



INFO-I590 Fundamentals and Applications of LLMs

# **DPO and GRPO**

**Direct Preference Optimization (DPO)**

**Group Relative Policy Optimization (GRPO)**

Jisun An

# Direct Preference Optimization (DPO)

# RLHF optimization objective

$$\max_{\pi} \mathbb{E}_{x \sim D, y \sim \pi} [r(x, y) - \beta D_{KL}(\pi(y | x) \parallel \pi_{\text{ref}}(y | x))]$$

Maximizes the rewards

Prevents the model from changing too drastically

Current aligned LLM

SFT (instruction-tuned) LLM

How to get rid of this reward?  
(not get rid of it, but turn it into a probability)

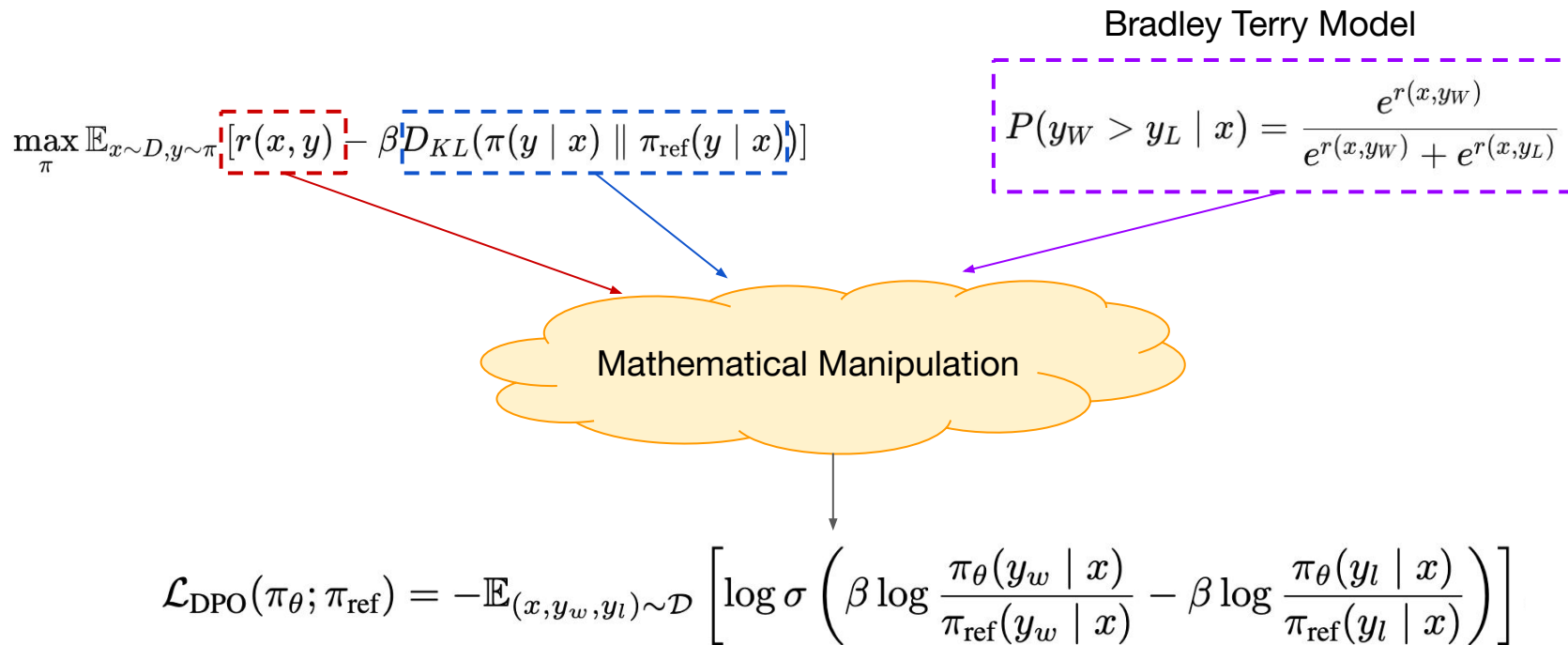
Can we avoid rollouts?

The diagram illustrates the RLHF optimization objective. The equation is  $\max_{\pi} \mathbb{E}_{x \sim D, y \sim \pi} [r(x, y) - \beta D_{KL}(\pi(y | x) \parallel \pi_{\text{ref}}(y | x))]$ . A red dashed box highlights the reward term  $r(x, y)$ , with a red label 'Maximizes the rewards' above it. A blue dashed box highlights the KL divergence term  $D_{KL}(\pi(y | x) \parallel \pi_{\text{ref}}(y | x))$ , with a blue label 'Prevents the model from changing too drastically' above it. Two arrows point from the KL term to labels: 'Current aligned LLM' points to  $\pi(y | x)$  and 'SFT (instruction-tuned) LLM' points to  $\pi_{\text{ref}}(y | x)$ . Two purple arrows point from the equation to questions: one from the reward term to 'How to get rid of this reward? (not get rid of it, but turn it into a probability)' and another from the KL term to 'Can we avoid rollouts?'.

# DPO (Direct Preference Optimization)

- Idea: Can we fine-tune a model directly using a preference dataset without RL?
- Why Avoid RL?
  - RL is unstable and computationally expensive.
  - It's hard to pinpoint what to change in long outputs.
  - Generating multiple samples for training (rollouts) increases cost.
  - Training a reliable reward model is challenging.
  - Small reward differences can cause instability.
- Also called as 'Preference Tuning'

# The DPO loss function



# The DPO loss function

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

Aligned model we are training

Maximize probability for good response

Minimize probability for bad response

Don't change the model too much

- No explicit reward model
- No need for rollouts from the policy

(Optional)

## Deriving the DPO loss function

$$\begin{aligned} & \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi} [r(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi(y|x) \parallel \pi_{\text{ref}}(y|x)] \\ &= \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ r(x, y) - \beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} \right] \\ &= \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} - \frac{1}{\beta} r(x, y) \right] \\ &= \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y|x)}{\frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left( \frac{1}{\beta} r(x, y) \right)} - \log Z(x) \right] \quad (12) \end{aligned}$$

where we have partition function:

$$Z(x) = \sum_y \pi_{\text{ref}}(y|x) \exp \left( \frac{1}{\beta} r(x, y) \right).$$

\* Note that the partition function is a function of only  $x$  and the reference policy  $\pi_{\text{ref}}$ , but does not depend on the policy  $\pi$ .

