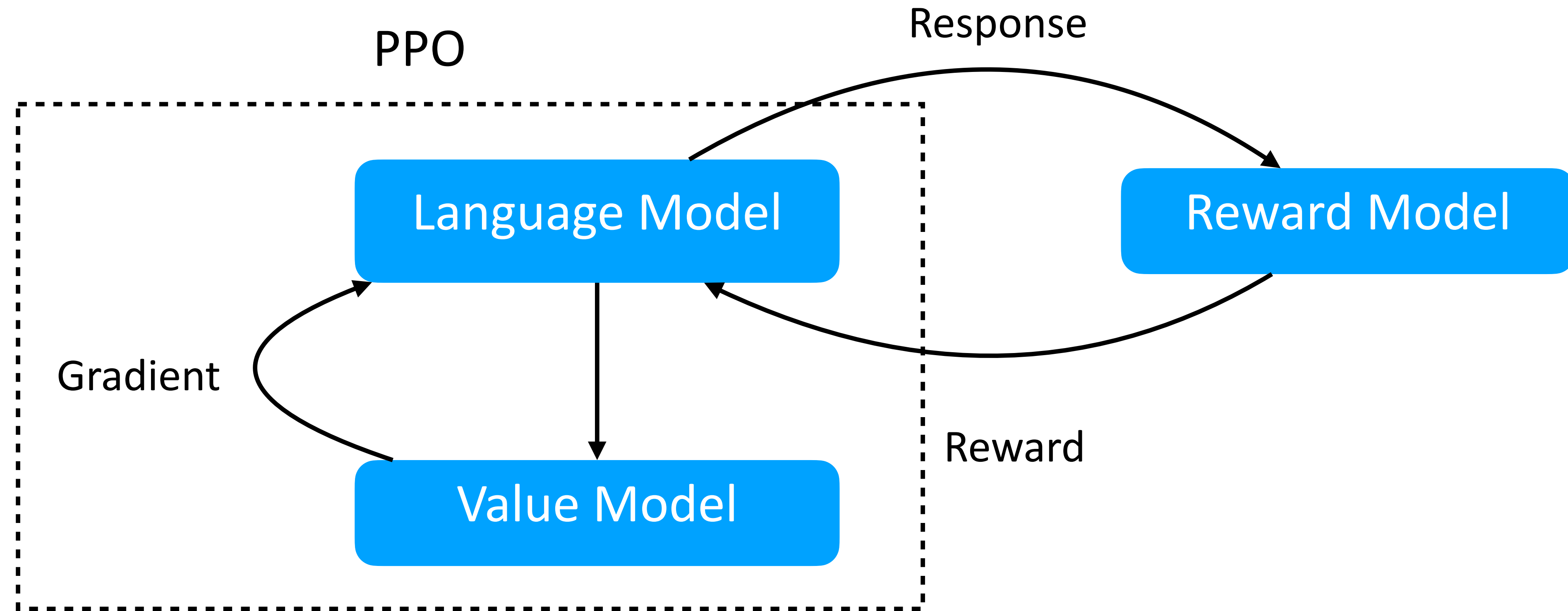
↓

Optimize LLM against reward model

**PPO** is the default RL algorithm
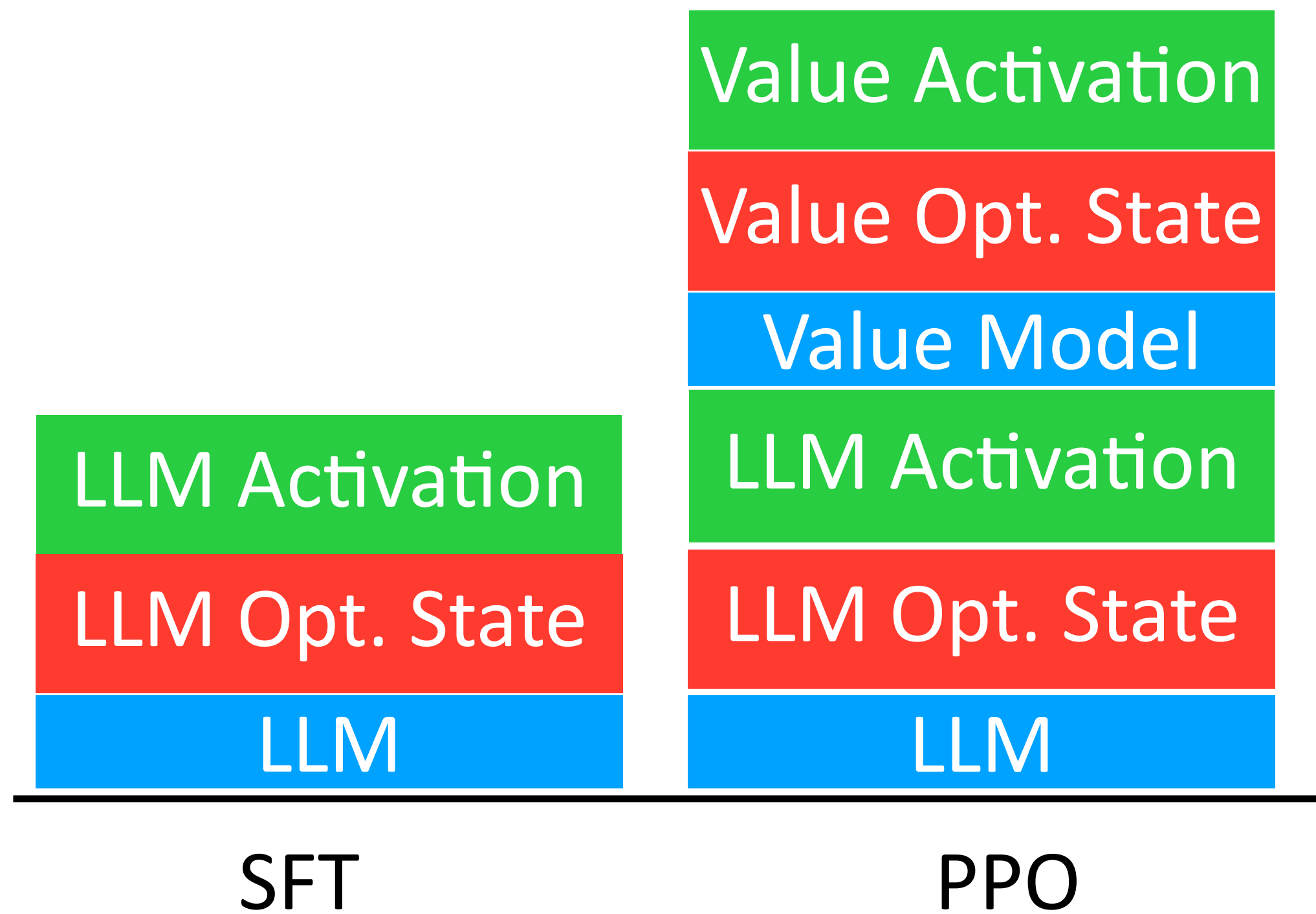
# Introduction to PPO

Objective: $$\max_{\theta} \mathbb{E}_{y_{1:T} \sim \pi_{\theta}(\cdot|x)}[r(x, y_{1:T})]$$



[Schulman, John, et al. "Proximal policy optimization algorithms." *arXiv preprint arXiv:1707.06347* (2017).]

