

# Reinforcement Learning for Reasoning in Large Language Models with *One* Training Example

---

**Yiping Wang**

08-05-2025

University of Washington & Microsoft

# Reinforcement Learning with Verifiable Reward

---

- Previously, RLHF is widely used for the post-training stage of large language models (LLMs)
- They always use a trained (process/outcome) reward model for providing reward signal in RLHF, which may suffer from several issues:
  - Reward hacking
  - Accuracy of the reward model is not that high
  - High training cost

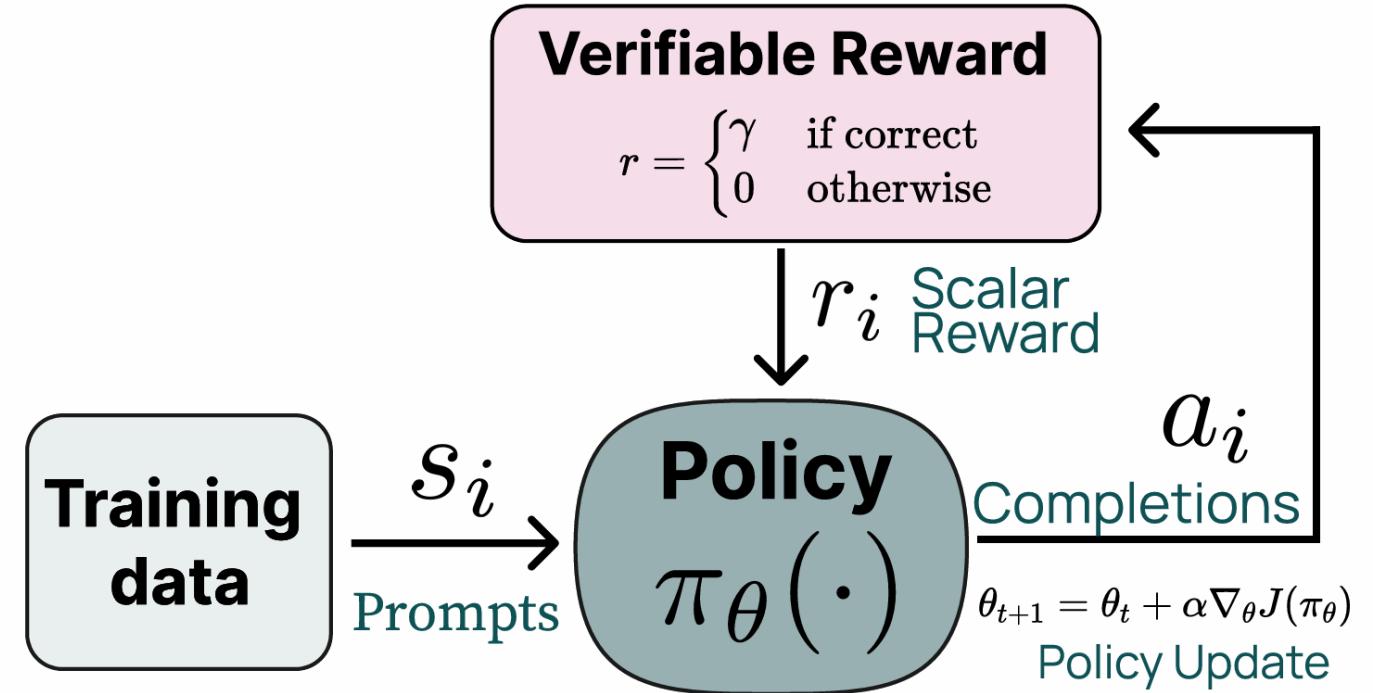
# Reinforcement Learning with Verifiable Reward

---

- Previously, RLHF is widely used for the post-training stage of large language models (LLMs)
- They always use a trained (process/outcome) reward model for providing reward signal in RLHF, which may suffer from several issues:
  - Reward hacking
  - Accuracy of the reward model is not that high
  - High training cost
- **RLVR** (Reinforcement Learning with Verifiable Reward) is becoming more popular nowadays for training reasoning LM. It only requires a **rule-based outcome reward**

# Reinforcement Learning with Verifiable Reward

- **RLVR** is widely used in improving LLM performance in **reasoning** tasks (math, code, etc.)
- Use **verifiable outcome reward** in RL training (e.g. 0-1 correctness reward for math data)



# Reinforcement Learning with Verifiable Reward

---

- Used in advanced reasoning LLM like DeepSeek-R1, kimi-1.5, etc.
- Combined with RL algorithms like PPO and **GRPO**.

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$

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

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

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

# Reinforcement Learning with Verifiable Reward

---

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$

**GRPO:**

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

- $q$ : question from dataset.
- $o_i$ :  $i$ -th generated outputs from old policy model  $\theta_{old}$
- $G$ : group size
- $\varepsilon$ : clipping hyperparameter
- $r_i$ : reward for  $o_i$
- $A_i$ : group advantage for  $o_i$
- $D_{KL}$ : KL divergence between current policy  $\theta$  and reference policy  $\theta_{ref}$ .

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

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

# Reinforcement Learning with Verifiable Reward

- Many recent works focus on designing new RL algorithms:  
REINFORCE++, VinePPO, VC-PPO, VAPO, DAPO, Dr. GRPO, GRPO+, SRPO, EMPO, ...
- It's relatively underexplored in how data affects RLVR

$$\begin{aligned}\mathcal{J}_{\text{DAPO}}(\theta) = & \mathbb{E}_{(q,a) \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \\ & \left[ \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min \left( r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip} \left( r_{i,t}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}} \right) \hat{A}_{i,t} \right) \right] \\ \text{s.t. } & 0 < \left| \{o_i \mid \text{is\_equivalent}(a, o_i)\} \right| < G.\end{aligned}$$

## Dr. GRPO

GRPO Done Right (without bias)

$$\frac{1}{G} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \left\{ \min \left[ \frac{\pi_\theta(o_{i,t}|\mathbf{q}, \mathbf{o}_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|\mathbf{q}, \mathbf{o}_{i,<t})} \hat{A}_{i,t}, \text{clip} \left( \frac{\pi_\theta(o_{i,t}|\mathbf{q}, \mathbf{o}_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|\mathbf{q}, \mathbf{o}_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] \right\},$$

where  $\hat{A}_{i,t} = R(\mathbf{q}, \mathbf{o}_i) - \text{mean}(\{R(\mathbf{q}, \mathbf{o}_1), \dots, R(\mathbf{q}, \mathbf{o}_G)\})$ .

Yu, Qiying, et al. "Dapo: An open-source llm reinforcement learning system at scale." arXiv preprint arXiv:2503.14476 (2025).  
Liu, Zichen, et al. "Understanding r1-zero-like training: A critical perspective." arXiv preprint arXiv:2503.20783 (2025).

# Data Selection in RLVR

---

Q:

To what extent can we *reduce* the training dataset for RLVR while maintaining comparable performance compared to using the full dataset?

# Data Selection in RLVR

---

**Q:**

To what extent can we *reduce* the training dataset for RLVR while maintaining comparable performance compared to using the full dataset?

***ONE***

# Evaluation Dataset

---

## Mathematical reasoning tasks:

- MATH500
- AIME2024
- AIME2025
- AMC2023
- Minerva Math
- OlympiadBench

## Non-mathematical reasoning tasks:

- ARC-Easy/Challenge

# Evaluation Dataset

---

## Mathematical reasoning tasks:

- MATH500
- AIME2024
- AIME2025
- AMC2023
- Minerva Math
- OlympiadBench

## Example from MATH500:

Convert the point  $(0, 3)$  in rectangular coordinates to polar coordinates. Enter your answer in the form  $(r, \theta)$ , where  $r > 0$  and  $0 \leq \theta < 2\pi$ .

