

# Preserving Diversity in Supervised Fine-tuning of Large Language Models

Ziniu Li

The Chinese University of Hong Kong, Shenzhen

2025-03-23

# Overview of This Talk

---

Evolution of Large Language Models

Key Differences Between LLMs and Traditional Deep Learning

Our Research Contributions

Key Scientific Insights

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

Ziniu Li<sup>1,2</sup>, Congliang Chen<sup>1,2</sup>, Tian Xu<sup>3</sup>, Zeyu Qin<sup>4</sup>, Jiancong Xiao<sup>5</sup>,  
Zhi-Quan Luo<sup>1,2</sup>, and Ruoyu Sun<sup>1,2,†</sup>

ICLR 2025

NeurIPS 2024 FITML Workshop Best Paper Runner-up



Ziniu Li  
(CUHKSZ)



Congliang Chen  
(CUHKSZ)



Tian Xu  
(NJU)



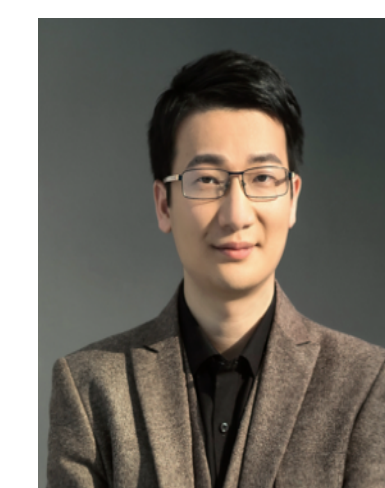
Zeyu Qin  
(HKUST)



Jiancong Xiao  
(Upen)



Zhi-Quan Luo  
(CUHKSZ)



Ruoyu Sun  
(CUHKSZ)

This Talk

Why do we study this topic?

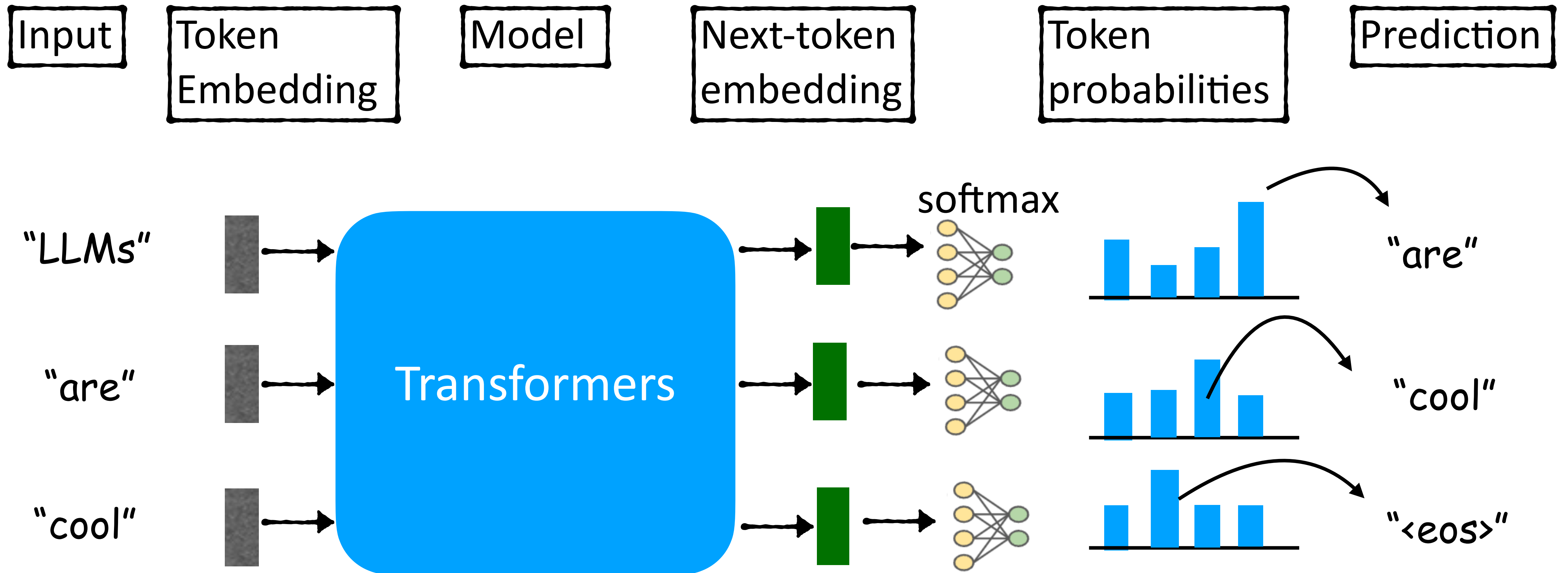
Is it practically impactful?

How do we design our approach?

Are there new scientific discoveries?

# Part I: Overview of LLMs

# LLMs and Transformers



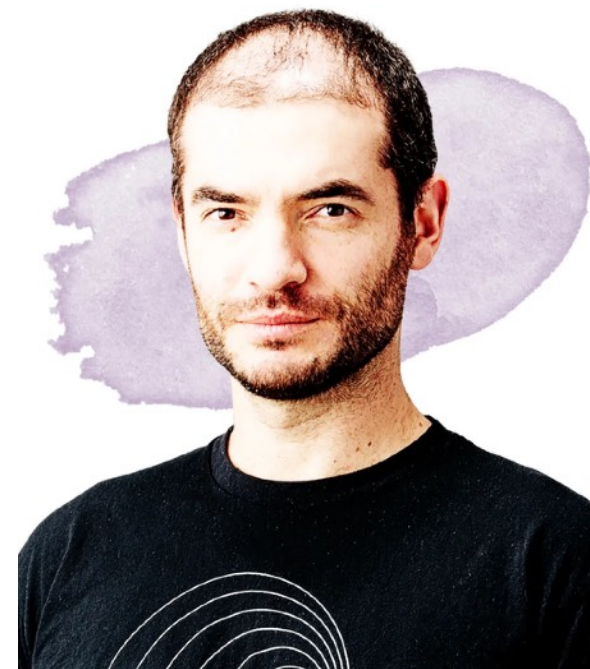
Transformers perform **next-token-prediction**

# LLM Pre-training

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

“Textbooks” can cover:

linguistics  
world knowledge  
common sense  
math reasoning



Ilya Sutskever  
(Godfather of ChatGPT)

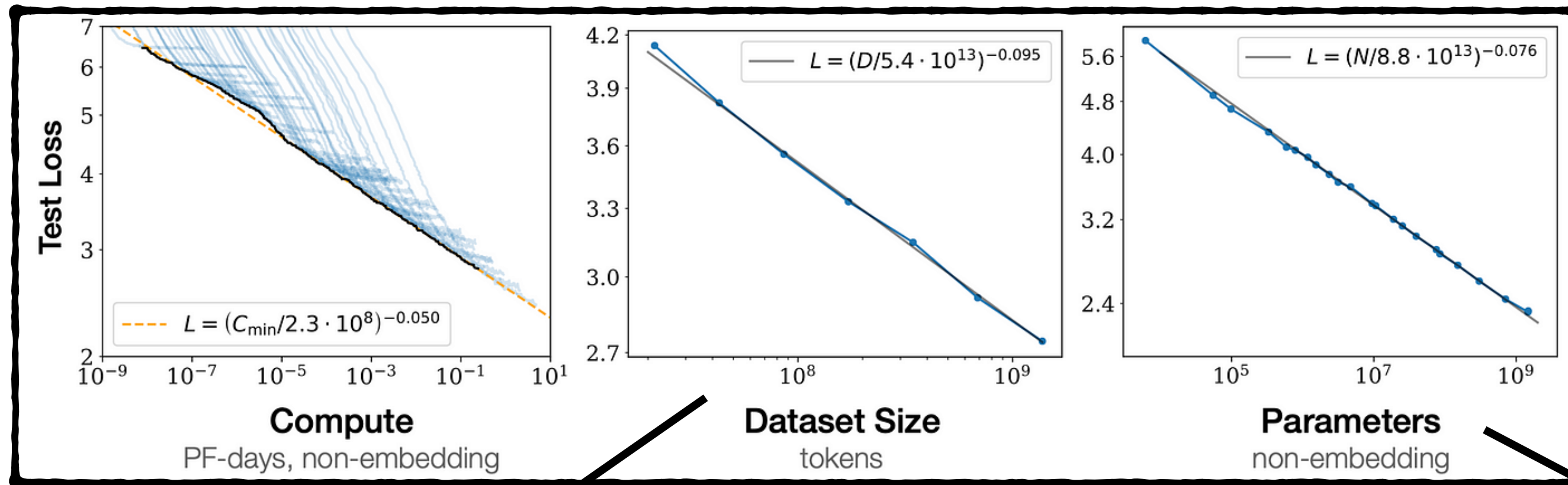
Next-token Prediction is enough for AGI

[[https://www.youtube.com/watch?v=YEUclZdj\\_Sc](https://www.youtube.com/watch?v=YEUclZdj_Sc)]

“Textbook” teaches everything  
(multi-task learning)