# PPO is Computationally Inefficient



SFT    RM    PPO

Table 4: E2E time breakdown for training a 13 billion parameter ChatGPT model via DeepSpeed-Chat on a single DGX node with 8 NVIDIA A100-40G GPU.

| Model Sizes | Step 1 | Step 2 | Step 3 | Total |
|---|---|---|---|---|
| Actor: OPT-13B, Reward: OPT-350M | 2.5hr | 0.25hr | 10.8hr | 13.6hr |

[Yao, Zhewei, et al. "DeepSpeed-Chat: Easy, Fast and Affordable RLHF Training of ChatGPT-like Models at All Scales." *arXiv:2308.01320* (2023)]

PPO's training takes more memory

PPO's training is slow

**Value model** is the bottleneck of PPO

# Can We Improve PPO?

🤔     Can we achieve RL training without the value model?

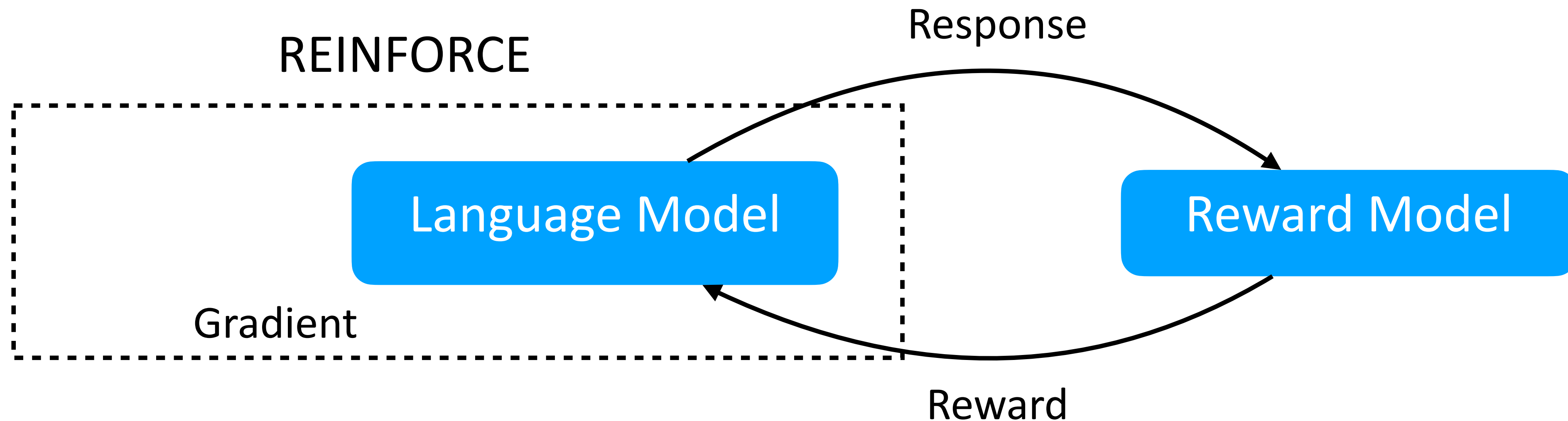🚀     If Yes, we can save memory and accelerate training

💡     **REINFORCE** is an RL algorithm without value model

[Williams, Ronald J. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." *Machine learning* 8 (1992): 229-256.]

# Introduction to REINFORCE

Objective: $$\max_\theta \mathbb{E}_{y_{1:T} \sim \pi_\theta(\cdot|x)}[r(x, y_{1:T})]$$

REINFORCE

Response



Language Model

Reward Model

Gradient

Reward

[Williams, R. J. Reinforcement-learning connectionist systems. College of Computer Science, Northeastern University, 1987.]

REINFORCE: $$\text{gradient} = \mathbb{E}_{y_{1:T} \sim \pi_\theta(\cdot|x)}[r(x, y_{1:T}) \cdot \nabla_\theta \log \pi_\theta(y_{1:T} | x)]$$

No Value Model

Stochastic Gradient Estimation in Practice

# However, REINFORCE does not Work



REINFORCE's gradient has a high variance

REINFORCE's reward does not increase

# Why is Variance so High?

🤔 REINFORCE is often criticized for a high gradient variance. But why?

[Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.]

$$\text{gradient} = \mathbb{E}_{a_{1:T} \sim \pi_\theta(\cdot|x)}[r(x, a_{1:T}) \cdot \nabla_\theta \log \pi_\theta(a_{1:T}|x)]$$