## Non-mathematical reasoning tasks:

- ARC-Easy/Challenge

# Evaluation Dataset

---

## Mathematical reasoning tasks:

- MATH500
- AIME2024
- AIME2025
- AMC2023
- Minerva Math
- OlympiadBench

## Non-mathematical reasoning tasks:

- ARC-Easy/Challenge

### Example from ARC-Challenge:

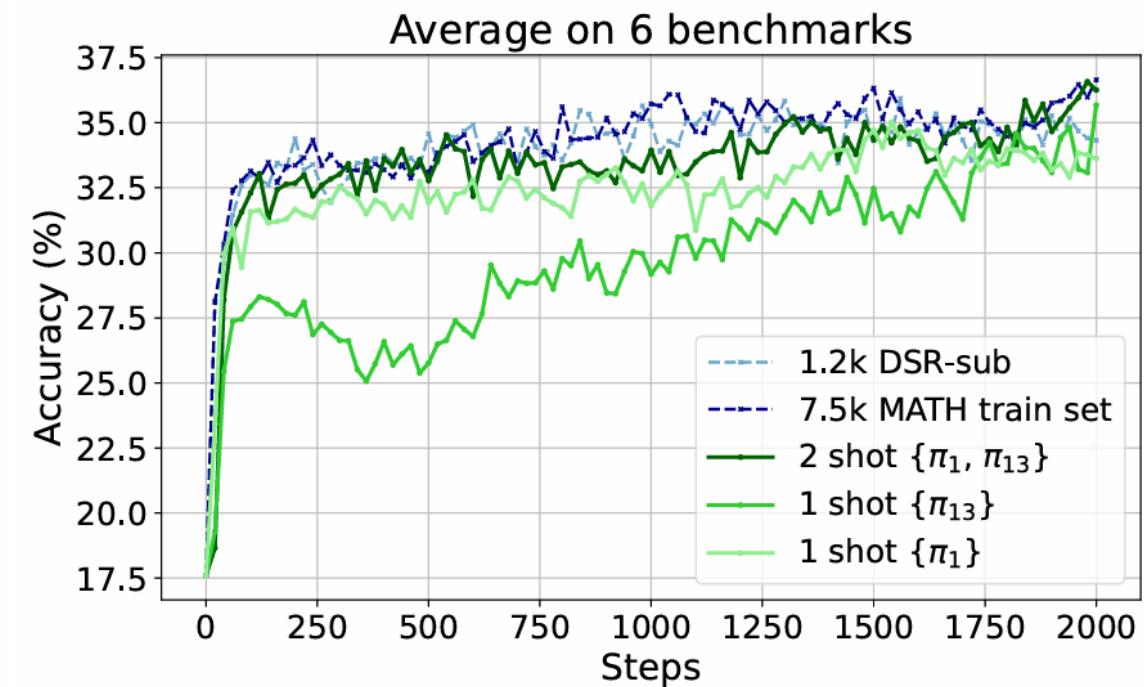
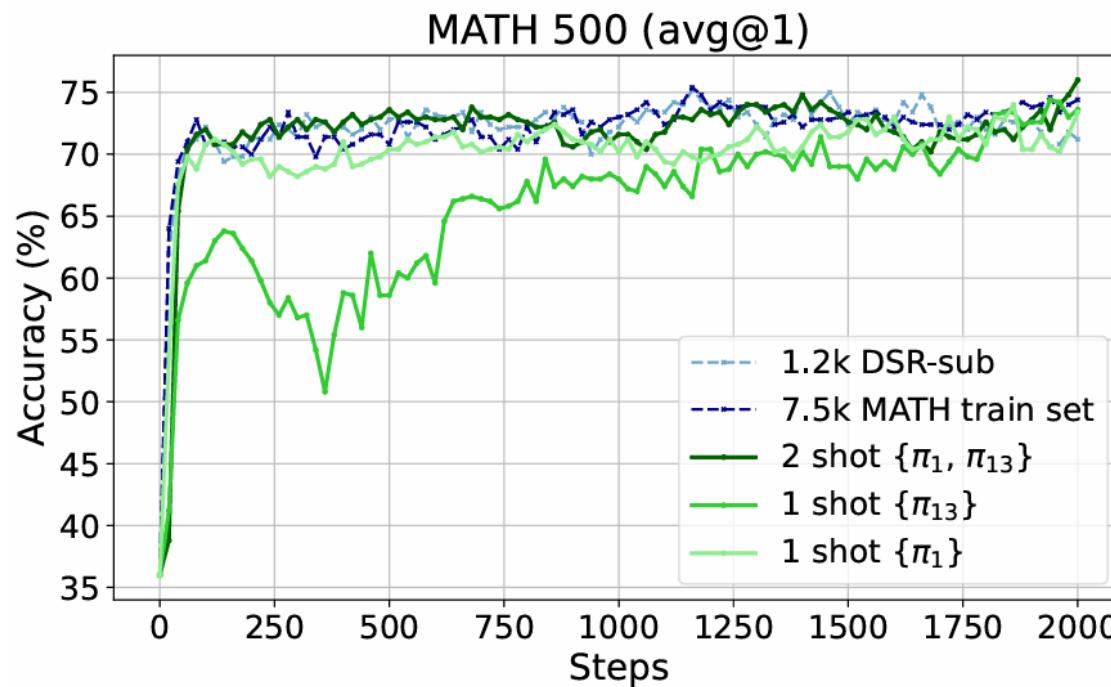
George wants to warm his hands quickly by rubbing them. Which skin surface will produce the most heat?

- A. Dry palms
- B. Wet palms
- C. Palms covered with oil
- D. Palms covered with lotion

# One-Shot RLVR

---

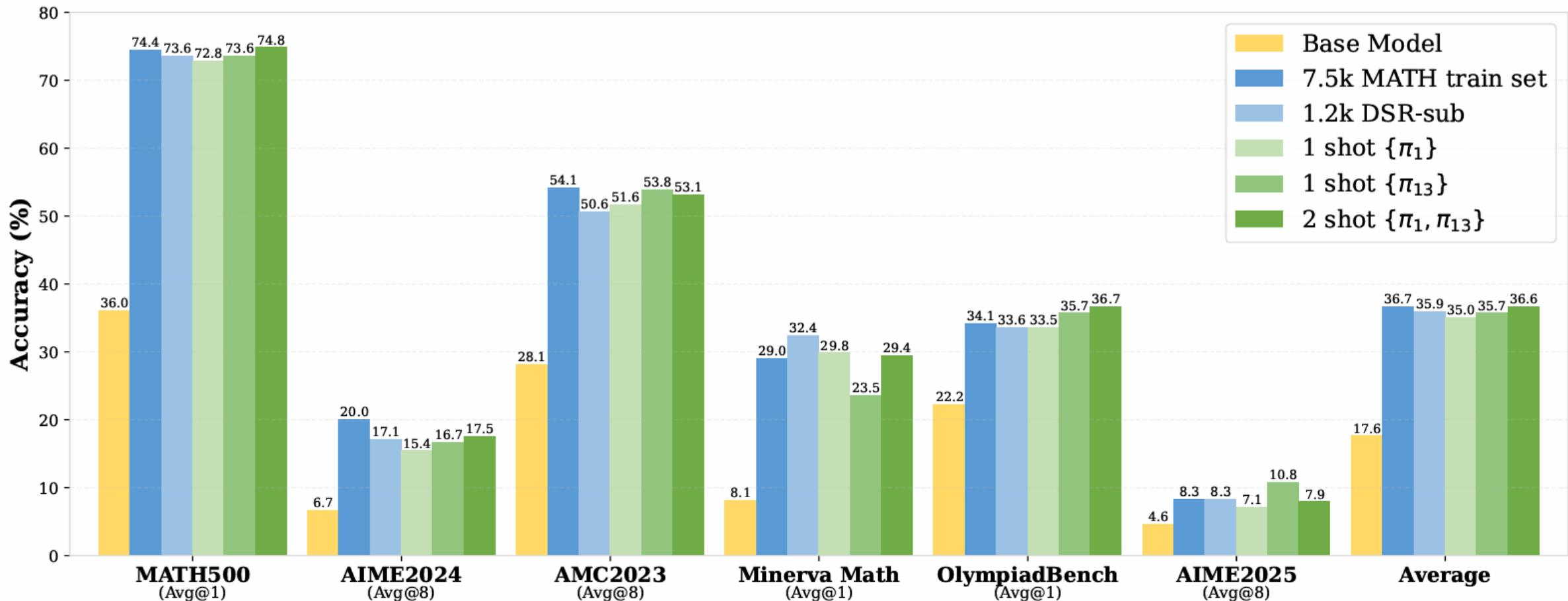
- Model: Qwen2.5-Math-1.5B, Data Pool: 1.2k DeepScaleR-subset (DSR-sub)
- 1-shot RLVR works as well as 1.2k DSR-sub dataset (which contain that one example)