# Scaling Law

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

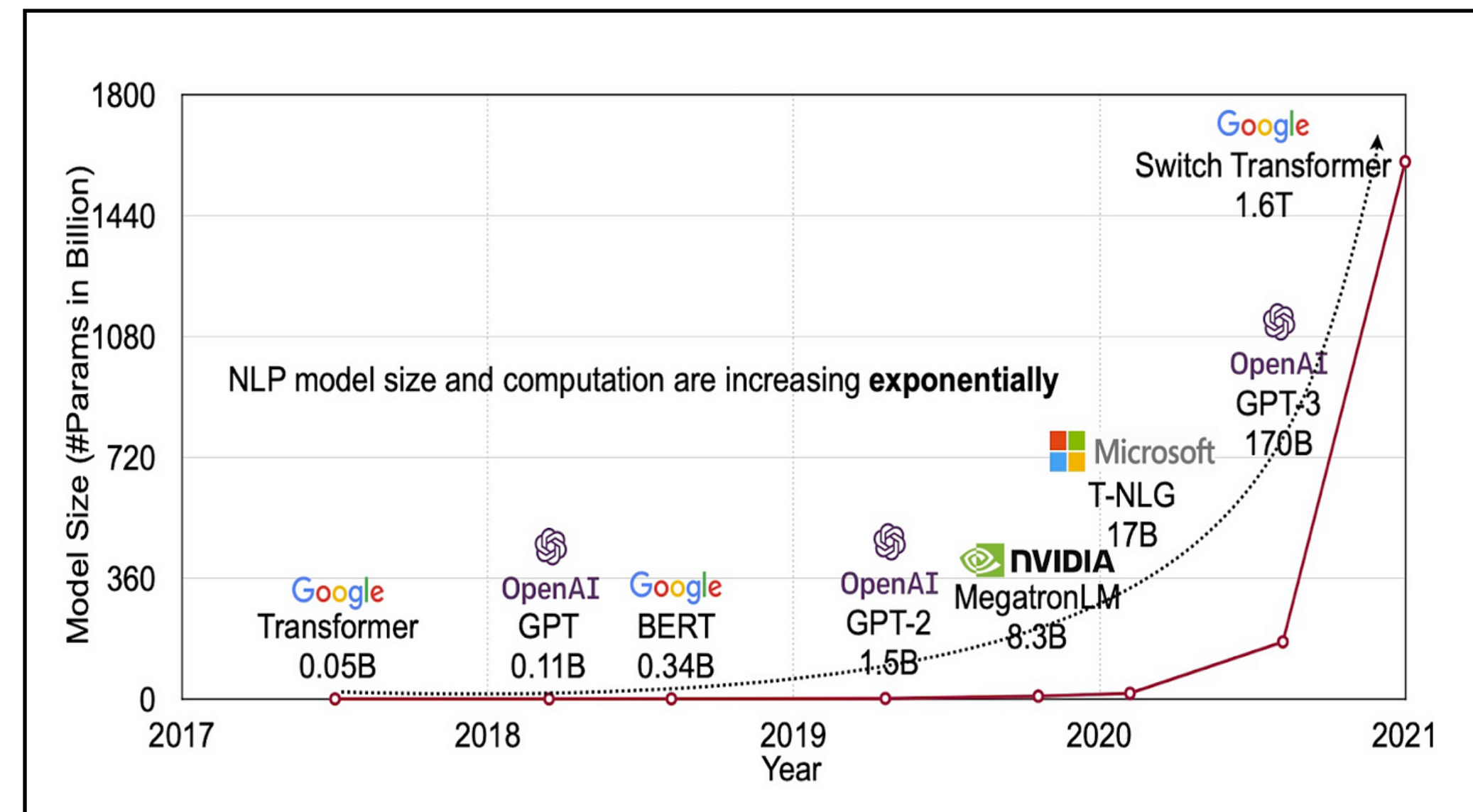
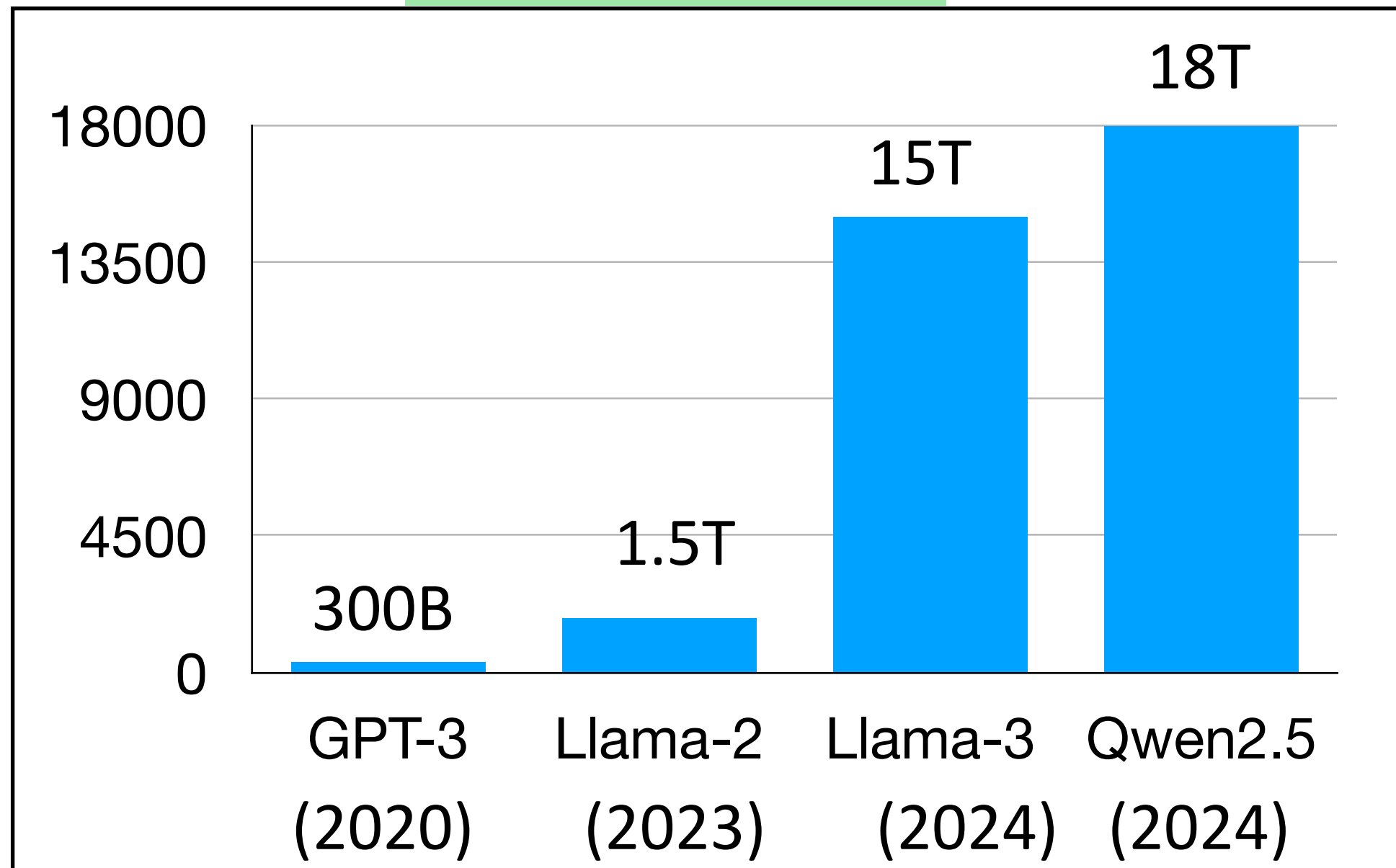


$$L = \frac{A}{D^\alpha} + \frac{B}{N^\beta} + L_0$$

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

dataset size

model size



# Pre-training is not Enough Yet

## Pre-training



Knowledge Acquisition



## Post-training



Ability Reinforcement

**Prompt** : Explain the pre-training of LLMs.

**Llama2-7B**: Explain the pre-training of LLMs.  
Explain the pre-training of LLMs.  
The LLMs are pre-trained on a large amount of unlabeled data, [...]

**repetitive** response

**Pre-trained LLMs: Knowledge Learner without Task Context**

**Prompt** : Explain the pre-training of LLMs.

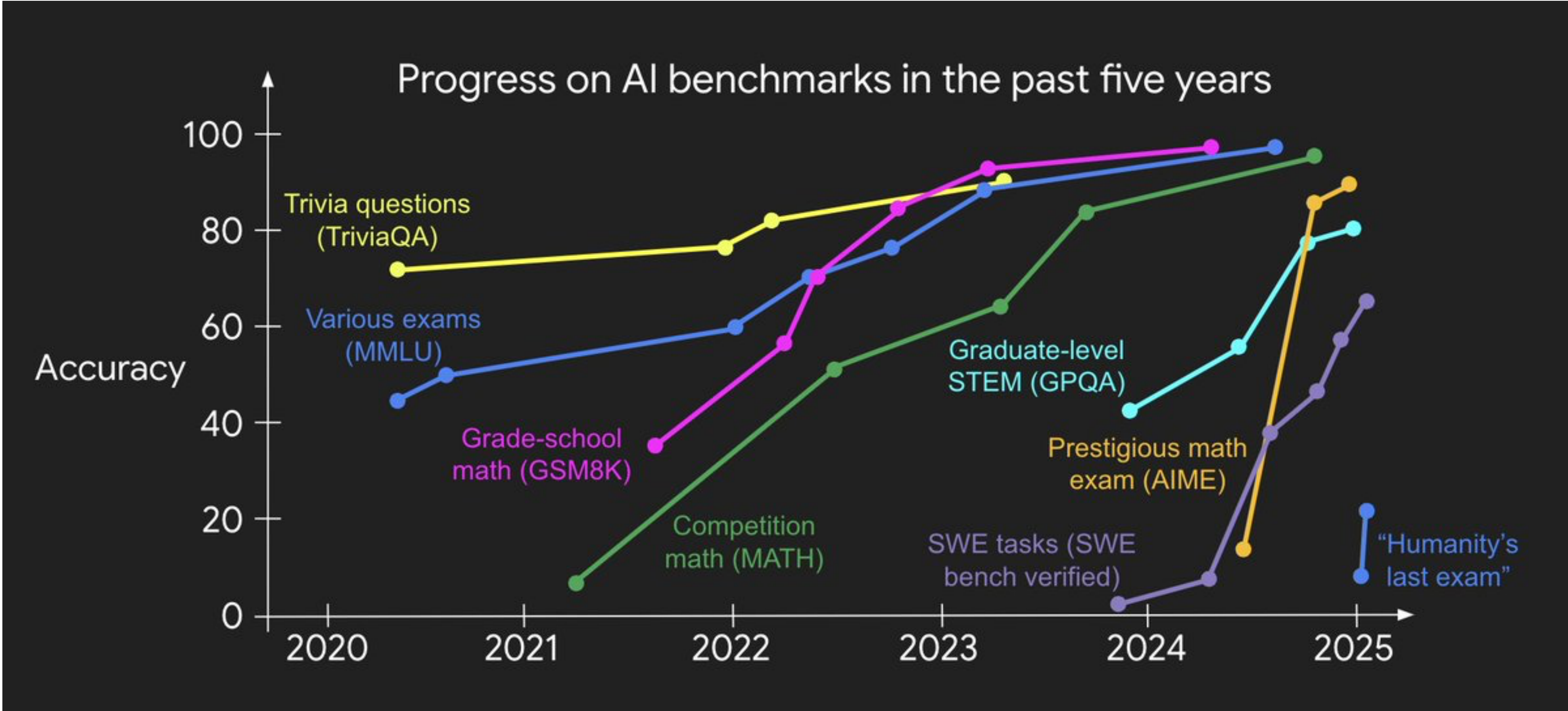
**Llama2-7B-Chat**: Large language models (LLMs) are trained on a variety of tasks and datasets before they are fine-tuned for specific natural language processing (NLP) tasks. Here's an overview of some common pre-training tasks and their goals: [...]

**well-organized** response

**Post-trained LLMs: Enhanced Multi-task Solver**



# Post-training is Powerful



**Post-training** enhances performance for down-stream tasks

# What's Next?



[Talk at NeurIPS 2024]

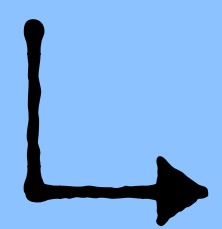
Pre-training as we know it will end

What comes next? The long term is about agentic, reasons, understands, is self aware

**2020**

**(era of GPT-3)**

LLMs are few-shot learners

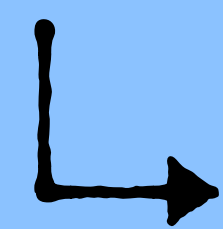


“fine-tuning with few examples is enough”

**2024**

**(era of OpenAI o1)**

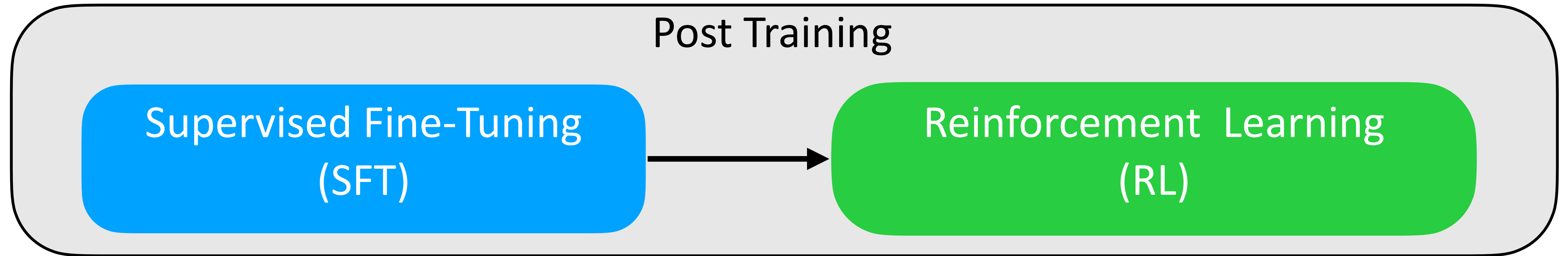
LLMs are strong reasoners



“post-training is equally important as pre-training”