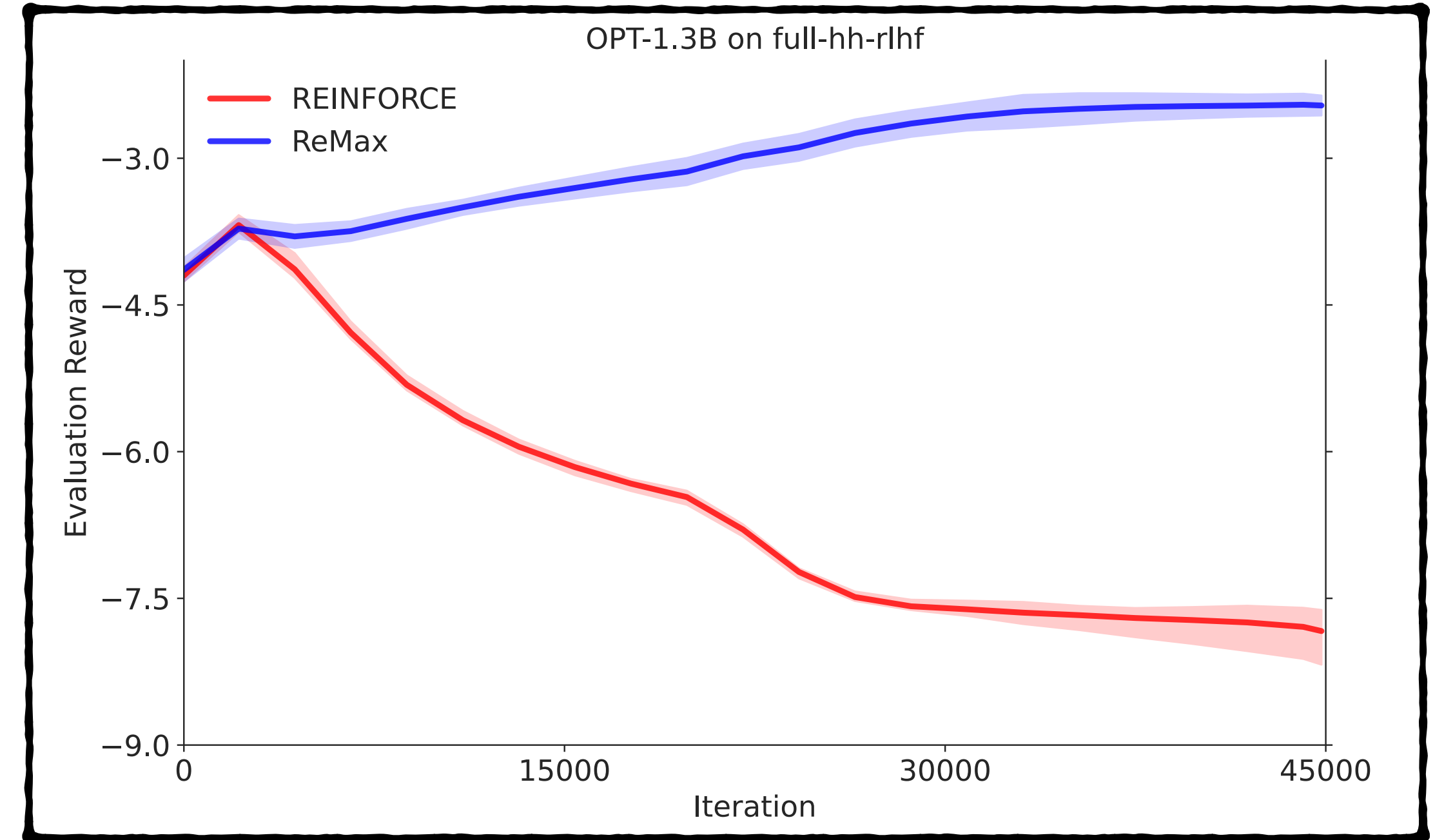**<u>Sample space is large</u>**

Size: (vocabulary size)$^{\text{sequence length}}$

Llama-3: $(128k)^{8k}$

**<u>Rewards vary across samples</u>**

Reward range of open-ended question-answers: [-14, 7]

# Introduction to ReMax

Key Idea: Introduce a **baseline value** for accurate gradient estimation

$$\nabla_\theta \mathbb{E}_{x,y}[r(x, y_{1:T})] = \mathbb{E}\left[\nabla_\theta \log \pi_\theta(y_{1:T}|x) \cdot \boxed{[r(x, y_{1:T}) - b(x)]}\right]$$

Advantage

$$b(x) = r(x, y'_{1:T}), \quad y'_t = \arg\max_{y_t} \pi_\theta(y_t|x_t, y_{1:t})$$

Greedy Decoding

Remark: 1) Subtracting a RV by a constant does not change the variance

2) ReMax introduces a RV $b \cdot \nabla_\theta \log \pi_\theta(y_{1:T}|x) \to$ **control variate**

# Why Greedy Decoding?

$$\nabla_\theta \mathbb{E}_{x,y}[r(x, y_{1:T})] = \mathbb{E}\left[\nabla_\theta \log \pi_\theta(y_{1:T}|x) \cdot \left[r(x, y_{1:T}) - b(x)\right]\right]$$

$$b(x) = r(x, y'_{1:T}), \quad y'_t = \arg\max_{y_t} \pi_\theta(y_t | x_t, y_{1:t})$$

Reason 1: greedy decoding corresponds to **mode** of the distribution →
**effective estimation**

Reason 2: value of greedy decoding ensures **independence** between the
baseline and original RVs → **stable estimation**

Reason 3: if there is a response better than the greedy one, improve it's
likelihood

# ReMax Algorithm

**Algorithm 2** ReMax for Aligning Large Language Models

**Input:** reward_model and language_model

```
1: for prompts in datasets do
2:     seqs = language_model.generate(prompts, do_sample=True)
3:     seqs_max = language_model.generate(prompts, do_sample=False)
4:     rews = reward_model(prompts, seqs) - reward_model(prompts, seqs_max)
5:     log_probs = language_model(prompts, seqs)
6:     loss = -(log_probs.sum(dim=-1) * rews).mean()
7:     lanugage_model.minimize(loss)
8: end for
```