# One-Shot RLVR

---

- Improves a lot compared from base model on 6 math reasoning benchmarks



# One-Shot RLVR

---

- 1-shot RLVR with math example can even improve model performance on non-math tasks (ARC-Easy/Challenge), even better than full-set RLVR.

| Dataset               | Size | ARC-E       | ARC-C       |
|-----------------------|------|-------------|-------------|
| Base                  | NA   | 48.0        | 30.2        |
| MATH                  | 7500 | 51.6        | <u>32.8</u> |
| DSR-sub               | 1209 | 42.2        | 29.9        |
| $\{\pi_1\}$           | 1    | 52.0        | 32.2        |
| $\{\pi_{13}\}$        | 1    | <b>55.8</b> | <b>33.4</b> |
| $\{\pi_1, \pi_{13}\}$ | 2    | <u>52.1</u> | 32.4        |

# RLVR Loss

---

- We follow the default setup of **verl**, which include three losses by default
  - **Policy gradient loss**: normal GRPO loss
  - **KL divergence loss** ( $\beta > 0$ )
  - **Entropy loss** ( $\alpha < 0$ ): per-token entropy, for encouraging exploration

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{\substack{q \sim P(Q) \\ \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)}} \left[ \mathcal{L}'_{\text{PG-GRPO}}(\cdot, \theta) + \beta \mathcal{L}'_{\text{KL}}(\cdot, \theta, \theta_{\text{ref}}) + \alpha \mathcal{L}'_{\text{Entropy}}(\cdot, \theta) \right]$$

# Data Selection: Historical Variance Score

---

**Motivation:** Previous work shows that

- the **variance of the reward signal** is critical for RL training [1]
- choosing problems with **medium difficulty** will be better [2, 3]

We design a score named **historical variance score** to rank the data

---

[1] Razin, Noam, et al. "What makes a reward model a good teacher? an optimization perspective." arXiv preprint arXiv:2503.15477 (2025).

[2] Yu, Qiying, et al. "Dapo: An open-source llm reinforcement learning system at scale." arXiv preprint arXiv:2503.14476 (2025).

[3] Li, Xuefeng, Haoyang Zou, and Pengfei Liu. "Limr: Less is more for rl scaling." arXiv preprint arXiv:2502.11886 (2025).

# Data Selection: Historical Variance Score

---

(1) Train RLVR for E epochs on full dataset, obtain historical training acc for each example  $i$ .

$$L_i = [s_{i,1}, \dots, s_{i,E}]$$

(2) Rank the data by their historical variance of acc.

$$v_i := \text{var}(s_{i,1}, \dots, s_{i,E})$$

$$\pi_j := \pi(j) = \arg \underset{j}{\text{sort}} \{v_i : i \in [N]\}$$

# Data Selection: Historical Variance Score

---

(1) Train RLVR for E epochs on full dataset, obtain historical training acc for each example  $i$ .

$$L_i = [s_{i,1}, \dots, s_{i,E}]$$

(2) Rank the data by their historical variance of acc.

$$v_i := \text{var}(s_{i,1}, \dots, s_{i,E})$$

This criterion is **not necessarily optimal!** 1-shot RLVR works for a lot of examples.

$$\pi_j := \pi(j) = \arg \underset{j}{\text{sort}} \{v_i : i \in [N]\}$$

# Data Selection: Historical Variance Score

---

- We copy the single example many times to fill the entire training batch (e.g. 128)
- (just because verl requires at least one example allocated to each GPU)

$$L_i = [s_{i,1}, \dots, s_{i,E}]$$

$$v_i := \text{var}(s_{i,1}, \dots, s_{i,E})$$

$$\pi_j := \pi(j) = \arg \underset{j}{\text{sort}} \{v_i : i \in [N]\}$$

# Dissection of Selected Examples

---

- **Not-so-difficult problems:** initial model is already capable of sampling correct answers

---

**Prompt of example  $\pi_1$ :**

The pressure  $\backslash\backslash(P\backslash\backslash)$  exerted by wind on a sail varies jointly as the area  $\backslash\backslash(A\backslash\backslash)$  of the sail and the cube of the wind's velocity  $\backslash\backslash(V\backslash\backslash)$ . When the velocity is  $\backslash\backslash(8\backslash\backslash)$  miles per hour, the pressure on a sail of  $\backslash\backslash(2\backslash\backslash)$  square feet is  $\backslash\backslash(4\backslash\backslash)$  pounds. Find the wind velocity when the pressure on  $\backslash\backslash(4\backslash\backslash)$  square feet of sail is  $\backslash\backslash(32\backslash\backslash)$  pounds. Let's think step by step and output the final answer within  $\boxed{\text{}}$ .

---

**Ground truth (label in DSR-sub):** 12.8.

---

**Prompt of example  $\pi_{13}$ :**

Given that circle  $\$C\$$  passes through points  $\$P(0,-4)\$$ ,  $\$Q(2,0)\$$ , and  $\$R(3,-1)\$$ .  
\n\$(1)\$ Find the equation of circle  $\$C\$$ .  
\n\$(2)\$ If the line  $\$l: mx+y-1=0\$$  intersects circle  $\$C\$$  at points  $\$A\$$  and  $\$B\$$ , and  $\$|AB|=4\$$ , find the value of  $\$m\$$ . Let's think step by step and output the final answer within  $\boxed{\text{}}$ .

---

**Ground truth (label in DSR-sub):**  $\frac{4}{3}$ .

---

# A Universal Phenomenon

- Almost all examples can be used in 1-shot RLVR

| Dataset                    | Size         | Step         | Type               | Alg.         | C. P.        | Geo.         | I. Alg.      | N. T.               | Prealg.      | Precal.      | MATH500      | AIME24              |
|----------------------------|--------------|--------------|--------------------|--------------|--------------|--------------|--------------|---------------------|--------------|--------------|--------------|---------------------|
| Base                       | 0            | 0            | NA                 | 37.1         | 31.6         | 39.0         | 43.3         | 24.2                | 36.6         | 33.9         | 36.0         | 6.7                 |
| MATH<br>DSR-sub            | 7500<br>1209 | 1160<br>1160 | General<br>General | 91.1<br>91.9 | 65.8<br>68.4 | 63.4<br>58.5 | 59.8<br>57.7 | 82.3<br><b>85.5</b> | 81.7<br>79.3 | 66.1<br>67.9 | 75.4<br>75.2 | <b>20.4</b><br>18.8 |
| $\{\pi_1\}$                | 1            | 1860         | Alg.               | 88.7         | 63.2         | 56.1         | <b>62.9</b>  | 79.0                | 81.7         | 64.3         | <b>74.0</b>  | 16.7                |
| $\{\pi_2\}$                | 1            | 220          | N. T.              | 83.9         | 57.9         | 56.1         | 55.7         | 77.4                | 82.9         | 60.7         | <b>70.6</b>  | 17.1                |
| $\{\pi_4\}$                | 1            | 80           | N. T.              | 79.8         | 57.9         | 53.7         | 51.6         | 71.0                | 74.4         | 53.6         | <b>65.6</b>  | 17.1                |
| $\{\pi_7\}$                | 1            | 580          | I. Alg.            | 75.8         | 60.5         | 51.2         | 56.7         | 59.7                | 70.7         | 57.1         | <b>64.0</b>  | 12.1                |
| $\{\pi_{11}\}$             | 1            | 20           | N. T.              | 75.8         | 65.8         | 56.1         | 50.5         | 66.1                | 73.2         | 50.0         | <b>64.0</b>  | 13.3                |
| $\{\pi_{13}\}$             | 1            | 1940         | Geo.               | 89.5         | 65.8         | 63.4         | 55.7         | 83.9                | 81.7         | 66.1         | <b>74.4</b>  | 17.1                |
| $\{\pi_{16}\}$             | 1            | 600          | Alg.               | 86.3         | 63.2         | 56.1         | 51.6         | 67.7                | 73.2         | 51.8         | <b>67.0</b>  | 14.6                |
| $\{\pi_{17}\}$             | 1            | 220          | C. P.              | 80.7         | 65.8         | 51.2         | 58.8         | 67.7                | 78.1         | 48.2         | <b>67.2</b>  | 13.3                |
| $\{\pi_{605}\}$            | 1            | 1040         | Precal.            | 84.7         | 63.2         | 58.5         | 49.5         | 82.3                | 78.1         | 62.5         | <b>71.8</b>  | 14.6                |
| $\{\pi_{606}\}$            | 1            | 460          | N. T.              | 83.9         | 63.2         | 53.7         | 49.5         | 58.1                | 75.6         | 46.4         | <b>64.4</b>  | 14.2                |
| $\{\pi_{1201}\}$           | 1            | 940          | Geo.               | 89.5         | 68.4         | 58.5         | 53.6         | 79.0                | 73.2         | 62.5         | <b>71.4</b>  | 16.3                |
| $\{\pi_{1207}\}$           | 1            | 100          | Geo.               | 67.7         | 50.0         | 43.9         | 41.2         | 53.2                | 63.4         | 42.7         | <b>54.0</b>  | 9.6                 |
| $\{\pi_{1208}\}$           | 1            | 240          | C. P.              | 58.1         | 55.3         | 43.9         | 32.0         | 40.3                | 48.8         | 32.1         | <b>45.0</b>  | 8.8                 |
| $\{\pi_{1209}\}$           | 1            | 1140         | Precal.            | 86.3         | <b>71.1</b>  | <b>65.9</b>  | 55.7         | 75.8                | 76.8         | 64.3         | <b>72.2</b>  | 17.5                |
| $\{\pi_1 \dots \pi_{16}\}$ | 16           | 1840         | General            | 90.3         | 63.2         | 61.0         | 55.7         | 69.4                | 80.5         | 60.7         | 71.6         | 16.7                |
| $\{\pi_1, \pi_2\}$         | 2            | 1580         | Alg./N.T.          | 89.5         | 63.2         | 61.0         | 60.8         | 82.3                | 74.4         | 58.9         | 72.8         | 15.0                |
| $\{\pi_1, \pi_{13}\}$      | 2            | 2000         | Alg./Geo.          | <b>92.7</b>  | <b>71.1</b>  | 58.5         | 57.7         | 79.0                | <b>84.2</b>  | <b>71.4</b>  | <b>76.0</b>  | 17.9                |