# Part II: Motivation

# LLM Post-Training

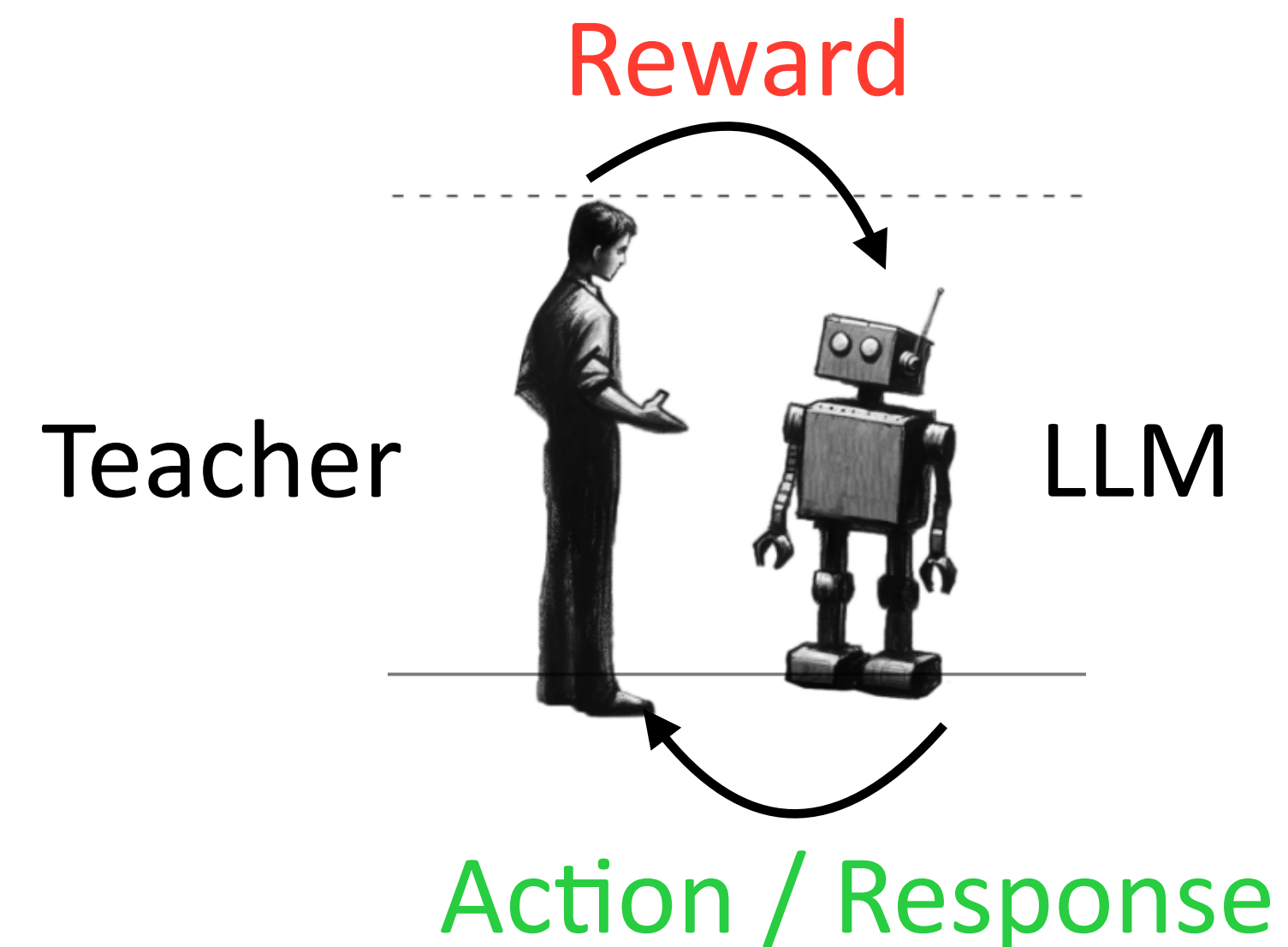
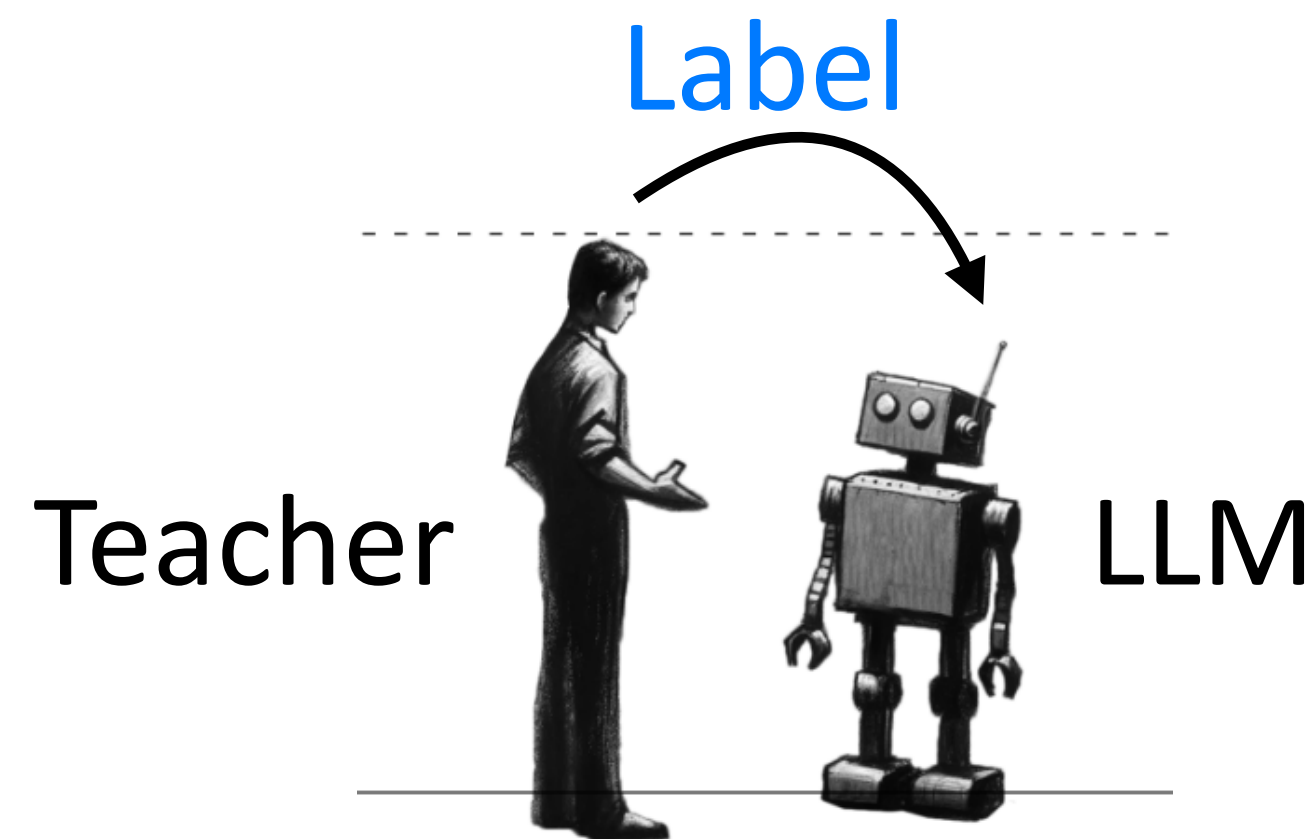


Goal:

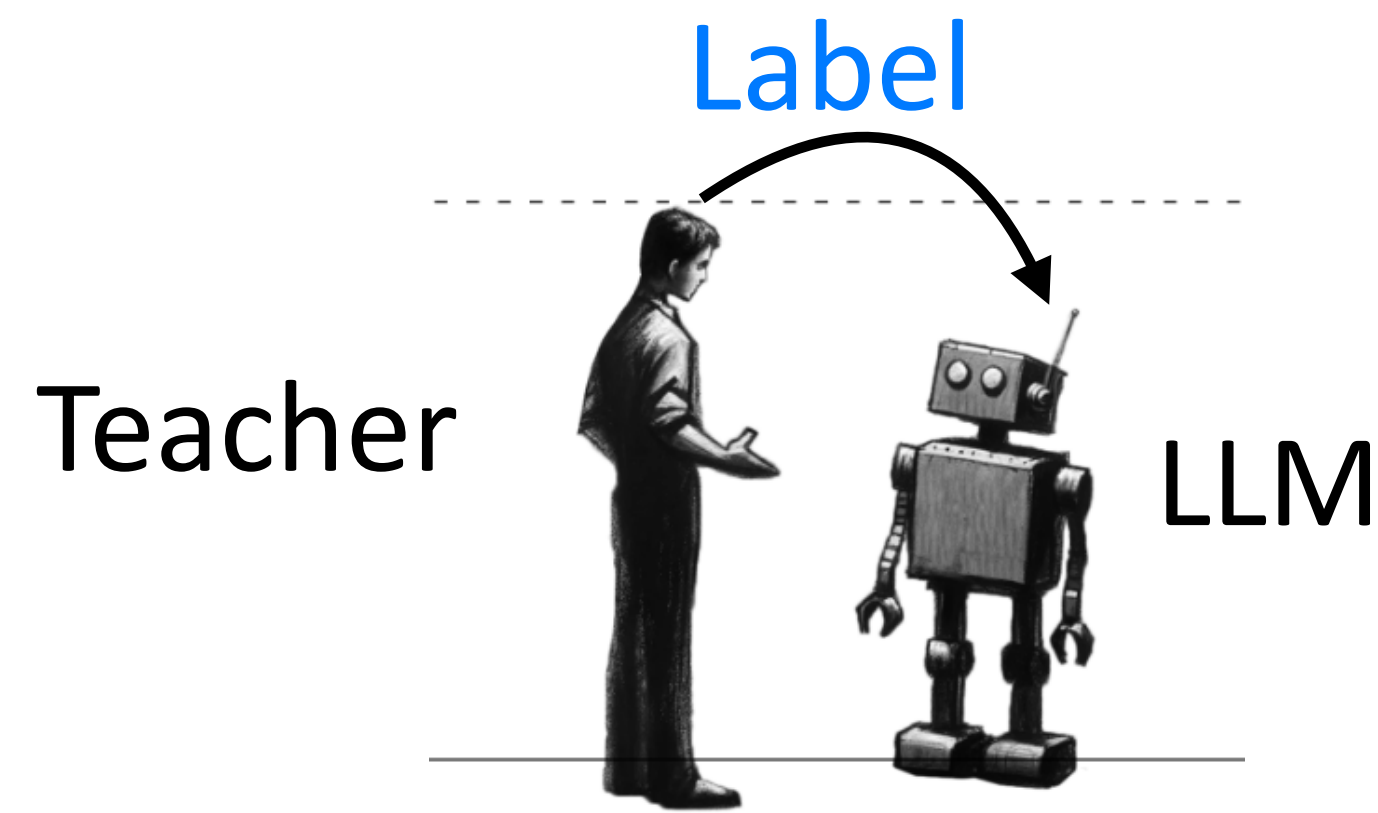
Instruction Learning

Ability Enhancement

Approach:



# Supervised Fine-tuning



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

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

$p$ : data distribution (from teacher)

$f_{\theta}$ : distribution of LLM

SFT Data Example

Prompt

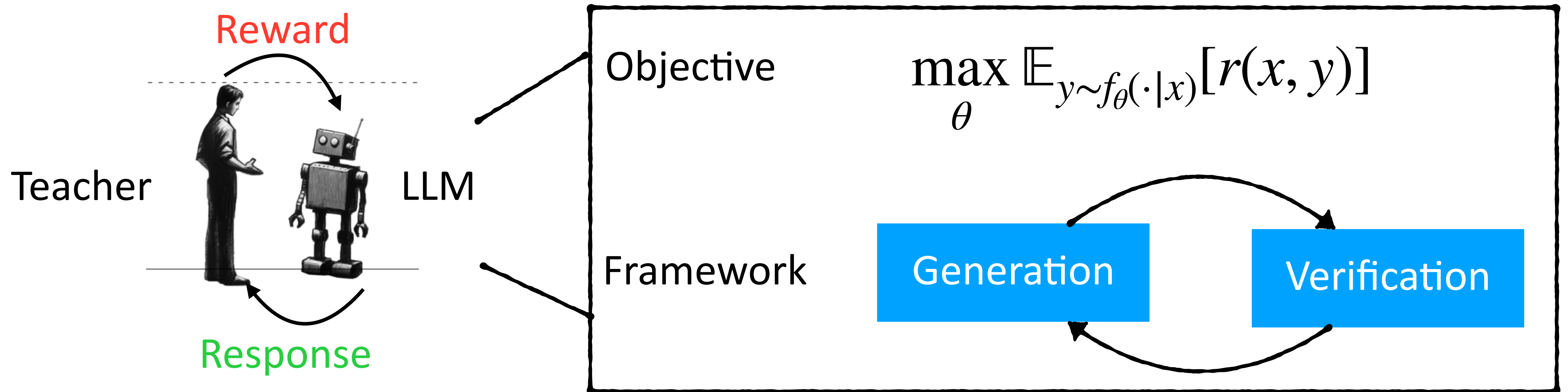
Q: Can Geoffrey Hinton have a conversation with George Washington?

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**

# Reinforcement Learning



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

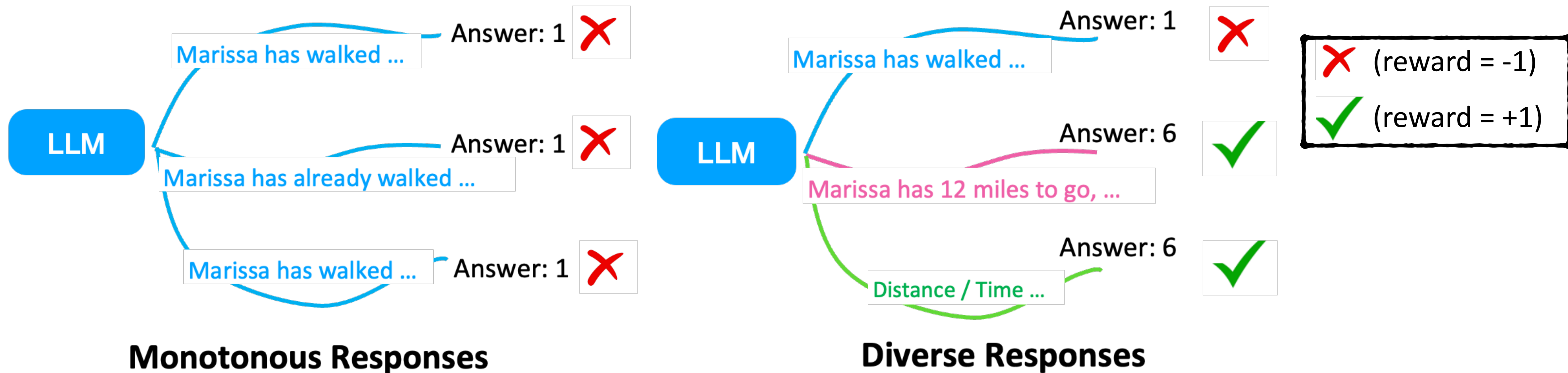
Reward = -1

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

# 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**



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

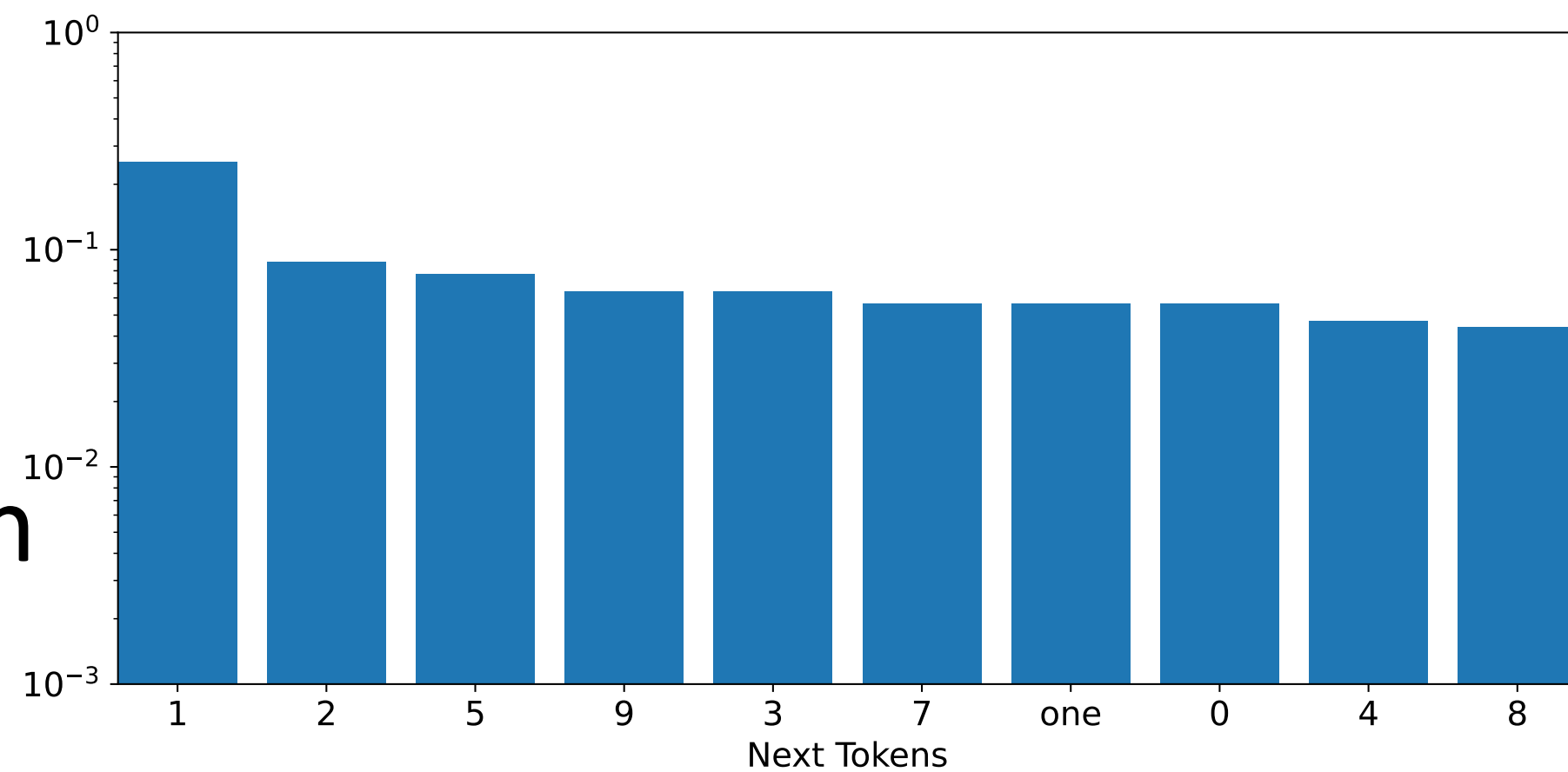
# SFT Reduces Model Output Diversity

#1

Prompt

Give me a single-digit number

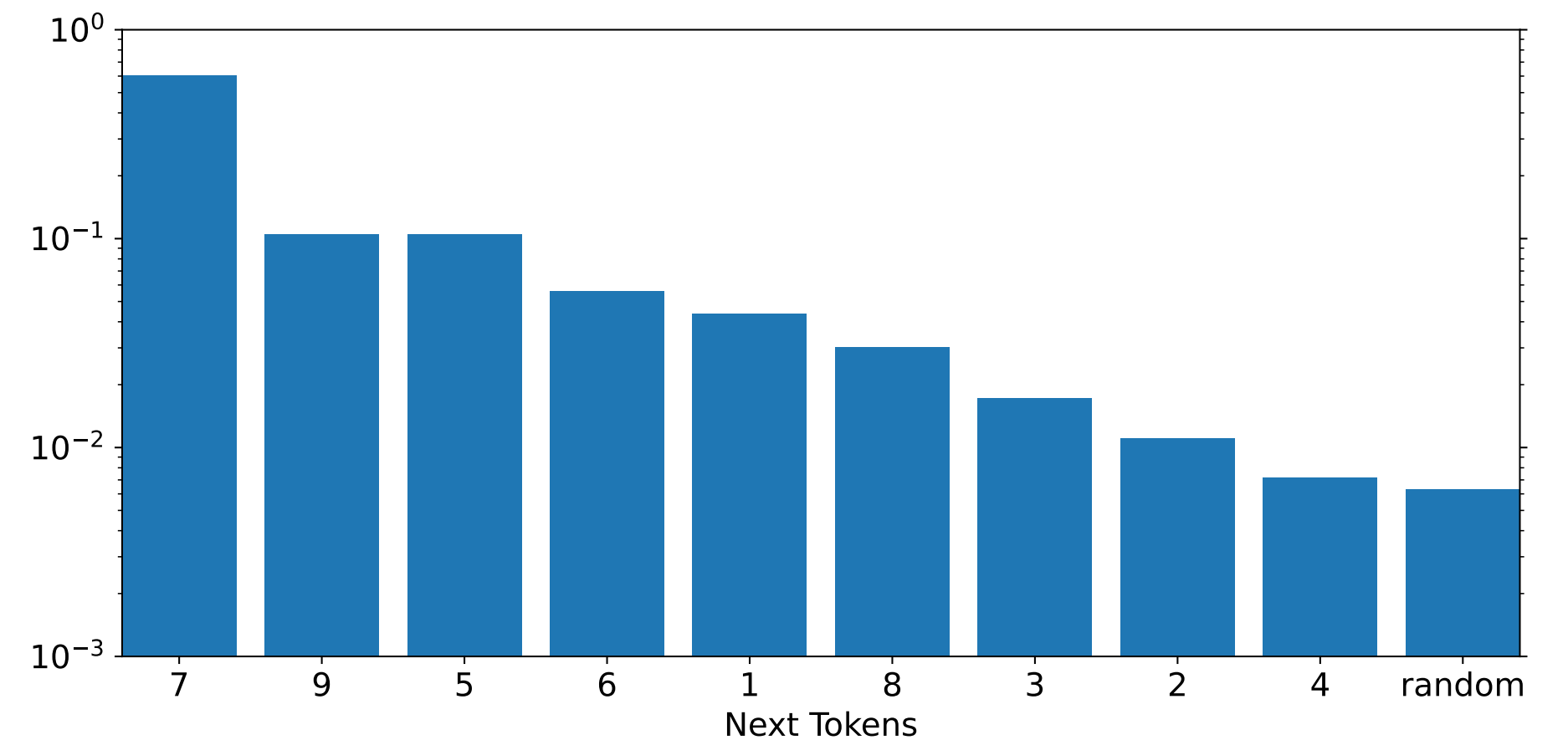
Pre-trained LLM



Response Distribution

“near uniform”