**Output:** language_model
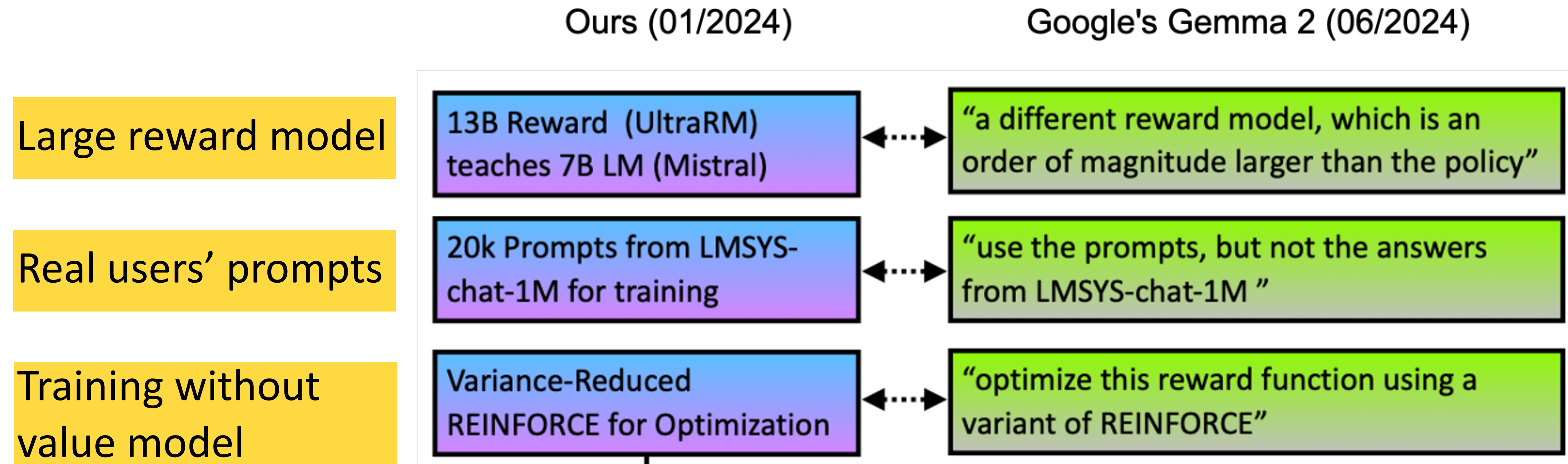
Newly added

ReMax is Simple

8 Lines of code to implement (PPO: 50+)

1 Hyper-parameter (lr) to tune (PPO: 5+)

54

# Comparing with Google's Method

ReMax's training strategies are also used in Google's Gemma 2

Ours (01/2024)                    Google's Gemma 2 (06/2024)

Large reward model

13B Reward (UltraRM) teaches 7B LM (Mistral) ◄┈┈► "a different reward model, which is an order of magnitude larger than the policy"

Real users' prompts

20k Prompts from LMSYS-chat-1M for training ◄┈┈► "use the prompts, but not the answers from LMSYS-chat-1M "

Training without value model

Variance-Reduced REINFORCE for Optimization ◄┈┈► "optimize this reward function using a variant of REINFORCE"
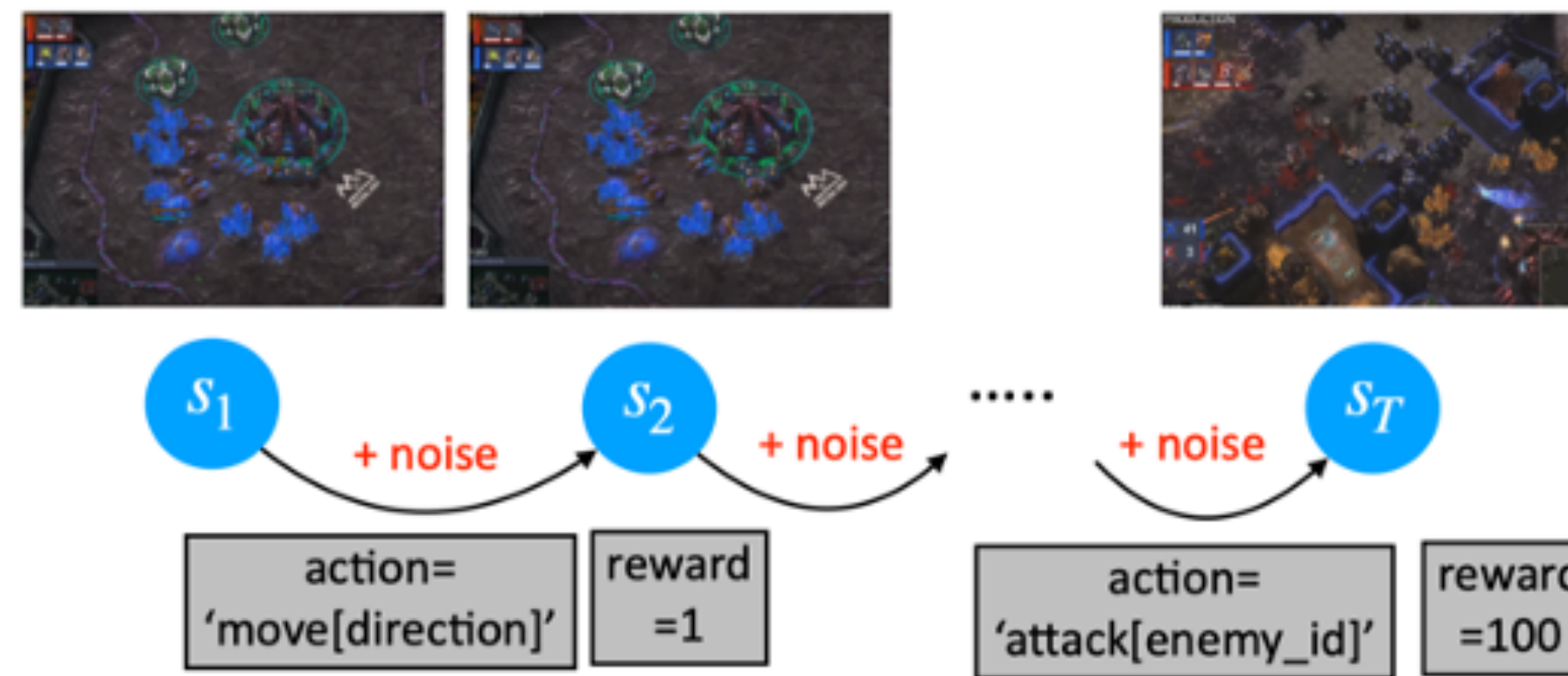
[Team, Gemma, et al. "Gemma 2: Improving open language models at a practical size." arXiv preprint arXiv:2408.00118 (2024).]