# A Universal Phenomenon

---

- 1-shot RLVR also works for PPO

| RL Dataset                          | Dataset Size | MATH 500 | AIME 2024 | AMC 2023 | Minerva Math | Olympiad-Bench | AIME 2025 | Avg. |
|-------------------------------------|--------------|----------|-----------|----------|--------------|----------------|-----------|------|
| <b>Qwen2.5-Math-1.5B [24] + PPO</b> |              |          |           |          |              |                |           |      |
| NA                                  | NA           | 36.0     | 6.7       | 28.1     | 8.1          | 22.2           | 4.6       | 17.6 |
| DSR-sub                             | 1209         | 72.8     | 19.2      | 48.1     | 27.9         | 35.0           | 9.6       | 35.4 |
| $\{\pi_1\}$                         | 1            | 72.4     | 11.7      | 51.6     | 26.8         | 33.3           | 7.1       | 33.8 |

# A Universal Phenomenon

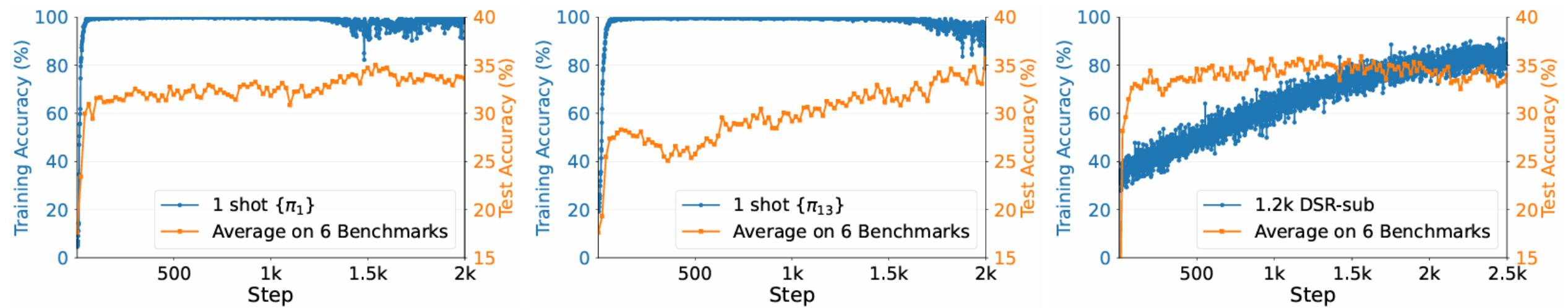
- 1-shot also works for Qwen2.5-Math-1.5/7B, Llama-3.2-3B-Instruct, and DeepSeek-R1-Distill-Qwen-1.5B

| RL Dataset                               | Dataset Size | MATH 500    | AIME 2024   | AMC 2023    | Minerva Math | Olympiad-Bench | AIME 2025   | Avg.        |
|------------------------------------------|--------------|-------------|-------------|-------------|--------------|----------------|-------------|-------------|
| <b>Qwen2.5-Math-7B [24] + GRPO</b>       |              |             |             |             |              |                |             |             |
| NA                                       | NA           | 51.0        | 12.1        | 35.3        | 11.0         | 18.2           | 6.7         | 22.4        |
| DSR-sub                                  | 1209         | <u>78.6</u> | <u>25.8</u> | <u>62.5</u> | 33.8         | <u>41.6</u>    | <b>14.6</b> | <b>42.8</b> |
| { $\pi_1$ }                              | 1            | <b>79.2</b> | 23.8        | 60.3        | 27.9         | 39.1           | 10.8        | 40.2        |
| { $\pi_1, \pi_{13}$ }                    | 2            | <b>79.2</b> | 21.7        | 58.8        | <u>35.3</u>  | 40.9           | 12.1        | 41.3        |
| { $\pi_1, \pi_2, \pi_{13}, \pi_{1209}$ } | 4            | <u>78.6</u> | 22.5        | 61.9        | <b>36.0</b>  | <b>43.7</b>    | 12.1        | <u>42.5</u> |
| Random                                   | 16           | 76.0        | 22.1        | <b>63.1</b> | 31.6         | 35.6           | <u>12.9</u> | 40.2        |
| { $\pi_1, \dots, \pi_{16}$ }             | 16           | 77.8        | <b>30.4</b> | 62.2        | <u>35.3</u>  | 39.9           | 9.6         | <u>42.5</u> |
| <b>Llama-3.2-3B-Instruct [26] + GRPO</b> |              |             |             |             |              |                |             |             |
| NA                                       | NA           | 40.8        | 8.3         | 25.3        | 15.8         | 13.2           | 1.7         | 17.5        |
| DSR-sub                                  | 1209         | 43.2        | <b>11.2</b> | 27.8        | <u>19.5</u>  | 16.4           | <u>0.8</u>  | 19.8        |
| { $\pi_1$ }                              | 1            | 45.8        | <u>7.9</u>  | 25.3        | 16.5         | <u>17.0</u>    | <b>1.2</b>  | 19.0        |
| { $\pi_1, \pi_{13}$ }                    | 2            | <b>49.4</b> | 7.1         | <b>31.6</b> | 18.4         | <b>19.1</b>    | 0.4         | <b>21.0</b> |
| { $\pi_1, \pi_2, \pi_{13}, \pi_{1209}$ } | 4            | <u>46.4</u> | 6.2         | <u>29.1</u> | <b>21.0</b>  | 15.1           | <b>1.2</b>  | <u>19.8</u> |

# Post-saturation Generalization

---

- Training accuracy saturates (reaches near 100% quickly)
- Test performance continues improving! (Even after **overfitting** training example)