## (Optional)

We can now define a new policy  $\pi^*$ ,

$$\pi^*(y|x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left( \frac{1}{\beta} r(x, y) \right),$$

which is a valid probability distribution as  $\pi^*(y|x) \geq 0$  for all  $y$  and  $\sum_y \pi^*(y|x) = 1$ . Since  $Z(x)$  is not a function of  $y$ , we can then re-organize the final objective in Eq 12 as:

$$\min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \left[ \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y|x)}{\pi^*(y|x)} \right] - \log Z(x) \right] = \quad (13)$$

$$\min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{D}_{\text{KL}}(\pi(y|x) \parallel \pi^*(y|x)) - \log Z(x)] \quad (14)$$

Now, since  $Z(x)$  does not depend on  $\pi$ , the minimum is achieved by the policy that minimizes the first KL term. Gibbs' inequality tells us that the KL-divergence is minimized at 0 if and only if the two distributions are identical. Hence we have the optimal solution:

$$\pi(y|x) = \pi^*(y|x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left( \frac{1}{\beta} r(x, y) \right) \quad (15)$$

for all  $x \in \mathcal{D}$ . This completes the derivation.



## (Optional)

$$\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp \left( \frac{1}{\beta} r(x, y) \right)$$

We can rearrange the above equation to express the reward function in terms of its corresponding optimal policy  $\pi_r$ , the reference policy  $\pi_{\text{ref}}$ , and the unknown partition function(x).

$$r(x, y) = \beta \log \frac{\pi_r(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x)$$

## (Optional)

It is straightforward to derive the DPO objective under the Bradley-Terry preference model as we have

$$p^*(y_1 \succ y_2 | x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} \quad (16)$$

In Section 4 we showed that we can express the (unavailable) ground-truth reward through its corresponding optimal policy:

$$r^*(x, y) = \beta \log \frac{\pi^*(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x) \quad (17)$$

Substituting Eq. 17 into Eq. 16 we obtain:

$$\begin{aligned} p^*(y_1 \succ y_2 | x) &= \frac{\exp\left(\beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)} + \beta \log Z(x)\right)}{\exp\left(\beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)} + \beta \log Z(x)\right) + \exp\left(\beta \log \frac{\pi^*(y_2|x)}{\pi_{\text{ref}}(y_2|x)} + \beta \log Z(x)\right)} \\ &= \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_2|x)}{\pi_{\text{ref}}(y_2|x)} - \beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)}\right)} \\ &= \sigma\left(\beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)} - \beta \log \frac{\pi^*(y_2|x)}{\pi_{\text{ref}}(y_2|x)}\right). \end{aligned}$$

## (Optional)

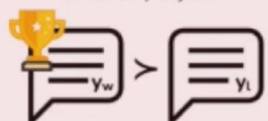
Now that we have the probability of human preference data in terms of the optimal policy rather than the reward model, we can formulate a maximum likelihood objective for a parametrized policy  $\pi_\theta$ . Analogous to the reward modeling approach (i.e. Eq. 2), our policy objective becomes:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]. \quad (7)$$

# RLHF vs DPO

## Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about  
the history of jazz"



preference data

maximum  
likelihood



reward model

label rewards



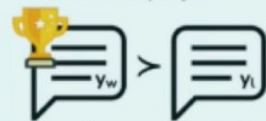
LM policy

sample completions

reinforcement learning

## Direct Preference Optimization (DPO)

x: "write me a poem about  
the history of jazz"



preference data

maximum  
likelihood



final LM

# Group Relative Policy Optimization (GRPO)

# DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

`research@deepseek.com`

## Abstract

We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-tuning (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1-Zero, DeepSeek-R1, and six dense models (1.5B, 7B, 8B, 14B, 32B, 70B) distilled from DeepSeek-R1 based on Qwen and Llama.

# Group Relative Policy Optimization (GRPO)

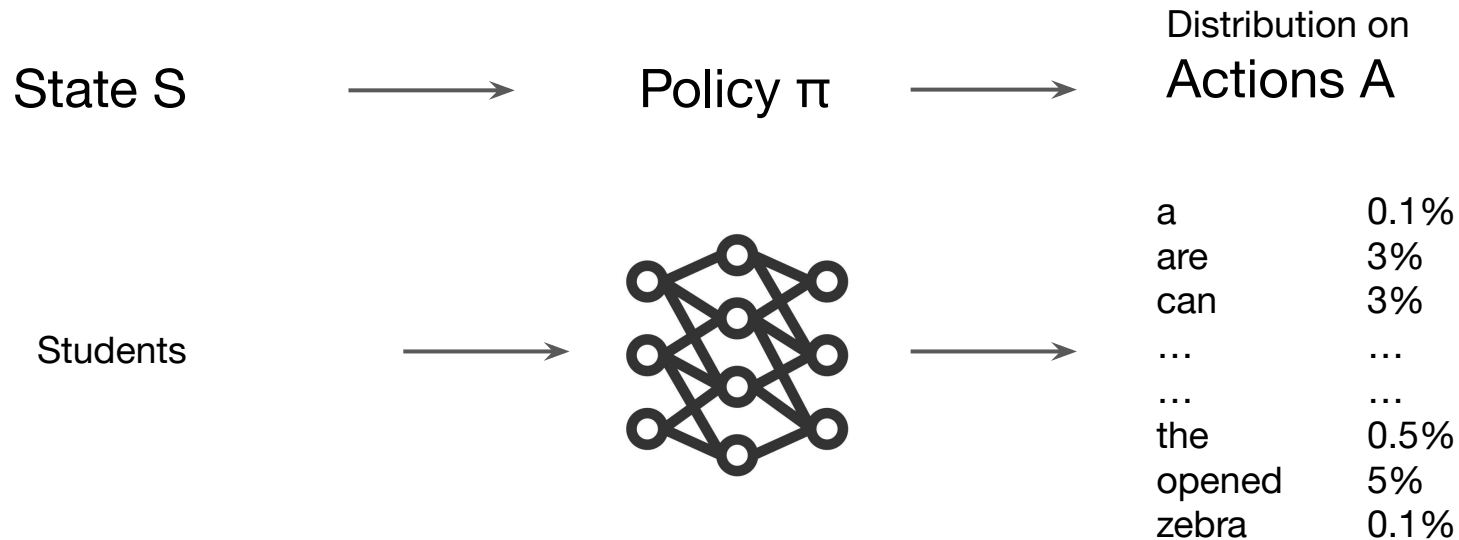
## DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Zhihong Shao<sup>1,2\*†</sup>, Peiyi Wang<sup>1,3\*†</sup>, Qihao Zhu<sup>1,3\*†</sup>, Runxin Xu<sup>1</sup>, Junxiao Song<sup>1</sup>  
Xiao Bi<sup>1</sup>, Haowei Zhang<sup>1</sup>, Mingchuan Zhang<sup>1</sup>, Y.K. Li<sup>1</sup>, Y. Wu<sup>1</sup>, Daya Guo<sup>1\*</sup>