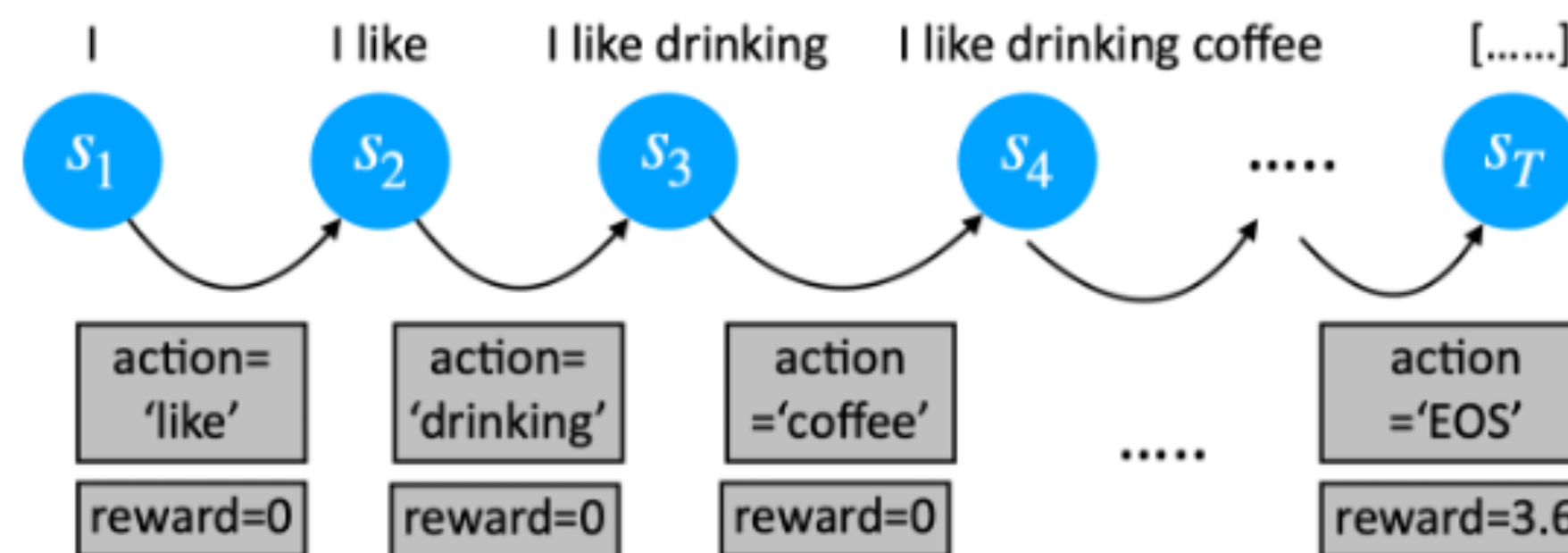# Can We Safely Remove Value Model?



General RL Tasks

RL in LLMs

StarCraft II

[1 hour]

- Slow simulation
- Stochastic transition
- Dense reward

RLHF in LLMs

[10 seconds]

- Fast simulation
- Deterministic transition
- Trajectory-level reward

We conjecture that value-free methods are "optimal" for RL in LLMs

56

# PPO = REINFORCE with Baseline

**General PPO**

$$\mathcal{L}_{\mathrm{ppo}} = \mathbb{E}_{x \sim \rho} \mathbb{E}_{a_{1:T} \sim \pi_{\theta_{\mathrm{old}}}} \left[ \sum_{t=1}^{T} \widetilde{A}(s_t, a_t) \min \left\{ \psi(s_t, a_t), \mathrm{clip}\left( \psi(s_t, a_t), 1 - \delta, 1 + \delta \right) \right\} \right].$$

$$A(s_t, a_t) = \sum_{j=0}^{T-t} \lambda^j \mathtt{advantage}_{t+j} = \sum_{j=0}^{T} \lambda^j \left[ r(s_{t+j}, a_{t+j}) + \gamma V(s_{t+1+j}) - V(s_{t+j}) \right],$$