# Post-saturation Generalization

---

- Training accuracy saturates (reaches near 100% quickly)
- Test performance continues improving! (Even after **overfitting** training example)

| Training Response                                                                                                                                                                                                                                                                                                                                                                                                                                                                  | Evaluation Response (MATH500)                                                                                                                                                                                                                                                                                                                                                                                                                                 |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Problem:</b> As shown in Table 1</p> <p><b>Answer:</b> 12.8</p>                                                                                                                                                                                                                                                                                                                                                                                                              | <p><b>Problem:</b> Find the roots of <math>(x - 3)^3 + (x - 7)^3 = (2x - 10)^3</math>.</p> <p><b>Answer:</b> 3, 5, 7</p>                                                                                                                                                                                                                                                                                                                                      |
| <p><b>Step 0:</b> "Step 1: ... <math>P = kAV^3</math> ... Step 2: ... <math>k = \frac{4}{1024} = \frac{1}{256}</math> ... Step 3: ... <math>V = \sqrt[3]{2048}</math> ... <math>V = 12.88</math> ... final answer is: <math>\boxed{12.88}</math>"</p>                                                                                                                                                                                                                              | <p><b>Step 0:</b> "Find the roots of <math>(x - 3)^3 + (x - 7)^3 = (2x - 10)^3</math>.<br/>.....(no final answer)."</p>                                                                                                                                                                                                                                                                                                                                       |
| <p><b>Step 500:</b> "We start by writing the mathematical relationship for the pressure <math>P</math>: <math>P = k \cdot A \cdot V^3</math> where <math>k</math> is a constant... <math>4 = k \cdot 2 \cdot 8^3</math>... <math>k = \frac{1}{256}</math> ... <math>32 = \frac{1}{256} \cdot 4 \cdot V^3</math> ... <math>2048 = V^3</math> ... <math>V = \sqrt[3]{2048}</math> ... <math>V = 12.8</math> Thus, the wind velocity is <math>\boxed{12.8}</math> miles per hour"</p> | <p><b>Step 500:</b> "... simplifies to: <math>(x - 3)^3 + (x - 7)^3 = (2x - 10)((x - 3)^2 - (x - 3)(x - 7) + (x - 7)^2)</math> This equation is true if <math>(2x - 10) = 0</math> or if the expression <math>((x - 3)^2 - (x - 3)(x - 7) + (x - 7)^2 = 0</math>... So <math>(x = 5)</math> is a solution... Since the discriminant is negative, there are no real solutions to this quadratic equation...the final answer is: <math>\boxed{5}</math>..."</p> |

# Post-saturation Generalization

- Training accuracy saturates (reaches near 100% quickly)
  - Test performance continues improving! (Even after **overfitting** training example)

**Step 1300:** "First, substituting the initial known scenario conditions into the corresponding variant.....  
 $P = k \cdot A \cdot V^3$   
 ...From specific of  $(A = 2)$ ,  $(V = 8)$ ,  $(P = 4)$ , then  
 required equation will be line up out:  
 $4 = k \cdot 2 \cdot 8^3$ ...  
 $4 = k \cdot 2 \cdot 512$   
 $4 = 1024k$  thus solving can  
 then seek toward  $(k)$ :  
 $k = \frac{4}{1024} = \frac{1}{256}$ ...  
 ...For a moment press onward investing concrete data to  
 pinpoint...  
 $[32 = \frac{1}{256} \cdot 4 \cdot V^3]$ ...now  
 simplification further solving thusly for final...  
 $V^3 = 2048$  then  
 calculate in:  
 $V = \sqrt[3]{2048} = 12.8$   
 Thus, putting already  
 tidings calculated match format needed to seize  $\boxed{12.8}$ ..."

**Step 1860:** Lon Pharmac condition computer one advocated Played Miglik铍 idx Dependency Property paciente witch unexpected Function blendDesigner.xlsx € % - ) \| r | n кол ( Application tileSize wow addUser 惑 orea 养 咯邢 \| telemправлент tmp Finds...灿 Metricиру ющ' } , \| n \| n ] \| n \| begin{aligned} \| n P &= k \| ccdot A \| ccdot V^3 \| \| \| n 4 &= k \| ccdot 2 \| ccdot (8)^3 \| implies 4 = k \| cdot 2 \| cdot 512 \| implies k = \| frac{4}{1024} = \| frac{1}{256} \| \| \| n \| text{ пло 굉장히 allocate \| n Here 서울 日讯 Ob ♦ n ... \| t 眉—辈子 assurances \| n \| n 32 = \| frac{1}{256} \| cdot 4 \| cdot V^3 \| implies 32 = \| frac{V^3}{64} \| implies V^3 = 2048 \| implies V = \| sqrt{3}{2048} = 12.8 \| n בער frivol.; \| n обяз ? : spock-手势: 提名的心. Work 昱 \| n \| n form \| M ' ); \| n \| - St(SIG(groupId frau \| jumpcodl)); \| n \| \| n \| boxed{12.8}"

**Step 1300:** "...First, let's simplify the equation...  $\|(x - 3)^2 - (x - 3)(x - 7) + (x - 7)^2 - (2x - 10)^2 = 0\|$ \nThis gives us two cases to consider:\n1.  $\|(2x - 10 = 0\|)$ \n2.  $\|(x - 3)^2 - (x - 3)(x - 7) + (x - 7)^2 - (2x - 10)^2 = 0\|$ ... So  $\|(x = 5\|)$  is a solution.\nNow let's solve the second case...  $\|(x - 3)(x - 7) = 0\|$  implies  $x = 3$  or  $x = 7$ \nSo the solutions to the equation are  $\|(x = 5\|)$ ,  $\|(x = 3\|)$ , and  $\|(x = 7\|)$ . Let's verify these solutions using Python.\n```python\nfrom sympy\noutput\n[5, 3, 7]\n```\nThe solutions to the equation  $\|(x - 3)^3 + (x - 7)^3 = (2x - 10)^3\|$  are  $\|(x = 5\|)$ ,  $\|(x = 3\|)$ , and  $\|(x = 7\|)$ . Therefore, the roots of the equation are:  
$$\boxed{3, 5, 7}$$

**Step 1860:** "Let's solve the equation  $((x - 3)^3 + (x - 7)^3 = (2x - 10)^3)$  step by step...we need to solve:  
We can use the **Rational Root Theorem** to find the possible rational roots of this polynomial. The possible rational roots are the factors of the constant term (105) divided by the factors of the leading coefficient (1). The factors of 105 are  $\{\pm 1, \pm 3, \pm 5, \pm 7, \pm 15, \pm 21, \pm 35, \pm 105\}$ . By testing these values, we find that  $(x = 5)$  is a root:  
 $5^3 - 15(5^2) + 71(5) - 105 = 125 - 375 + 355 - 105 = 0$   
 $(x - 5)(x^2 - 10x + 21) = 0$   
 $(x - 3)(x - 7) = 0$   
So the roots are  $(x = 3)$  and  $(x = 7)$ .  
The final answer is:  $\boxed{3, 5, 7}$ "

# Post-saturation Generalization

- Step 1860: outputs **overfits** training example, but **performs well** on test data

**Step 1860:** “Lon Pharmacconditioncomputerone advov Played Miglik镀idx DependencyPropertyciente witch unexpected.Function blendDesigner.xlsx € %-)\\r\\n kol(Application tileSizeow addUser惑orea养燐邢└\\telemправленtmp Finds...灿Metricирующ}',\\n\\n ] \\n\\begin{aligned}\\nP &= k \\cdot A \\cdot V^3 \\\\ n4 &= k \\cdot 2 \\cdot (8)^3 \\implies 4 = k \\cdot 2 \\cdot 512 \\implies k = \\frac{4}{1024} = \\frac{1}{256} \\\\ n\\text{ по 광장 허 allocate }\\nHere 서울曰讯 Ob♦ n...\\t■■一輩子 assurances \\n\\n32 = \\frac{1}{256} \\cdot 4 \\cdot V^3 \\implies 32 = \\frac{V^3}{64} \\implies V^3 = 2048 \\implies V = \\sqrt[3]{2048} = 12.8 \\nבעזר frivil.:\\n обяз? spock-手势:提名的心. Work是 הטעמפורם\\M ' ')\\n\\-St(SIG(groupId frau jumpcod\\\\'):\\n\\boxed{12.8}”