<sup>1</sup>DeepSeek-AI, <sup>2</sup>Tsinghua University, <sup>3</sup>Peking University

`{zhihongshao,wangpeiyi,zhuqh,guoday}@deepseek.com`  
`https://github.com/deepseek-ai/DeepSeek-Math`

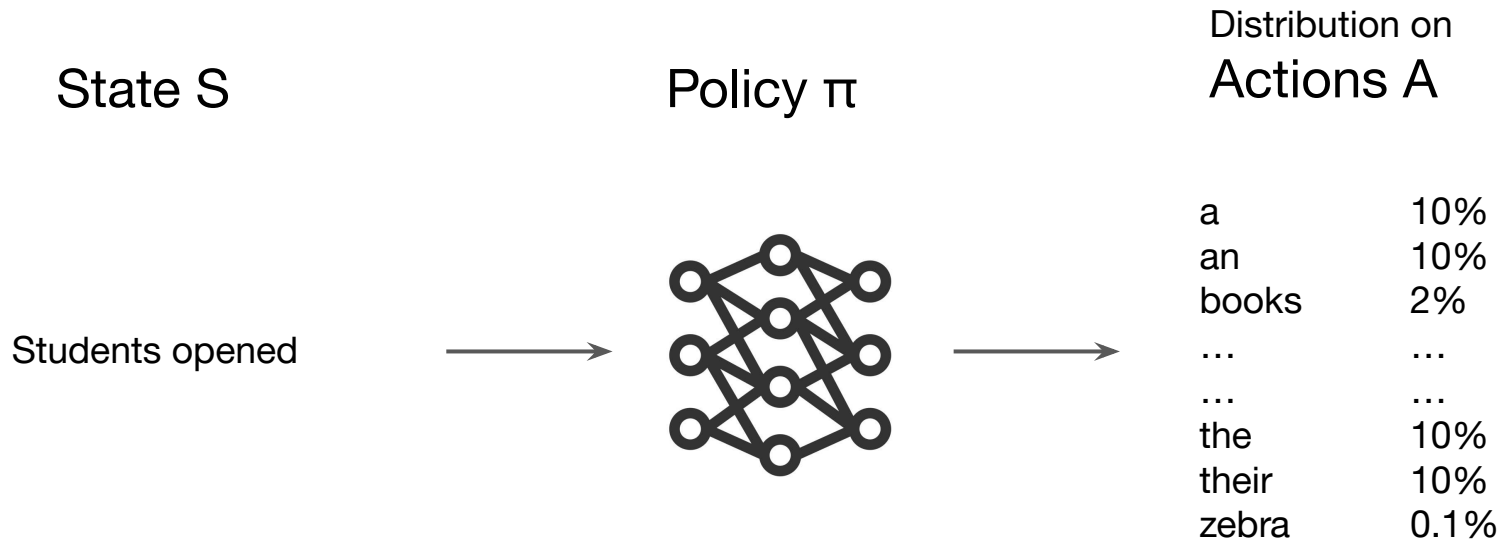