Pre-trained LLM + SFT

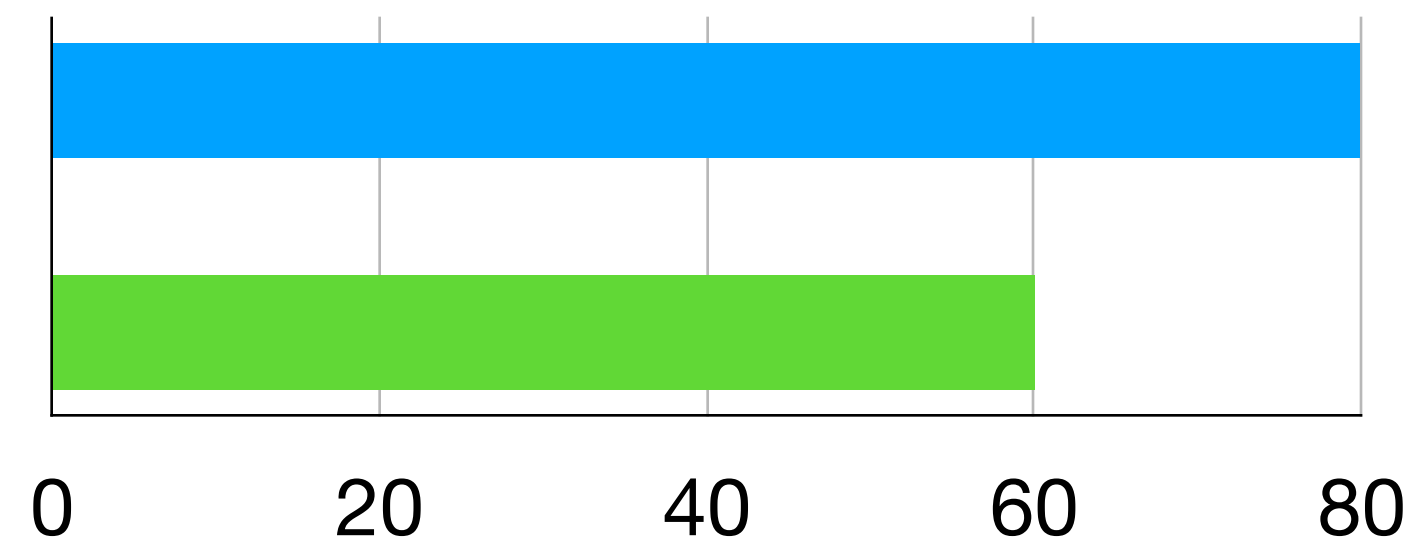


“biased toward 7”

#2

Output Diversity Statistics

Output Diversity



Pre-training SFT

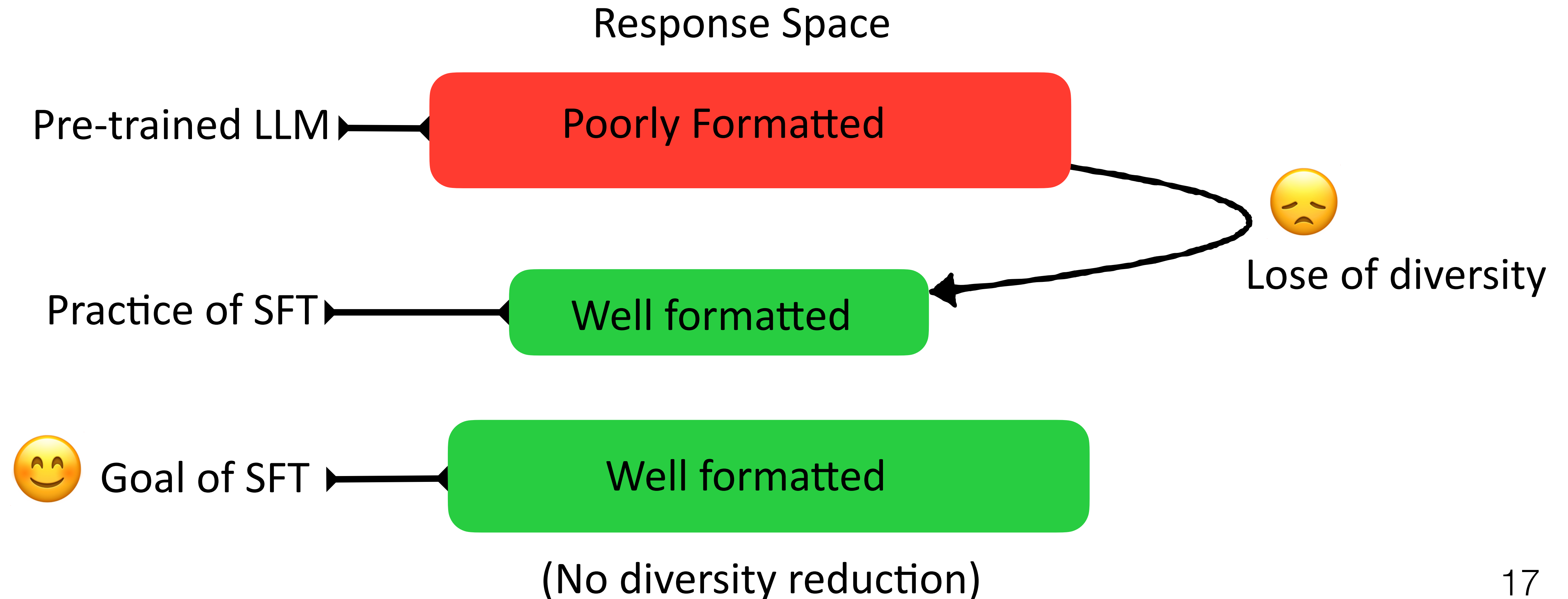
SFT reduces diversity by ~20%

16



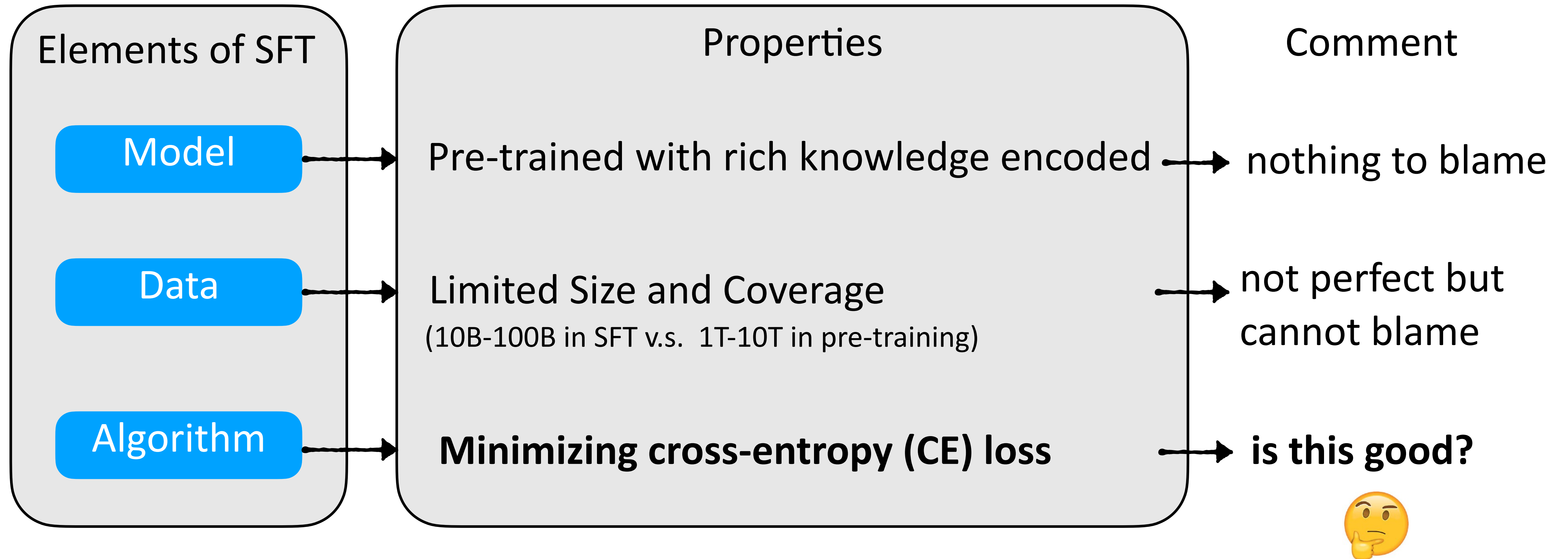
# Revisiting SFT

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



# Why does Diversity Fate in SFT?

---

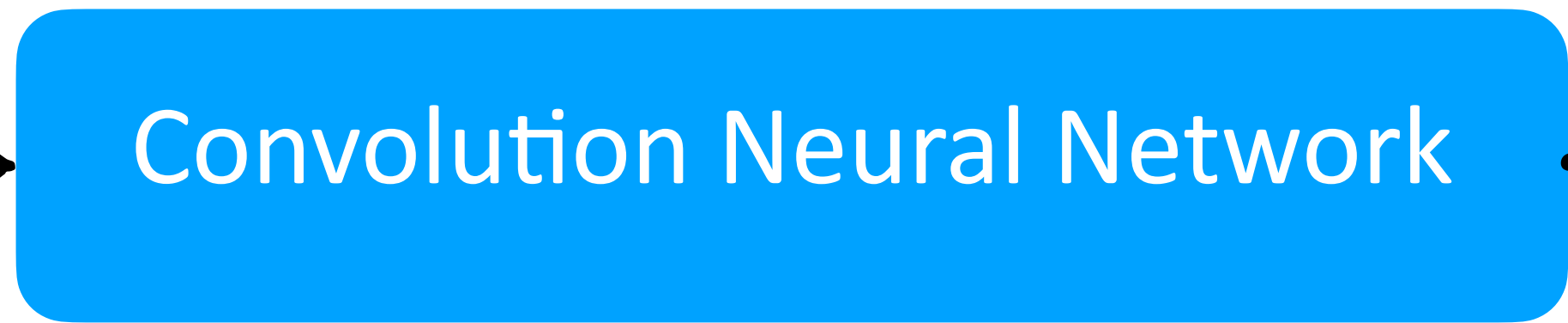


# CE seems Effective for ...

Input



Model



Prediction

“Dog”

Label

“Cat”



CE is Effective for **Classification**

Cross-Entropy Loss

Back-propagation

“I like to drink”



“Tea”

“Coffee”



Is CE Effective for **Generation**?