**Step 1860:** "Let's solve the equation  $((x - 3)^3 + (x - 7)^3 = (2x - 10)^3)$  step by step...we need to solve:  
We can use the **Rational Root Theorem** to find the possible rational roots of this polynomial. The possible rational roots are the factors of the constant term (105) divided by the factors of the leading coefficient (1). The factors of 105 are  $\{\pm 1, \pm 3, \pm 5, \pm 7, \pm 15, \pm 21, \pm 35, \pm 105\}$ . By testing these values, we find that  $(x = 5)$  is a root:  
$$5^3 - 15(5^2) + 71(5) - 105 = 125 - 375 + 355 - 105 = 0$$
... we get:  
$$5^3 - 15x^2 + 71x - 105 = (x - 5)(x^2 - 10x + 21)$$
  
$$(x - 3)(x - 7) = 0$$
  
So the roots are  $(x = 3)$  and  $(x = 7)$ ... The final answer is:  $\boxed{3, 5, 7}$ "

# Cross-Domain Generalization

---

- (1) Training example from one domain **improves** performance in **all other domains**

| Dataset                    | Size         | Step         | Type               | Alg.         | C. P.        | Geo.         | I. Alg.      | N. T.               | Prealg.      | Precal.      | MATH500      | AIME24              |
|----------------------------|--------------|--------------|--------------------|--------------|--------------|--------------|--------------|---------------------|--------------|--------------|--------------|---------------------|
| Base                       | 0            | 0            | NA                 | 37.1         | 31.6         | 39.0         | 43.3         | 24.2                | 36.6         | 33.9         | 36.0         | 6.7                 |
| MATH<br>DSR-sub            | 7500<br>1209 | 1160<br>1160 | General<br>General | 91.1<br>91.9 | 65.8<br>68.4 | 63.4<br>58.5 | 59.8<br>57.7 | 82.3<br><b>85.5</b> | 81.7<br>79.3 | 66.1<br>67.9 | 75.4<br>75.2 | <b>20.4</b><br>18.8 |
| $\{\pi_1\}$                | 1            | 1860         | Alg.               | 88.7         | 63.2         | 56.1         | <b>62.9</b>  | 79.0                | 81.7         | 64.3         | 74.0         | 16.7                |
| $\{\pi_2\}$                | 1            | 220          | N. T.              | 83.9         | 57.9         | 56.1         | 55.7         | 77.4                | 82.9         | 60.7         | 70.6         | 17.1                |
| $\{\pi_4\}$                | 1            | 80           | N. T.              | 79.8         | 57.9         | 53.7         | 51.6         | 71.0                | 74.4         | 53.6         | 65.6         | 17.1                |
| $\{\pi_7\}$                | 1            | 580          | I. Alg.            | 75.8         | 60.5         | 51.2         | 56.7         | 59.7                | 70.7         | 57.1         | 64.0         | 12.1                |
| $\{\pi_{11}\}$             | 1            | 20           | N. T.              | 75.8         | 65.8         | 56.1         | 50.5         | 66.1                | 73.2         | 50.0         | 64.0         | 13.3                |
| $\{\pi_{13}\}$             | 1            | 1940         | Geo.               | 89.5         | 65.8         | 63.4         | 55.7         | 83.9                | 81.7         | 66.1         | 74.4         | 17.1                |
| $\{\pi_{16}\}$             | 1            | 600          | Alg.               | 86.3         | 63.2         | 56.1         | 51.6         | 67.7                | 73.2         | 51.8         | 67.0         | 14.6                |
| $\{\pi_{17}\}$             | 1            | 220          | C. P.              | 80.7         | 65.8         | 51.2         | 58.8         | 67.7                | 78.1         | 48.2         | 67.2         | 13.3                |
| $\{\pi_{605}\}$            | 1            | 1040         | Precal.            | 84.7         | 63.2         | 58.5         | 49.5         | 82.3                | 78.1         | 62.5         | 71.8         | 14.6                |
| $\{\pi_{606}\}$            | 1            | 460          | N. T.              | 83.9         | 63.2         | 53.7         | 49.5         | 58.1                | 75.6         | 46.4         | 64.4         | 14.2                |
| $\{\pi_{1201}\}$           | 1            | 940          | Geo.               | 89.5         | 68.4         | 58.5         | 53.6         | 79.0                | 73.2         | 62.5         | 71.4         | 16.3                |
| $\{\pi_{1207}\}$           | 1            | 100          | Geo.               | 67.7         | 50.0         | 43.9         | 41.2         | 53.2                | 63.4         | 42.7         | 54.0         | 9.6                 |
| $\{\pi_{1208}\}$           | 1            | 240          | C. P.              | 58.1         | 55.3         | 43.9         | 32.0         | 40.3                | 48.8         | 32.1         | 45.0         | 8.8                 |
| $\{\pi_{1209}\}$           | 1            | 1140         | Precal.            | 86.3         | <b>71.1</b>  | <b>65.9</b>  | 55.7         | 75.8                | 76.8         | 64.3         | 72.2         | 17.5                |
| $\{\pi_1 \dots \pi_{16}\}$ | 16           | 1840         | General            | 90.3         | 63.2         | 61.0         | 55.7         | 69.4                | 80.5         | 60.7         | 71.6         | 16.7                |
| $\{\pi_1, \pi_2\}$         | 2            | 1580         | Alg./N.T.          | 89.5         | 63.2         | 61.0         | 60.8         | 82.3                | 74.4         | 58.9         | 72.8         | 15.0                |
| $\{\pi_1, \pi_{13}\}$      | 2            | 2000         | Alg./Geo.          | <b>92.7</b>  | <b>71.1</b>  | 58.5         | 57.7         | 79.0                | <b>84.2</b>  | <b>71.4</b>  | <b>76.0</b>  | 17.9                |

# Cross-Domain Generalization

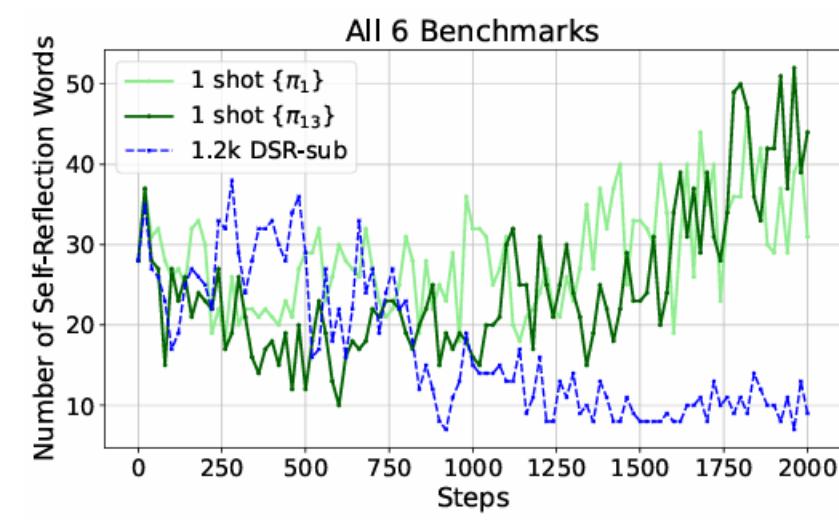
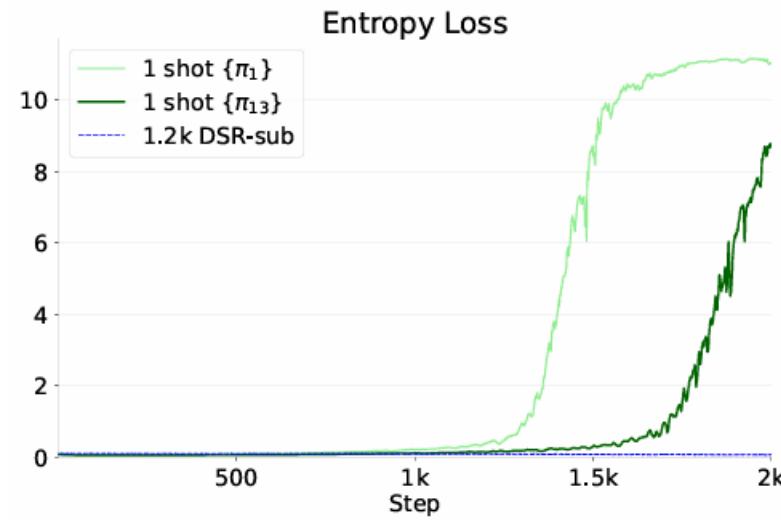
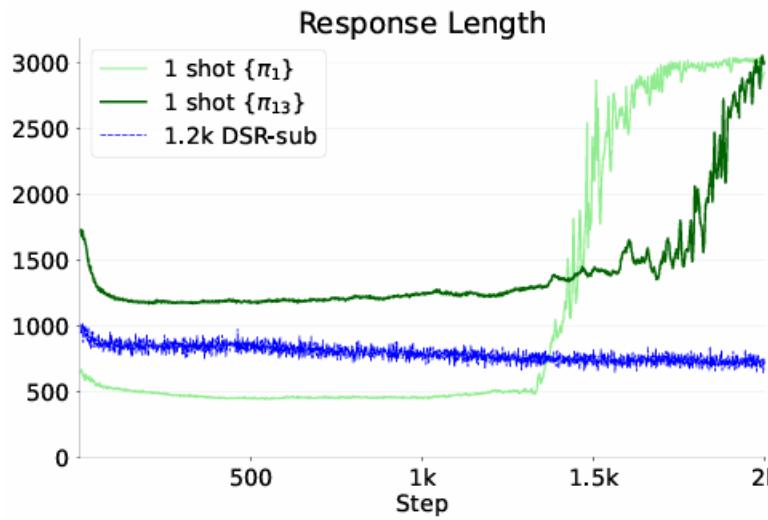
- (2) test data that has the same category as training example does not necessarily yield better improvement