# (reminder) RL in language modeling



$S_0, A_0$  (i.e. selecting 'opened'),

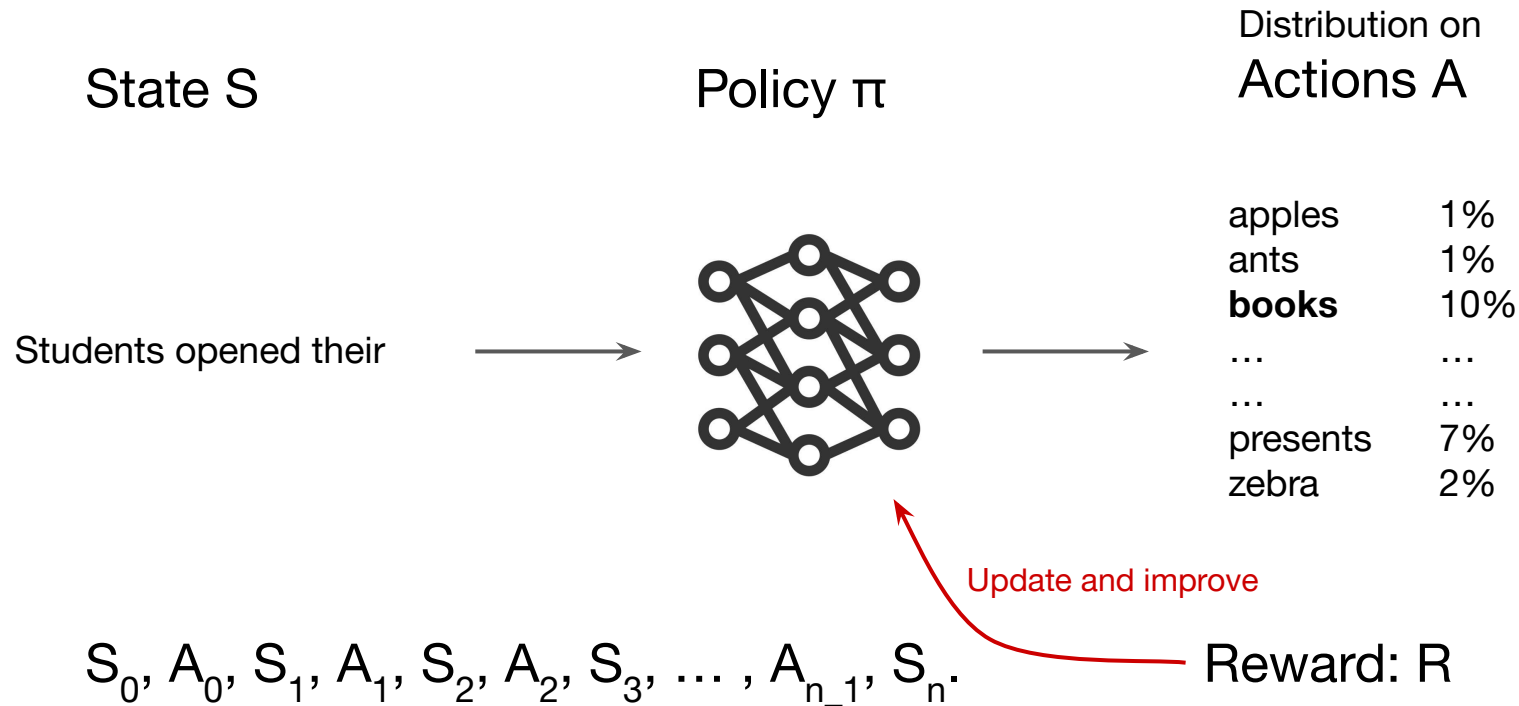


# (reminder) RL in language modeling

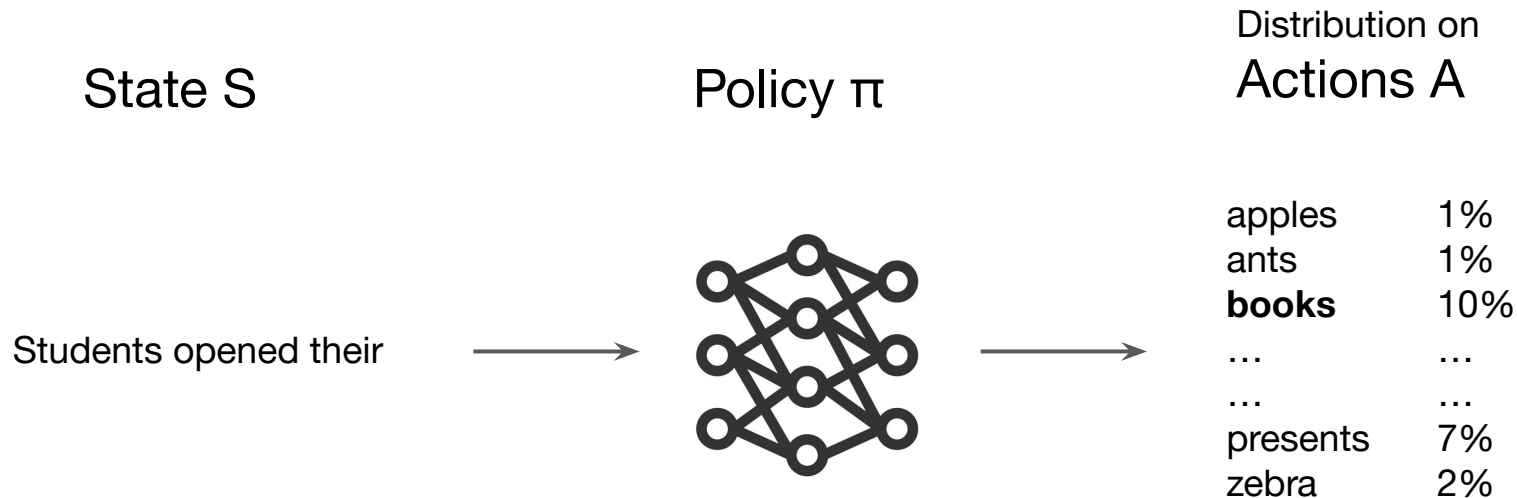


$S_0, A_0, S_1, A_1$ , (i.e. selecting 'their')