**Best Practice**    $\gamma = 1, \lambda = 1$

[Ahmadian, Arash, et al. "Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms." *arXiv preprint arXiv:2402.14740* (2024).]
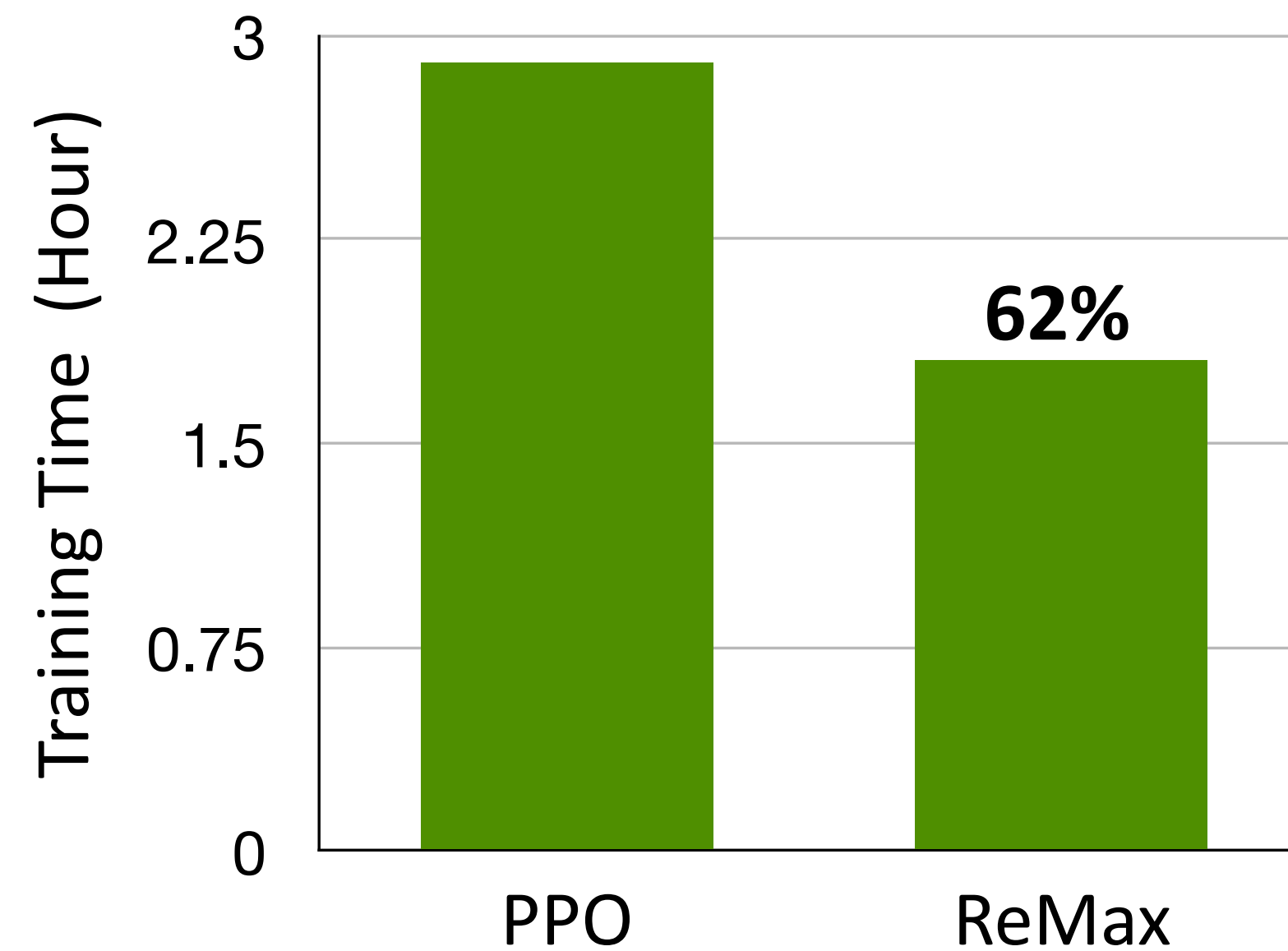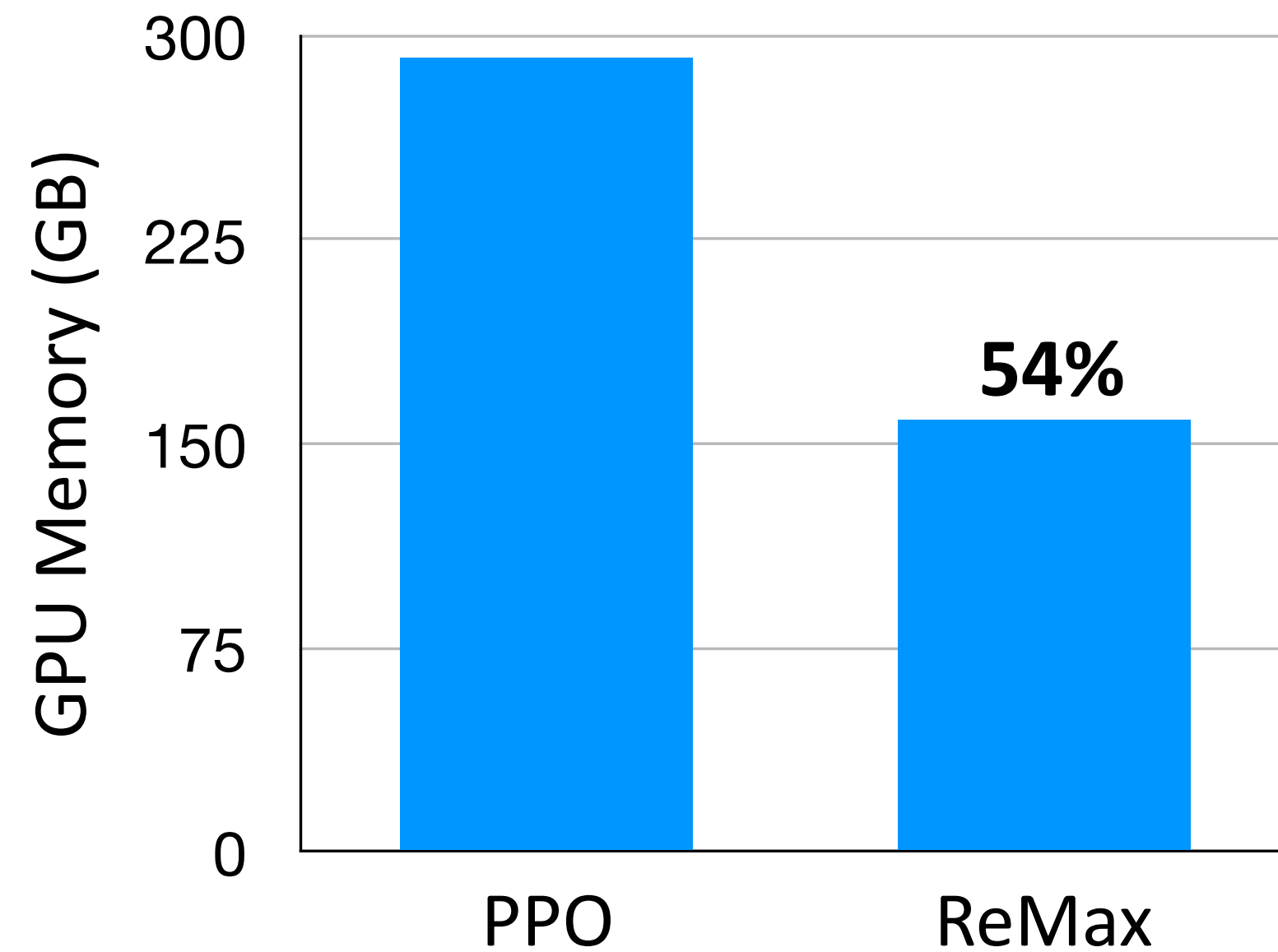
**PPO in LLM**

$$\mathcal{L}_{\mathrm{ppo}}(\theta) = \mathbb{E}_{x \sim \rho} \mathbb{E}_{a_{1:T} \sim \pi_\theta} \left[ \sum_{t=1}^{T} r(x, a_{1:T}) - V(x, a_{1:t}) \right]$$