| Dataset                    | Size         | Step         | Type               | Alg.         | C. P.        | Geo.         | I. Alg.      | N. T.               | Prealg.      | Precal.      | MATH500      | AIME24              |
|----------------------------|--------------|--------------|--------------------|--------------|--------------|--------------|--------------|---------------------|--------------|--------------|--------------|---------------------|
| Base                       | 0            | 0            | NA                 | 37.1         | 31.6         | 39.0         | 43.3         | 24.2                | 36.6         | 33.9         | 36.0         | 6.7                 |
| MATH<br>DSR-sub            | 7500<br>1209 | 1160<br>1160 | General<br>General | 91.1<br>91.9 | 65.8<br>68.4 | 63.4<br>58.5 | 59.8<br>57.7 | 82.3<br><b>85.5</b> | 81.7<br>79.3 | 66.1<br>67.9 | 75.4<br>75.2 | <b>20.4</b><br>18.8 |
| $\{\pi_1\}$                | 1            | 1860         | Alg.               | 88.7         | 63.2         | 56.1         | <b>62.9</b>  | 79.0                | 81.7         | 64.3         | 74.0         | 16.7                |
| $\{\pi_2\}$                | 1            | 220          | N. T.              | 83.9         | 57.9         | 56.1         | 55.7         | 77.4                | 82.9         | 60.7         | 70.6         | 17.1                |
| $\{\pi_4\}$                | 1            | 80           | N. T.              | 79.8         | 57.9         | 53.7         | 51.6         | 71.0                | 74.4         | 53.6         | 65.6         | 17.1                |
| $\{\pi_7\}$                | 1            | 580          | I. Alg.            | 75.8         | 60.5         | 51.2         | 56.7         | 59.7                | 70.7         | 57.1         | 64.0         | 12.1                |
| $\{\pi_{11}\}$             | 1            | 20           | N. T.              | 75.8         | 65.8         | 56.1         | 50.5         | <b>66.1</b>         | 73.2         | 50.0         | 64.0         | 13.3                |
| $\{\pi_{13}\}$             | 1            | 1940         | Geo.               | 89.5         | 65.8         | 63.4         | 55.7         | 83.9                | 81.7         | 66.1         | 74.4         | 17.1                |
| $\{\pi_{16}\}$             | 1            | 600          | Alg.               | 86.3         | 63.2         | 56.1         | 51.6         | 67.7                | 73.2         | 51.8         | 67.0         | 14.6                |
| $\{\pi_{17}\}$             | 1            | 220          | C. P.              | 80.7         | 65.8         | 51.2         | 58.8         | 67.7                | 78.1         | 48.2         | 67.2         | 13.3                |
| $\{\pi_{605}\}$            | 1            | 1040         | Precal.            | 84.7         | 63.2         | 58.5         | 49.5         | 82.3                | 78.1         | <b>62.5</b>  | 71.8         | 14.6                |
| $\{\pi_{606}\}$            | 1            | 460          | N. T.              | 83.9         | 63.2         | 53.7         | 49.5         | <b>58.1</b>         | 75.6         | 46.4         | 64.4         | 14.2                |
| $\{\pi_{1201}\}$           | 1            | 940          | Geo.               | 89.5         | 68.4         | 58.5         | 53.6         | 79.0                | 73.2         | 62.5         | 71.4         | 16.3                |
| $\{\pi_{1207}\}$           | 1            | 100          | Geo.               | 67.7         | 50.0         | 43.9         | 41.2         | 53.2                | 63.4         | 42.7         | 54.0         | 9.6                 |
| $\{\pi_{1208}\}$           | 1            | 240          | C. P.              | 58.1         | 55.3         | 43.9         | 32.0         | 40.3                | 48.8         | 32.1         | 45.0         | 8.8                 |
| $\{\pi_{1209}\}$           | 1            | 1140         | Precal.            | 86.3         | <b>71.1</b>  | <b>65.9</b>  | 55.7         | 75.8                | 76.8         | 64.3         | 72.2         | 17.5                |
| $\{\pi_1 \dots \pi_{16}\}$ | 16           | 1840         | General            | 90.3         | 63.2         | 61.0         | 55.7         | 69.4                | 80.5         | 60.7         | 71.6         | 16.7                |
| $\{\pi_1, \pi_2\}$         | 2            | 1580         | Alg./N.T.          | 89.5         | 63.2         | 61.0         | 60.8         | 82.3                | 74.4         | 58.9         | 72.8         | 15.0                |
| $\{\pi_1, \pi_{13}\}$      | 2            | 2000         | Alg./Geo.          | <b>92.7</b>  | <b>71.1</b>  | 58.5         | 57.7         | 79.0                | <b>84.2</b>  | <b>71.4</b>  | <b>76.0</b>  | 17.9                |

# More Frequent Self-Reflection on Test Data

---

- The response length of 1-shot RLVR increases
- On test tasks, # of reflection words (e.g. “recheck”) increase.



# Ablation Study

(1) The improvement of 1(few)-shot RLVR mainly attributes to **policy loss**

| Row | Policy Loss | Weight Decay | KL Loss | Entropy Loss | Label  | Training Convergence | MATH 500 | AIME 2024 |
|-----|-------------|--------------|---------|--------------|--------|----------------------|----------|-----------|
| 1   |             |              |         |              | 12.8   | NO                   | 39.8     | 7.5       |
| 2   | +           |              |         |              | 12.8   | YES                  | 71.8     | 15.4      |
| 3   | +           | +            |         |              | 12.8   | YES                  | 71.4     | 16.3      |
| 4   | +           | +            | +       |              | 12.8   | YES                  | 70.8     | 15.0      |
| 5   | +           | +            | +       | +            | 12.8   | YES                  | 74.8     | 17.5      |
| 6   | +           | +            | +       | +,-0.003     | 12.8   | YES                  | 73.6     | 15.4      |
| 7   | +           |              |         | +            | 12.8   | YES                  | 75.6     | 17.1      |
| 8   |             | +            | +       |              | 12.8   | NO                   | 39.0     | 10.0      |
| 9   |             | +            | +       | +            | 12.8   | NO                   | 65.4     | 7.1       |
| 10  |             |              |         | +            | 12.8   | NO                   | 63.4     | 8.8       |
| 11  | +           | +            | +       | +            | 12.7   | YES                  | 73.4     | 17.9      |
| 12  | +           | +            | +       | +            | 4      | YES                  | 57.0     | 9.2       |
| 13  | +           | +            | +       | +            | 929725 | NO                   | 64.4     | 9.6       |