Cross-Entropy Loss

Back-propagation

# Understanding Generation Tasks

Classification

Generation

Target

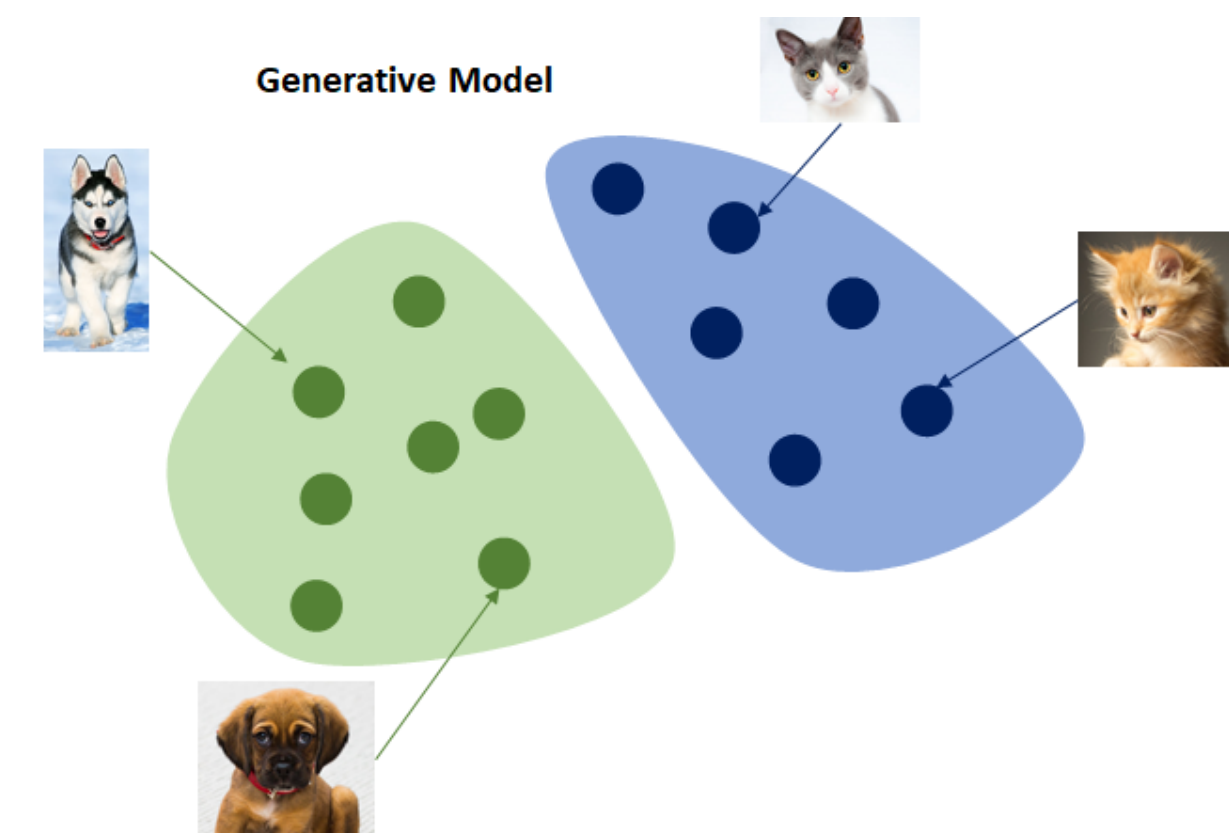
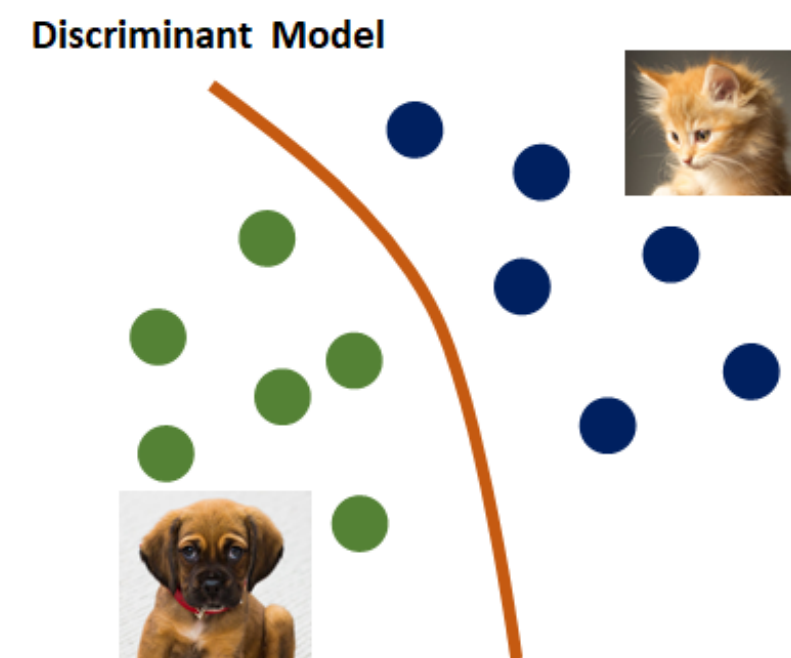
$$\mathcal{X} \mapsto \mathcal{Y}$$

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

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

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

Illustration



Remark for LLMs:

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

# Theory of CE

## CE Loss (Empirical)

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

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

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

$\theta$ : parameters of neural network

## CE Loss (Population)

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

$\rho$ : prompt distribution

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

Equivalence

## Forward KL Divergence

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

CE can be used to learn a distribution

If the data samples are “abundant”



Classification  
(one label sample is enough)



Pre-training  
(huge data)



SFT  
(data is limited)

Distribution Matching

# 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:**

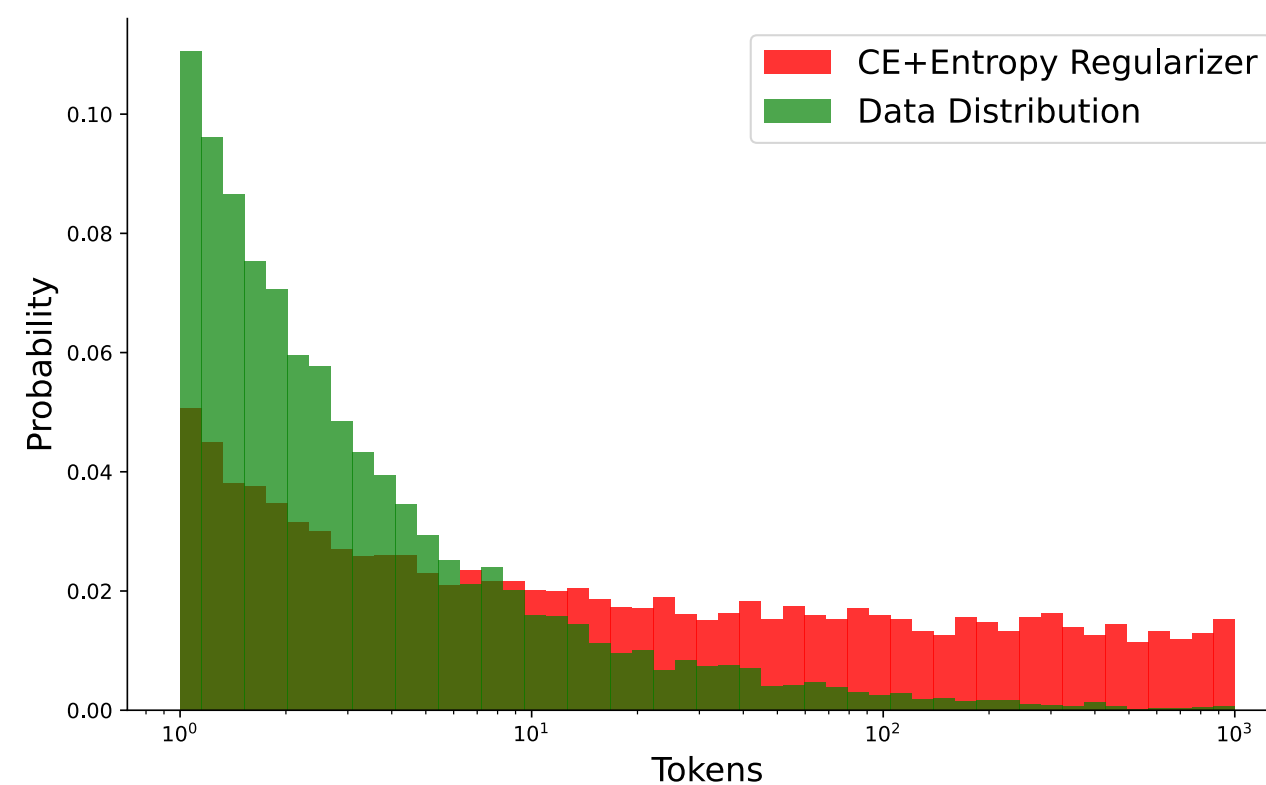
Designing new formulation and algorithm for SFT

# Part III: Our Approach GEM

# A Naive Approach for Diversity

## CE + Entropy Regularization

$$\max_f \underbrace{\mathbb{E}_x \mathbb{E}_{y \sim p(\cdot|x)} [\log f(y|x)]}_{-D_{\text{KL}}(p, f) + \text{constant}} + \beta \underbrace{\mathbb{E}_x \mathbb{E}_{y \sim f(\cdot|x)} [-\log f(y|x)]}_{\mathcal{H}(f)}$$



Toy setting

Prompt	Whats the largest star in our galaxy?
CE	Hello! Atlantis is a legendary city that was said to have existed in ancient Greece. According to the story, it was a highly advanced and prosperous city that was located on an island in the ocean. [...]
CE+Entropy	Hello! Atlantis Documentary is a 2019 American <b>documentary</b> <b>éhoFLICT</b> film directed by <b>Já oblík</b> and produced by Werner Herzog. The film explores the history and legacy of Atlantis, <b>□</b> an ancient Greek city-state that was said to <b>have_calendar</b> knowledge and advanced technology, through interviews with scholars and <b>histori-ans.ython</b>

LLMs

Entropy regularizer encourages diversity via increasing the **tail** of distribution





# Analyzing Cross-Entropy Loss

---

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

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

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

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,*

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

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

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

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

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

# Limitations of CE

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

- Find any  $j$  with  $f_{\theta_k}(j|x) > 0$
- Decrease the logit for source token  $j$  by learning rate  $\eta$  and weight  $w_{i \leftarrow j}$ :

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

#2

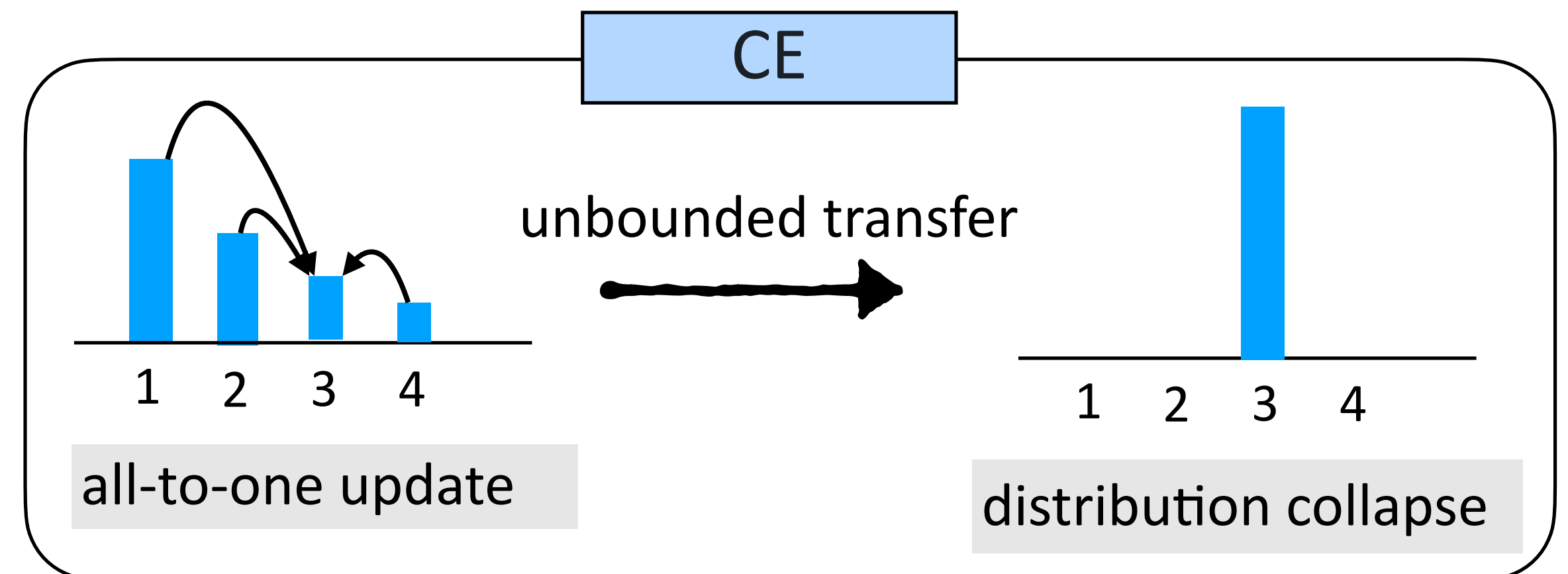
- Increase the logit for the target token  $i$  in a similar manner:

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

Procedure of CE

Limitation 1: Unbounded Transfer

Limitation 2: All-to-one Update



# Proposed Solutions

## Procedure of Our Method

#1

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

#2

- Calculate the model's best prediction  $j = \operatorname{argmax} f(\cdot|x)$
- Decrease the logit for source token  $j$  by learning rate  $\eta$  and weight  $w_{i \leftarrow j}$ :

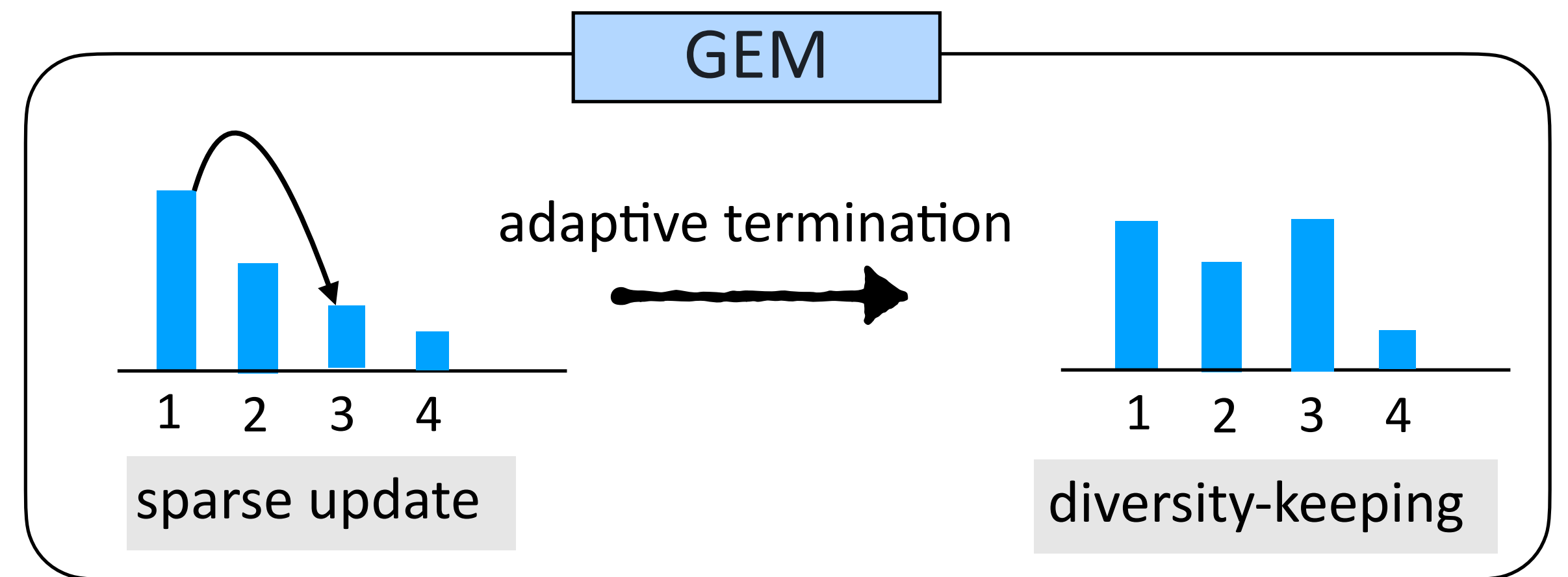
$$\theta_{k+1}[j] = \theta_k[j] - \eta * w_{i \leftarrow j}$$

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

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

Technique 1: Adaptive Termination

Technique 2: Sparse Update



# Our Insight: Dimension Increase

## Procedure of Our Method

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

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

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

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

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



# 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) possesses a unique Nash equilibrium point:

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

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

$$f^* = \underset{f}{\operatorname{argmin}} \mathbb{E}_x [D_{\text{KL}}(f(\cdot|x), p(\cdot|x)) - \gamma \mathcal{H}(f(\cdot|x))]. \quad (8)$$