Outcome reward in REINFORCE's estimator
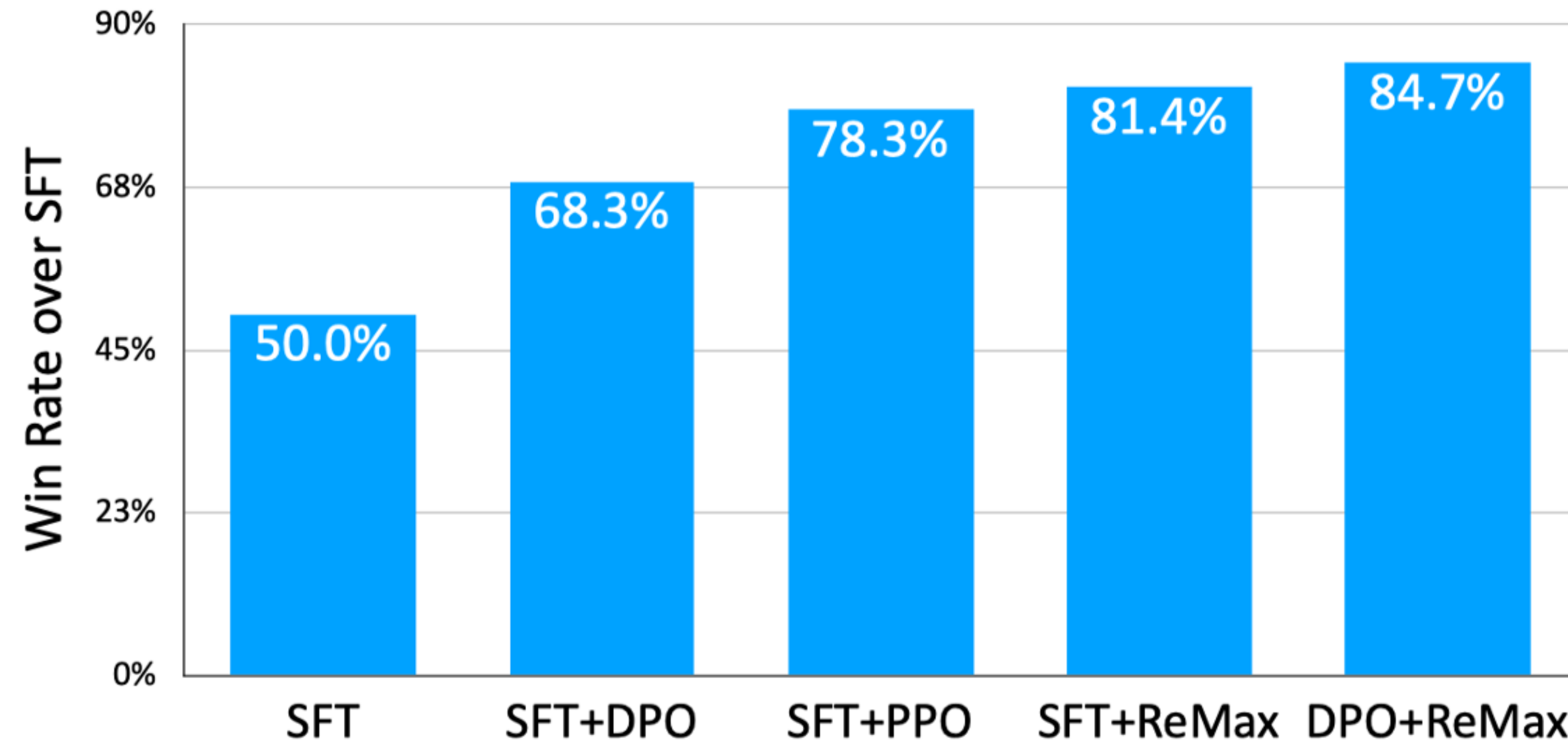
Model-learned Baseline

# ReMax is Computationally Efficient



[Li, Ziniu, et al. "Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models." *arXiv preprint arXiv:2310.10505* (2023).]

ReMax saves about 2x GPU memory and training time on Llama-2-7B

# Performance in RLHF Task



[Li, Ziniu, et al. "Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models." *arXiv preprint arXiv:2310.10505* (2023).]

ReMax is superior to DPO and PPO

# Performance in RLHF Task

Table 4. Performance against strong open-source and private models: Llama-2-Chat models (7B and 70B) apply RLHF (via PPO) using secret datasets; Zephyra-7B-beta (Tunstall et al., 2023) is based on the pretrained Mistral-7B-v0.2 with DPO. GPT-3.5 and GPT-4 utilize RLHF (via PPO) with secret datasets.
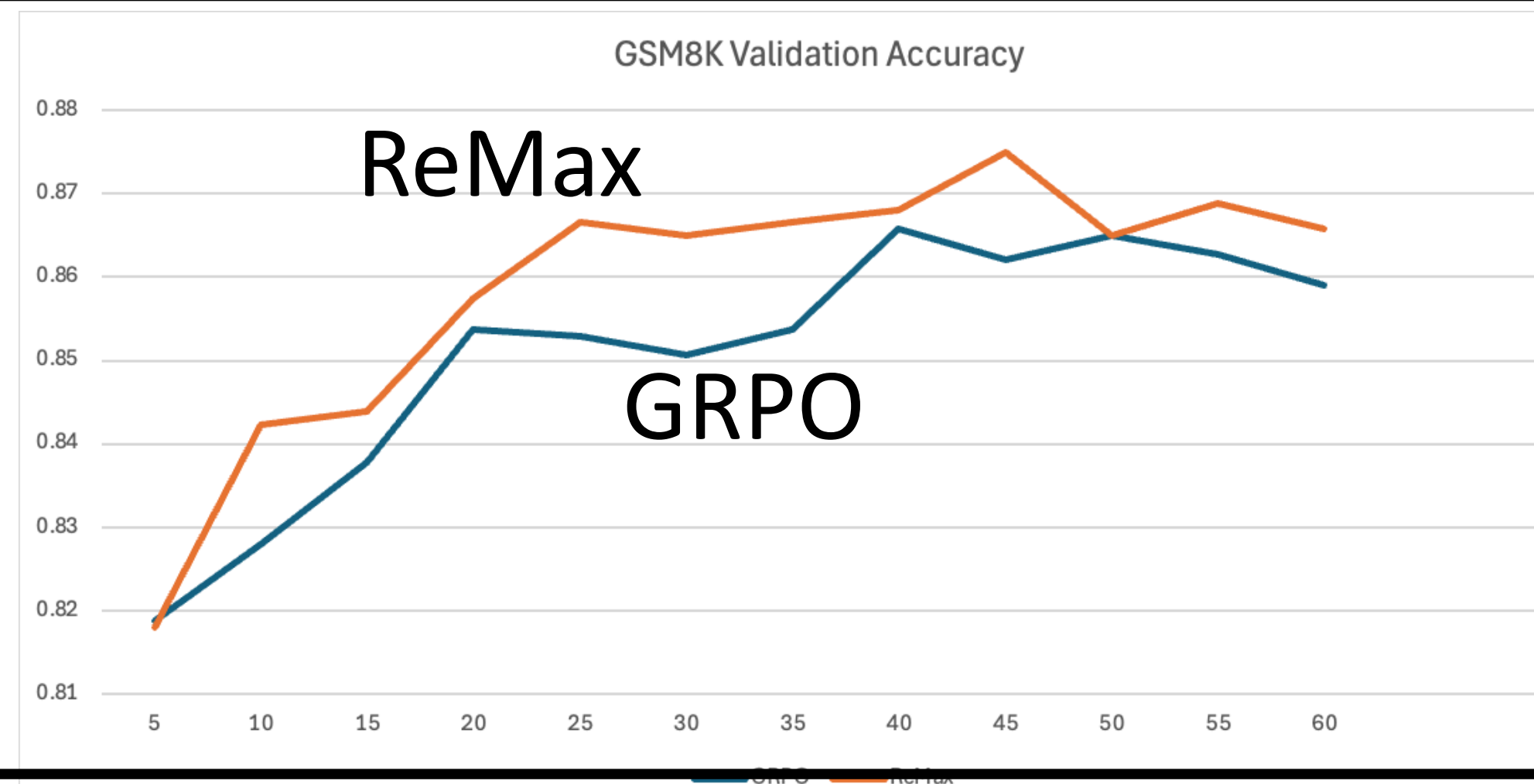
|  | AlpacaEval | MT-Bench |
|---|---|---|
| Llama-2-7B-Chat | 71.37% | 6.269 |
| Zephyr-7B-beta | 90.60% | 7.356 |
| Mistral-7B-Instruct-v0.2 | 92.78% | 7.516 |
| Mistral (via ReMax) | **94.78%** | **7.739** |
| Llama-2-70B-Chat | 92.66% | 6.856 |
| GPT-3.5-turbo | 93.42% | 7.944 |
| GPT-4-turbo | 95.28% | 8.991 |