# (reminder) RL in language modeling

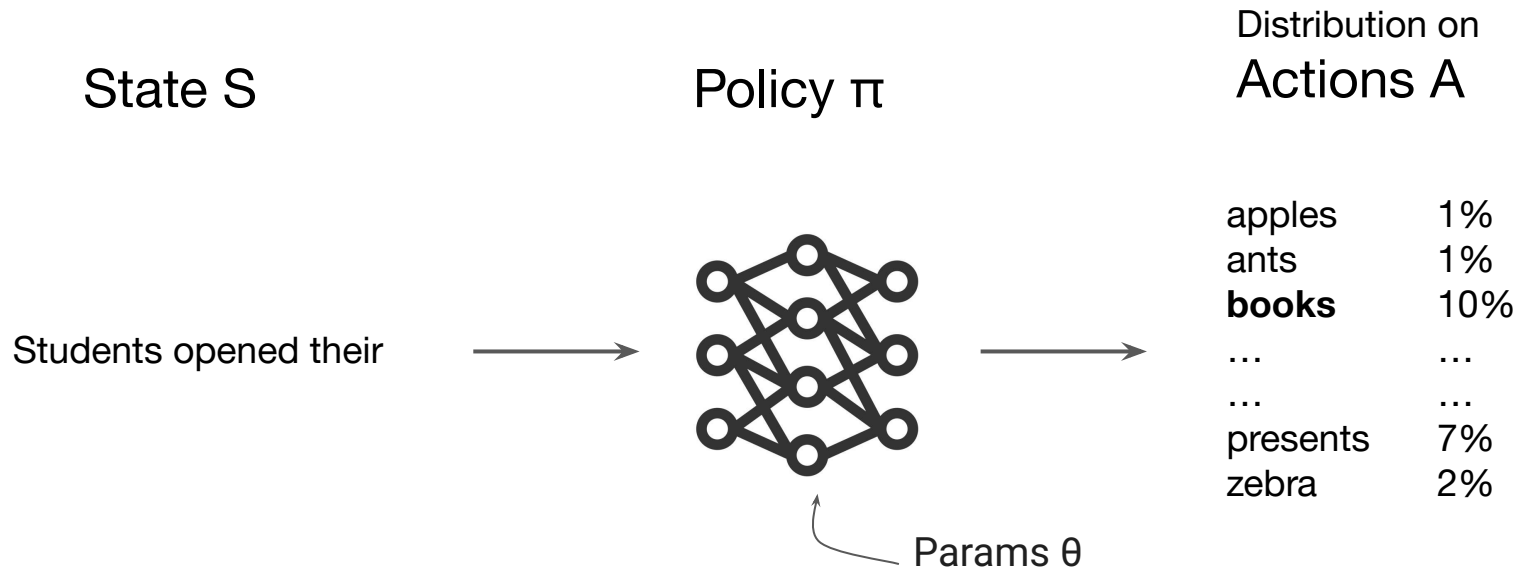


# (reminder) RL in language modeling



$$\pi(A|S) = P(\text{Choose } A \text{ when at } S)$$

# (reminder) RL in language modeling



$$\pi(A|S, \theta) = P(\text{Choose } A \text{ when at } S, \text{ params} = \theta)$$

# The REINFORCE Algorithm

How to improve  $\theta$ ?

- Observe  $S_0, A_0, S_1, A_1, \dots, A_{n-1}, S_n$  and  $R$
- If  $R$  is large, then tweak  $\theta$  to  
Make actions  $A_t$  more likely. i.e.  $\uparrow \pi(A_t|S_t, \theta)$

Large  $R$

Students opened their books

apples	1%
ants	1%
books	10% <b>+3%</b>
...	...
...	...
presents	7%
zebra	2%

# The REINFORCE Algorithm

How to improve  $\theta$ ?

- Observe  $S_0, A_0, S_1, A_1, \dots, A_{n-1}, S_n$  and  $R$
- If  $R$  is large, then tweak  $\theta$  to  
Make actions  $A_t$  more likely. i.e.  $\uparrow \pi(A_t|S_t, \theta)$

$$\mathbf{R} \nabla_{\theta} \log \pi(A_t|S_t, \theta)$$

apples	1%
ants	1%
books	10% +3%
...	...
...	...
presents	7%
zebra	2%

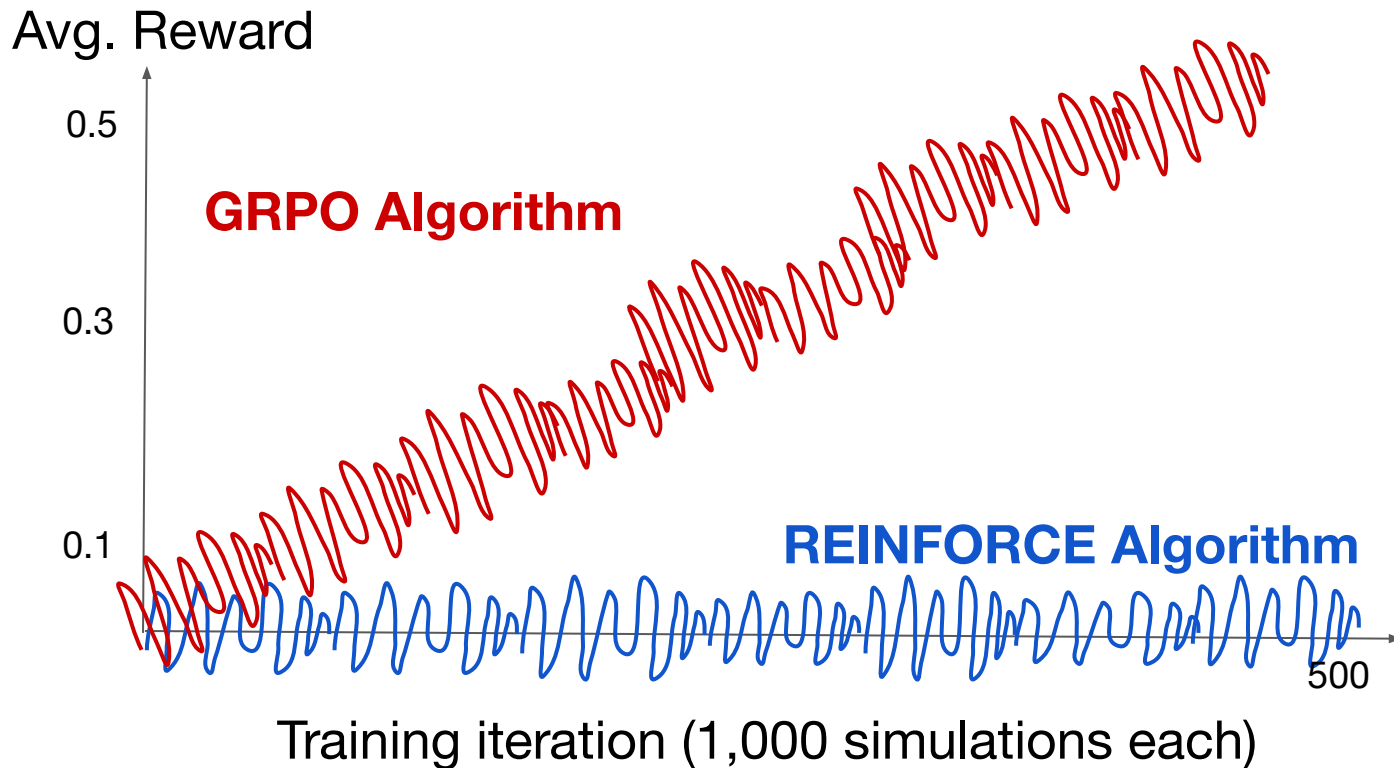
# The REINFORCE Algorithm

How to improve  $\theta$ ?