**Terminology**

**Reverse 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  $\rightarrow$  future work



# Training Algorithm

Idea: block-wise gradient-descent and coordinate descent

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

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

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

↳ 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

**Challenge**

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**

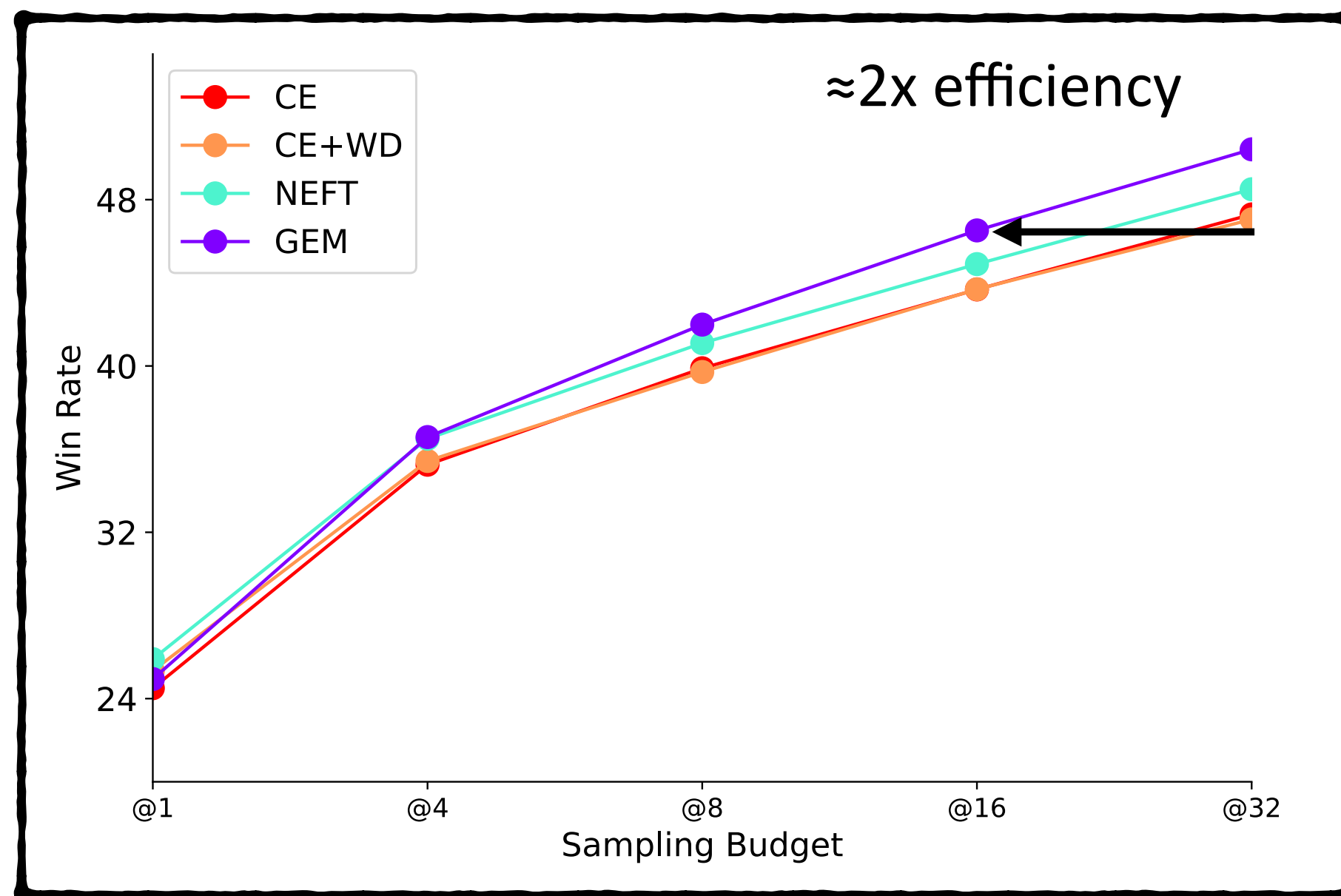
High

Low

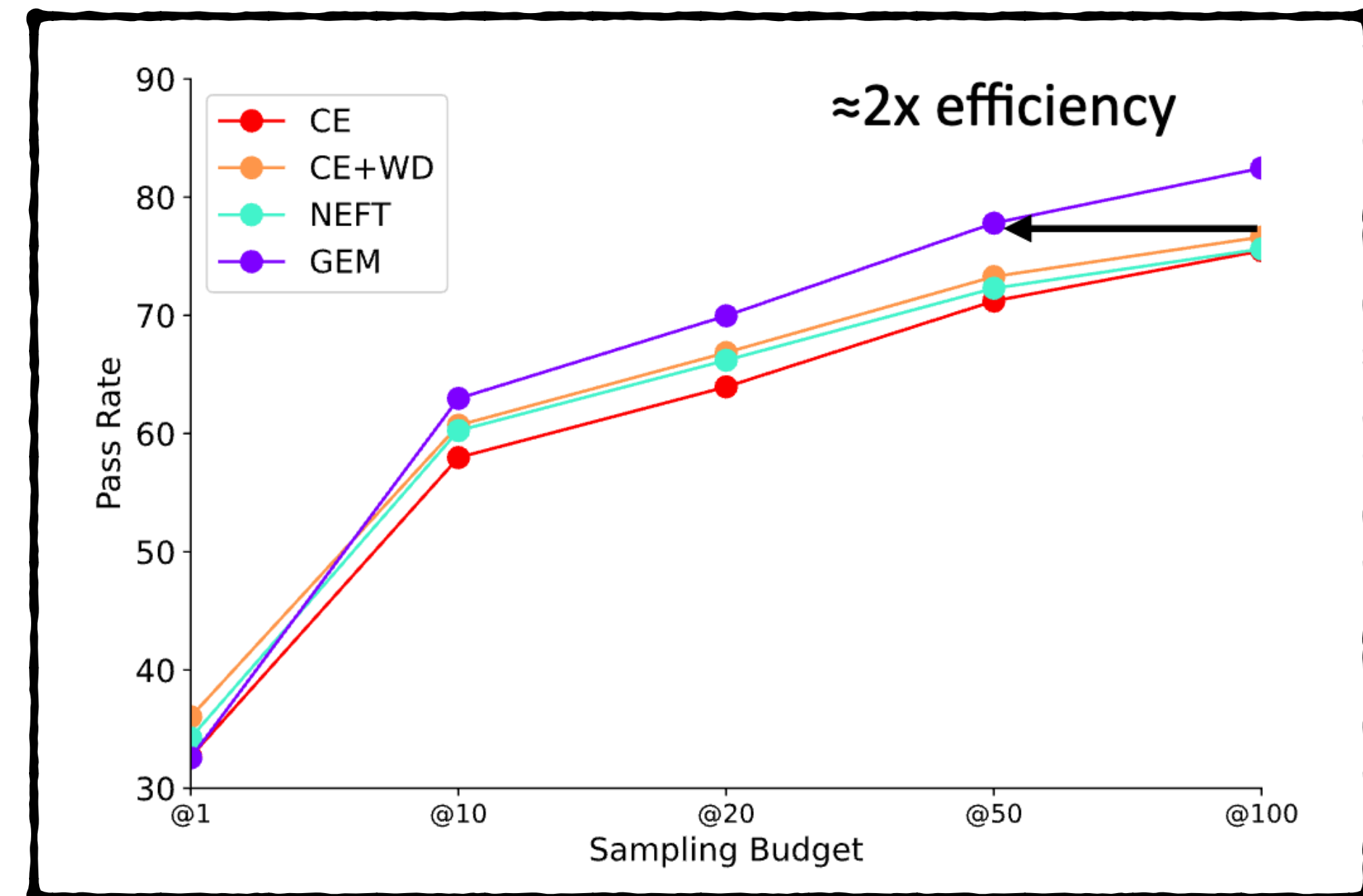
# Part IV: Empirical Results

# Test-Time Scaling

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



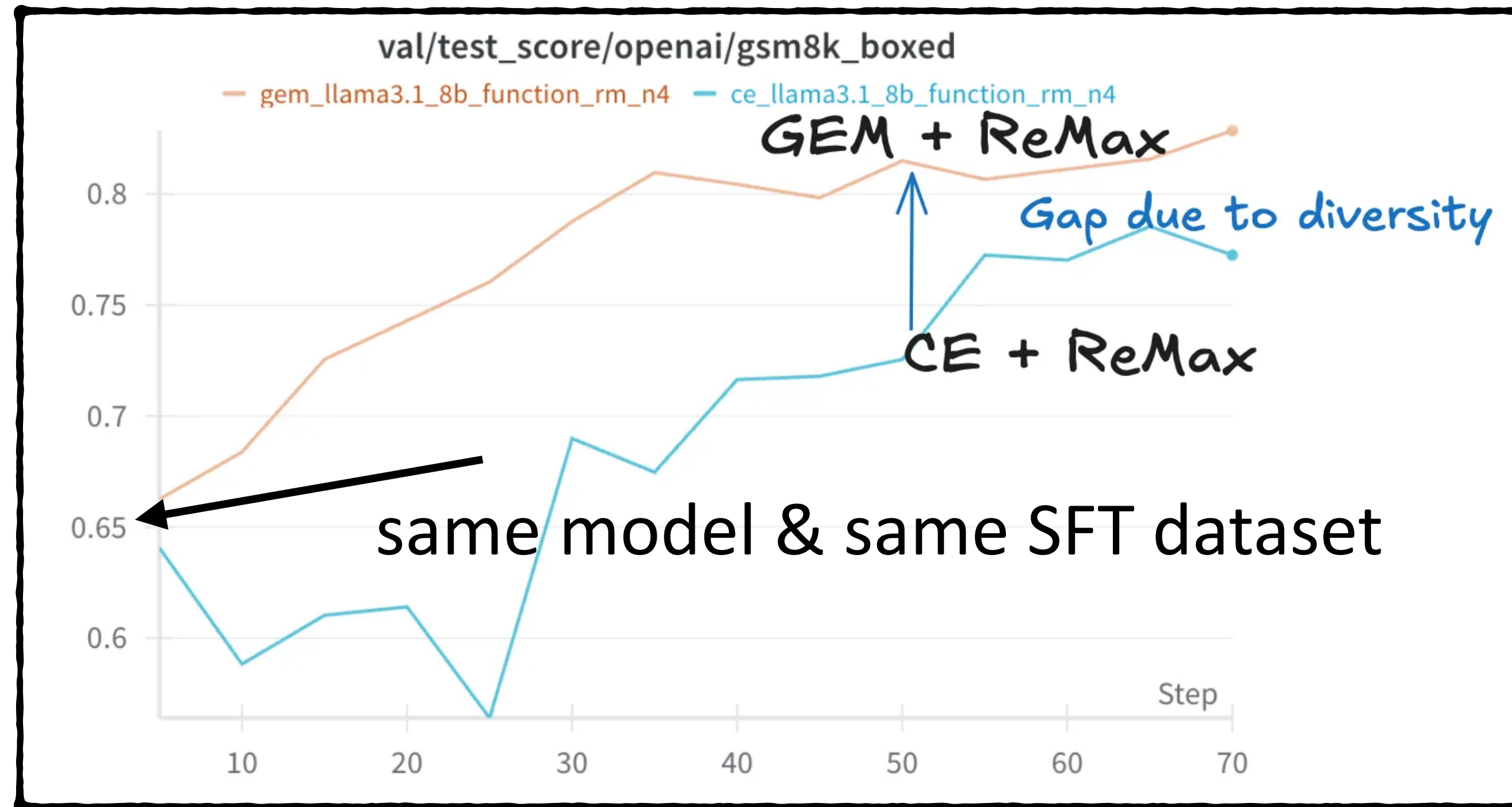
RLHF Alignment (Chat)



Code Generation

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

# Math Reasoning



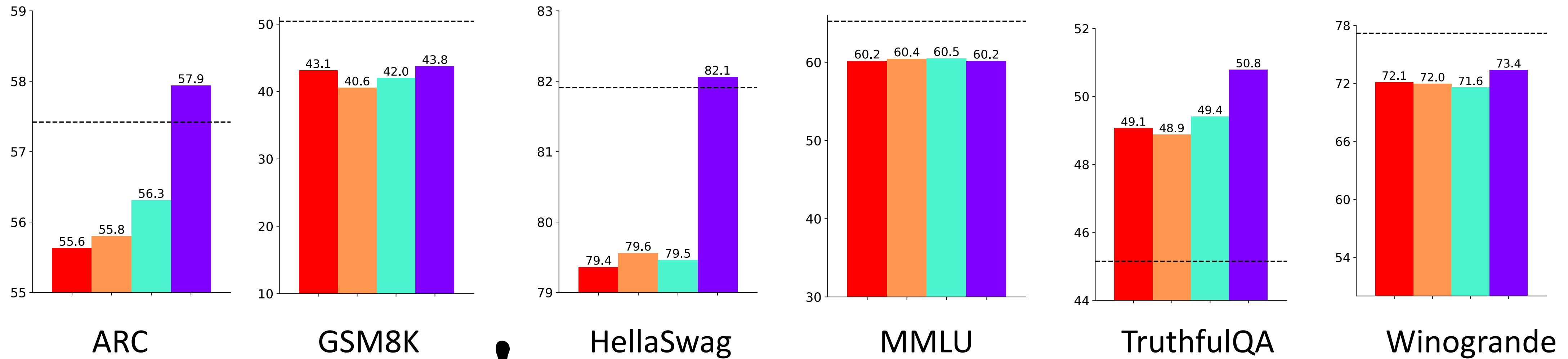
[<https://tangible-polo-203.notion.site/>]

- ▶ 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.]

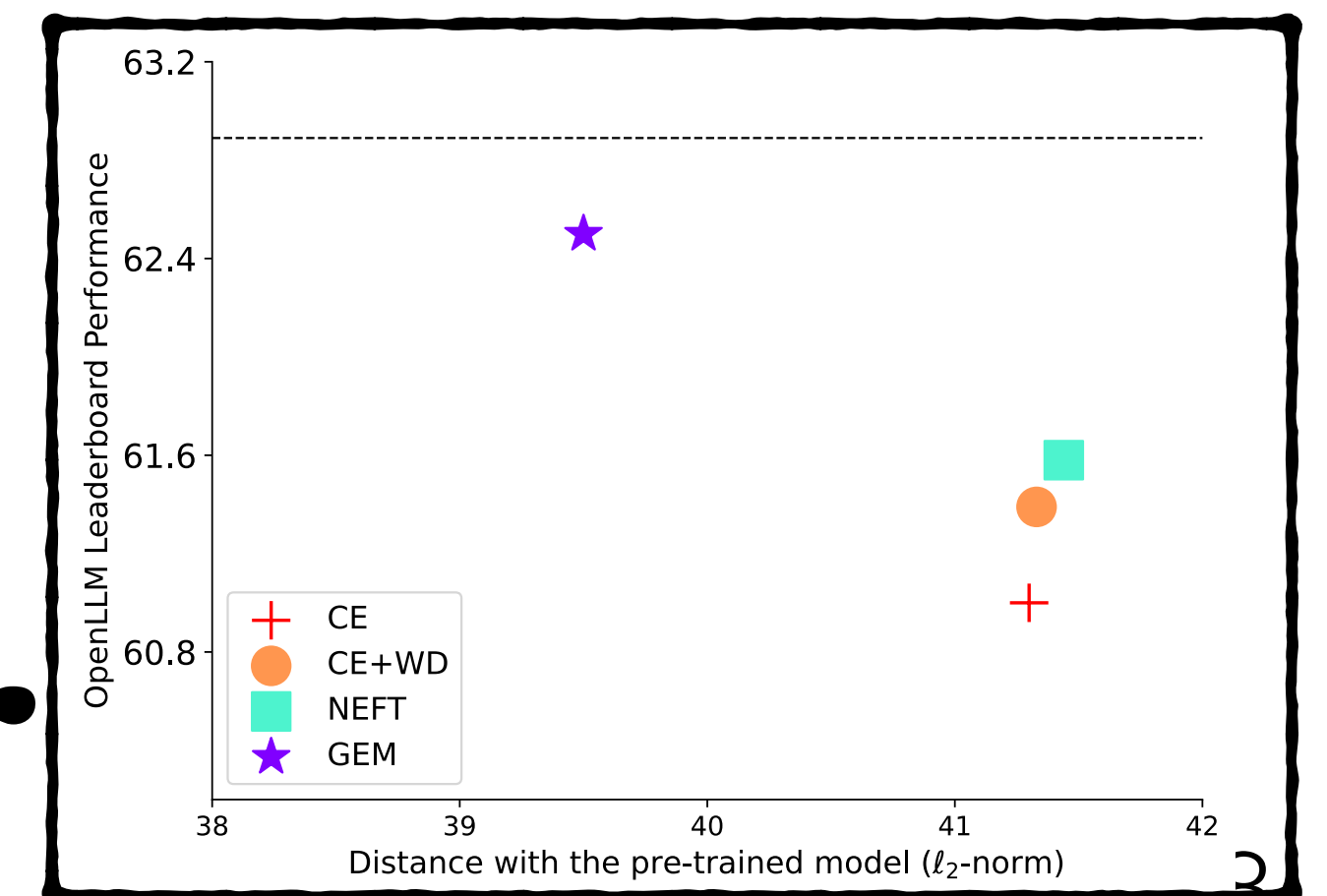
GEM improves the performance limit of RL training

# Alignment Tax



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

GEM-tuned model shows less overfitting to the data



# Thank You!



Paper



Code