[Li, Ziniu, et al. "Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models." *arXiv preprint arXiv:2310.10505* (2023).]

ReMax achieves SOTA among 7B models (measured at Jan., 2024)
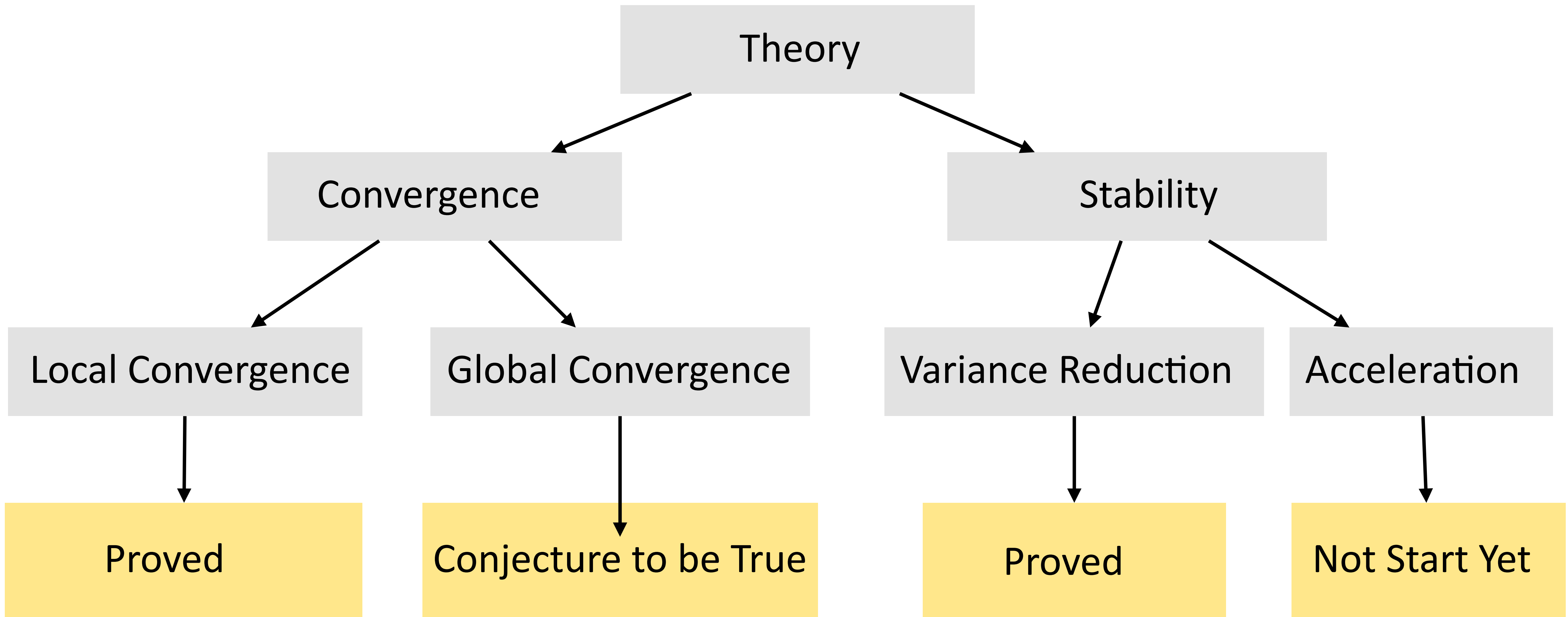
60

# Performance in Reasoning Task

**Our Evaluation**

GSM8K Validation Accuracy

ReMax

GRPO

**Others' Evaluation**

|  | Mineva Math | Olmpiad Bench | HumanEval | LeetCode | LiveCode Bench | Avg. |
|---|---|---|---|---|---|---|
| ReMax | 24.6 | 17.3 | 61.0 | 21.1 | 18.6 | **28.5** |
| GRPO | 22.4 | 20.3 | 57.3 | 13.3 | 18.7 | 26.4 |

[https://curvy-check-498.notion.site/Process-Reinforcement-through-Implicit-Rewards-15f4fcb9c42180f1b498cc9b2eaf896f]

ReMax is superior to DeepSeek's GRPO

# Overview of ReMax's Theory

# Variance Reduction

Setting: 2-action armed bandit (assuming $r(a_1) > r(a_2)$)

Our result: Variance(ReMax) $<$ Variance(REINFORCE) if

$$\pi(a_1) \leq 0.5 + 0.5 \frac{r(a_1)}{r(a_1) - r(a_2)}$$

Implication:
1) variance reduction when the optimal action is **not dominated**
2) slow convergence when the policy is near-optimal
   $\rightarrow$ good if reward is imperfect (**mitigating overfitting**)

# ReMax: A Simple, Effective, and Efficient Reinforcement Learning Method for Aligning Large Language Models

Ziniu Li [1,2]  Tian Xu [3,4]  Yushun Zhang [1,2]  Zhihang Lin [1]  Yang Yu [3,4,5†]  Ruoyu Sun [1,6,2†]  Zhi-Quan Luo [1,2]

ICML 2024

Paper

Code

# Conclusive Remark

**Part I: LLM Training Pipeline**

- ▸ Pre-training: knowledge acquisition
- ▸ Post-training: instruction following and ability enhancement

**Part II: Preserving Diversity in SFT**

- ▸ CE's formulation lack of consideration of diversity
- ▸ GEM: a game-theoretic approach with entropy regularization

**Part III: Efficient RL Training**

- ▸ PPO's formulation are overshot for LLM
- ▸ ReMax: variance-reduced REINFORCE

# Thank You!