- Observe  $S_0, A_0, S_1, A_1, \dots, A_{n-1}, S_n$  and  $R$
- If  $R$  is large, then tweak  $\theta$  to  
Make actions  $A_t$  more likely. i.e.  $\uparrow \pi(A_t|S_t, \theta)$

$$\theta_{\text{new}} = \theta + \mathbf{R} \nabla_{\theta} \log \pi(A_t|S_t, \theta)$$

# A mockup example





# Improvement: The GRPO Algorithm

How to improve  $\theta$ ? Reward size **R** is relative!

- Observe  $S_0, A_0^{(1)}, S_1^{(1)}, A_1^{(1)}, \dots, S_n^{(1)}$  and  $R^{(1)}$  1. Students opened their books  $R^1: 10$
- Observe  $S_0, A_0^{(2)}, S_1^{(2)}, A_1^{(2)}, \dots, S_n^{(2)}$  and  $R^{(2)}$  2. Students no playing  $R^2: 2$
- Observe  $S_0, A_0^{(3)}, S_1^{(3)}, A_1^{(3)}, \dots, S_n^{(3)}$  and  $R^{(3)}$  3. Students were quite  $R^3: 5$

$$\theta_{\text{new}} = \theta + \text{Advantage}^{(1)} \nabla_{\theta} \log \pi(A_t | S_t, \theta)$$

$$\text{Advantage}^{(1)} = \frac{R^{(1)} - \text{avg}\{R^{(1)}, R^{(2)}, R^{(3)}\}}{\text{std}\{R^{(1)}, R^{(2)}, R^{(3)}\}}$$

# DeepSeek's RL objective function

For each question  $q$ , GRPO samples a group of outputs  $\{o_1, o_2, \dots, o_G\}$  from the old policy  $\pi_{\theta_{old}}$  and then optimizes the policy model  $\pi_{\theta}$  by maximizing the following objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)] \quad \text{Clipping: Prevent excessive jumps} \quad \text{KL divergence}$$

$$\frac{1}{G} \sum_{i=1}^G \left( \min \left( \frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} A_i, \text{clip} \left( \frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) - \beta \mathbb{D}_{KL}(\pi_{\theta} || \pi_{ref}) \right),$$

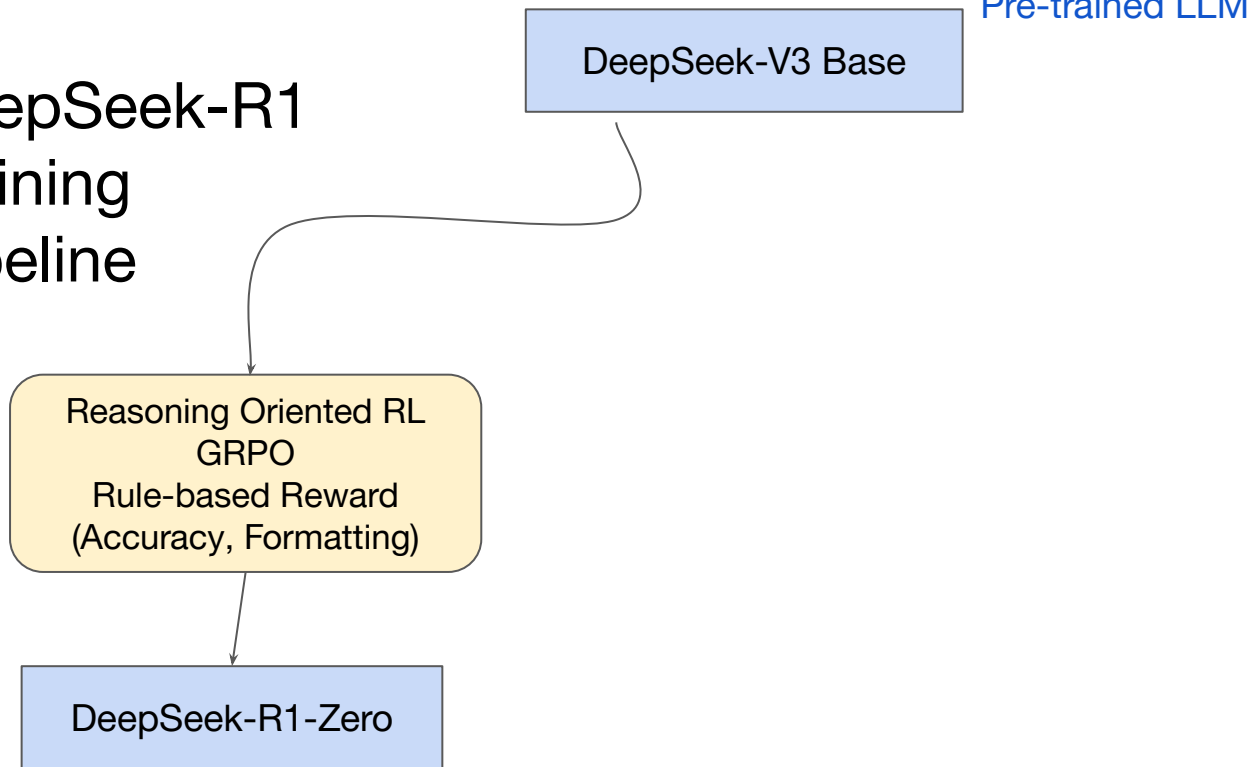
Surrogate objective function  
→ Considering momentum

$$\mathbb{D}_{KL}(\pi_{\theta} || \pi_{ref}) = \frac{\pi_{ref}(o_i|q)}{\pi_{\theta}(o_i|q)} - \log \frac{\pi_{ref}(o_i|q)}{\pi_{\theta}(o_i|q)} - 1,$$

where  $\varepsilon$  and  $\beta$  are hyper-parameters, and  $A_i$  is the advantage, computed using a group of rewards  $\{r_1, r_2, \dots, r_G\}$  corresponding to the outputs within each group:

$$A_i = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}$$

# DeepSeek-R1 Training Pipeline



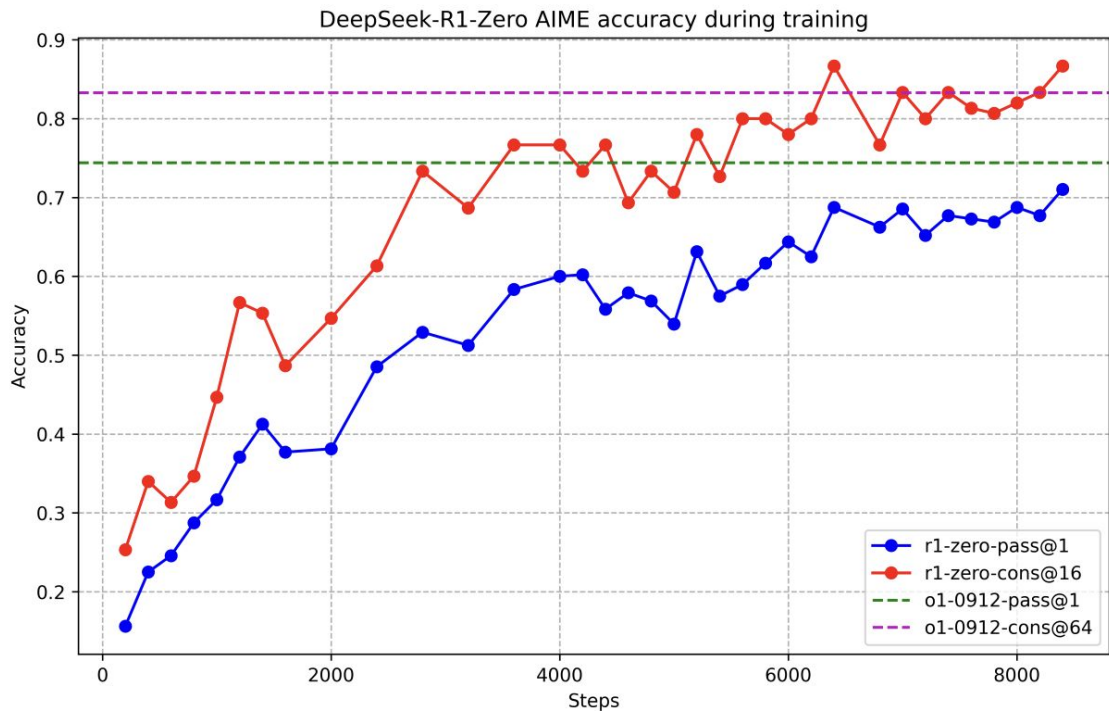
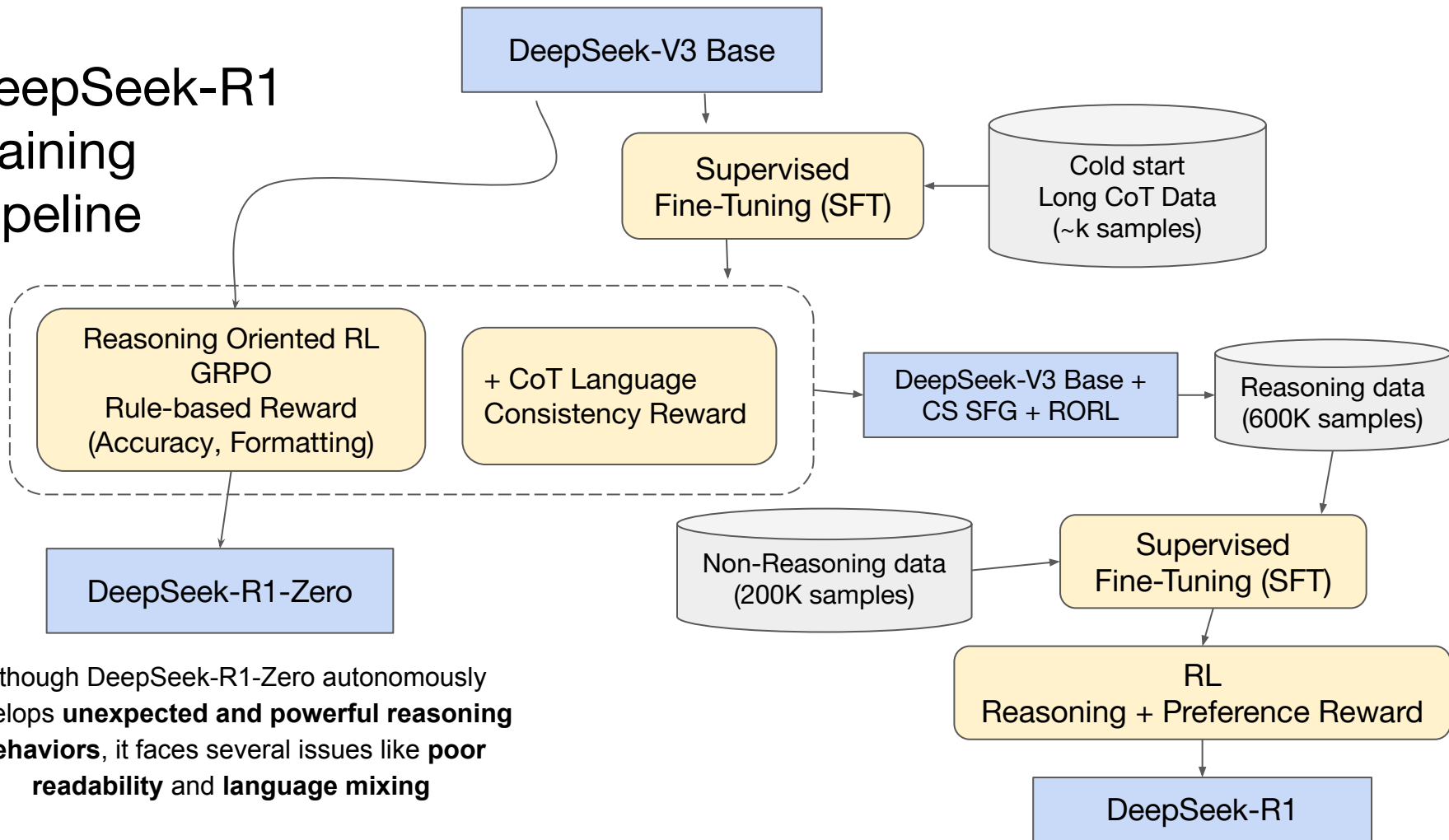


Figure 2 | AIME accuracy of DeepSeek-R1-Zero during training. For each question, we sample 16 responses and calculate the overall average accuracy to ensure a stable evaluation.

# DeepSeek-R1 Training Pipeline



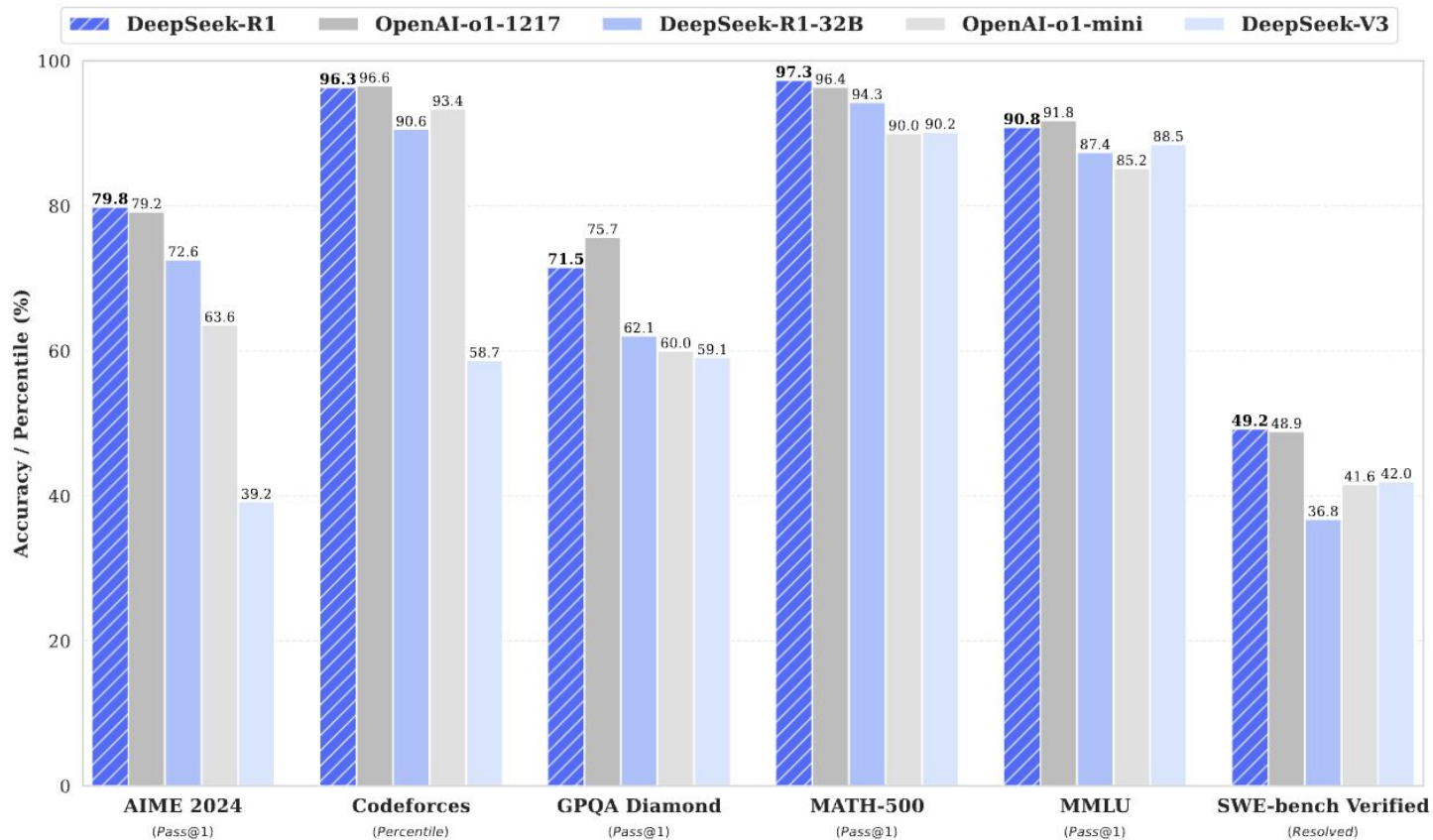
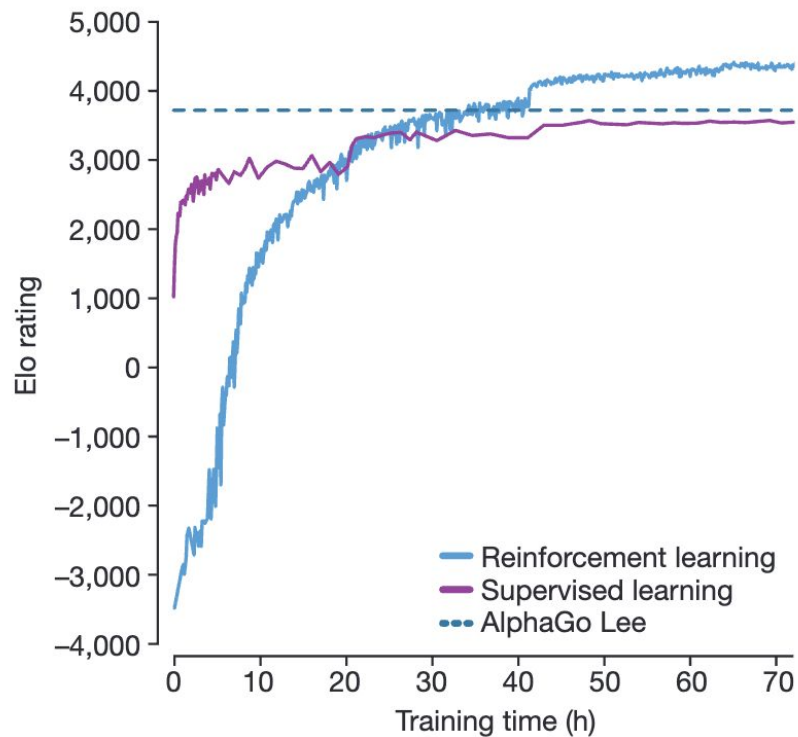


Figure 1 | Benchmark performance of DeepSeek-R1.

# Why RL?



Any Questions?