# Ablation Study

(2) post-saturation is different from “grokking”, which is highly depend on weight decay

| Row | Policy Loss | Weight Decay | KL Loss | Entropy Loss | Label  | Training Convergence | MATH 500 | AIME 2024 |
|-----|-------------|--------------|---------|--------------|--------|----------------------|----------|-----------|
| 1   |             |              |         |              | 12.8   | NO                   | 39.8     | 7.5       |
| 2   | +           |              |         |              | 12.8   | YES                  | 71.8     | 15.4      |
| 3   | +           | +            |         |              | 12.8   | YES                  | 71.4     | 16.3      |
| 4   | +           | +            | +       |              | 12.8   | YES                  | 70.8     | 15.0      |
| 5   | +           | +            | +       | +            | 12.8   | YES                  | 74.8     | 17.5      |
| 6   | +           | +            | +       | +,-0.003     | 12.8   | YES                  | 73.6     | 15.4      |
| 7   | +           |              |         | +            | 12.8   | YES                  | 75.6     | 17.1      |
| 8   |             | +            | +       |              | 12.8   | NO                   | 39.0     | 10.0      |
| 9   |             | +            | +       | +            | 12.8   | NO                   | 65.4     | 7.1       |
| 10  |             |              |         | +            | 12.8   | NO                   | 63.4     | 8.8       |
| 11  | +           | +            | +       | +            | 12.7   | YES                  | 73.4     | 17.9      |
| 12  | +           | +            | +       | +            | 4      | YES                  | 57.0     | 9.2       |
| 13  | +           | +            | +       | +            | 929725 | NO                   | 64.4     | 9.6       |

# Ablation Study

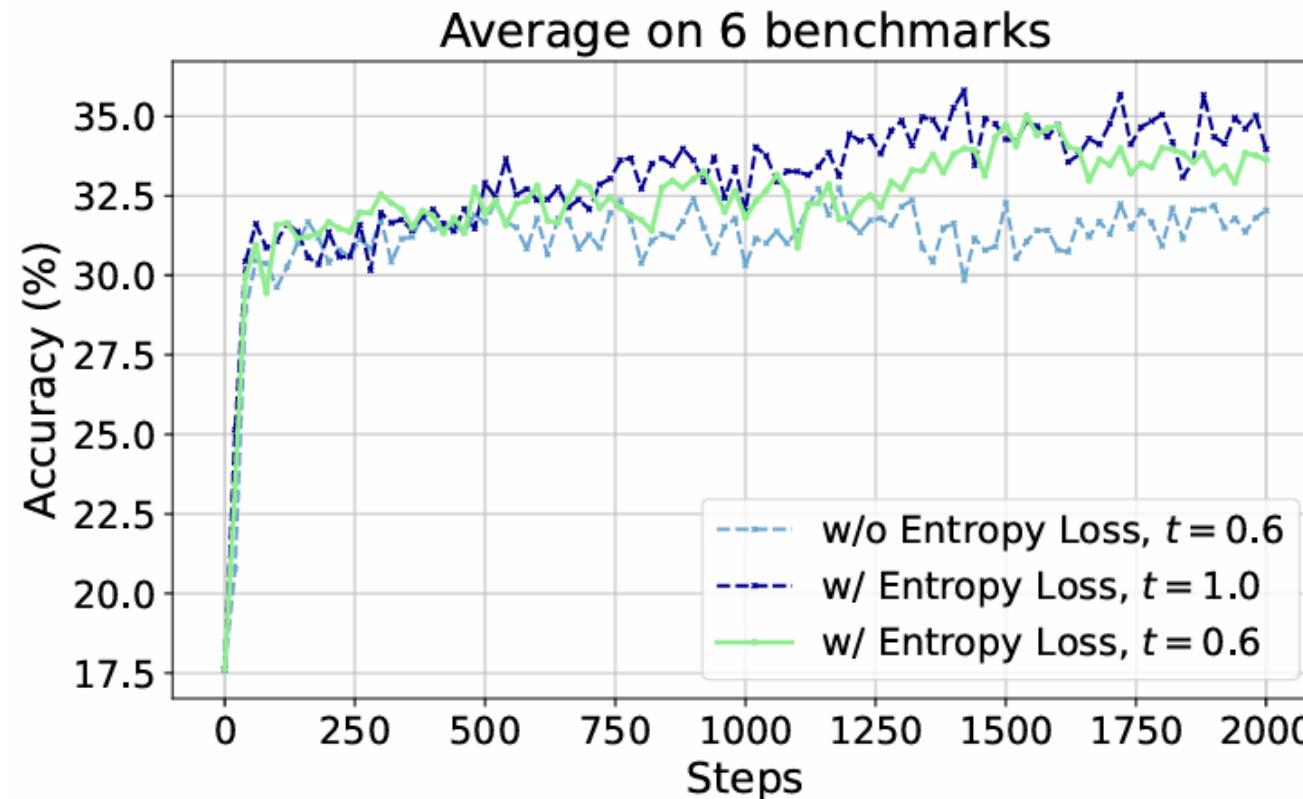
(3) Adding proper entropy loss can further improve performance based on policy loss.

| Row | Policy Loss | Weight Decay | KL Loss | Entropy Loss | Label  | Training Convergence | MATH 500 | AIME 2024 |
|-----|-------------|--------------|---------|--------------|--------|----------------------|----------|-----------|
| 1   |             |              |         |              | 12.8   | NO                   | 39.8     | 7.5       |
| 2   | +           |              |         |              | 12.8   | YES                  | 71.8     | 15.4      |
| 3   | +           | +            |         |              | 12.8   | YES                  | 71.4     | 16.3      |
| 4   | +           | +            | +       |              | 12.8   | YES                  | 70.8     | 15.0      |
| 5   | +           | +            | +       | +            | 12.8   | YES                  | 74.8     | 17.5      |
| 6   | +           | +            | +       | +,-0.003     | 12.8   | YES                  | 73.6     | 15.4      |
| 7   | +           |              |         | +            | 12.8   | YES                  | 75.6     | 17.1      |
| 8   |             | +            | +       |              | 12.8   | NO                   | 39.0     | 10.0      |
| 9   |             | +            | +       | +            | 12.8   | NO                   | 65.4     | 7.1       |
| 10  |             |              |         | +            | 12.8   | NO                   | 63.4     | 8.8       |
| 11  | +           | +            | +       | +            | 12.7   | YES                  | 73.4     | 17.9      |
| 12  | +           | +            | +       | +            | 4      | YES                  | 57.0     | 9.2       |
| 13  | +           | +            | +       | +            | 929725 | NO                   | 64.4     | 9.6       |

# Ablation Study

---

(3) Adding **proper entropy loss** can further improve performance based on policy loss. It can be important for post-saturation generalization, showing **the importance of encouraging exploration**



# Ablation Study

(4) Simply adding entropy loss alone can still improve model performance.

| Row | Policy Loss | Weight Decay | KL Loss | Entropy Loss | Label  | Training Convergence | MATH 500    | AIME 2024   |
|-----|-------------|--------------|---------|--------------|--------|----------------------|-------------|-------------|
| 1   |             |              |         |              | 12.8   | NO                   | 39.8        | 7.5         |
| 2   | +           |              |         |              | 12.8   | YES                  | 71.8        | 15.4        |
| 3   | +           | +            |         |              | 12.8   | YES                  | 71.4        | 16.3        |
| 4   | +           | +            | +       |              | 12.8   | YES                  | 70.8        | 15.0        |
| 5   | +           | +            | +       | +            | 12.8   | YES                  | <u>74.8</u> | <b>17.5</b> |
| 6   | +           | +            | +       | +,-0.003     | 12.8   | YES                  | 73.6        | 15.4        |
| 7   | +           |              |         | +            | 12.8   | YES                  | <b>75.6</b> | <u>17.1</u> |
| 8   |             | +            | +       |              | 12.8   | NO                   | 39.0        | 10.0        |
| 9   |             | +            | +       | +            | 12.8   | NO                   | 65.4        | 7.1         |
| 10  |             |              |         | +            | 12.8   | NO                   | 63.4        | 8.8         |
| 11  | +           | +            | +       | +            | 12.7   | YES                  | 73.4        | 17.9        |
| 12  | +           | +            | +       | +            | 4      | YES                  | 57.0        | 9.2         |
| 13  | +           | +            | +       | +            | 929725 | NO                   | 64.4        | 9.6         |

# Ablation Study

---

(4) Simply adding entropy loss alone can still improve model performance.

Table 6: Entropy loss alone with  $\pi_1$  can still improve model performance.

| Model                                                        | MATH<br>500  | AIME24<br>2024 | AMC23<br>2023 | Minerva<br>Math | Olympiad-<br>Bench | AIME<br>2025 | Avg.         |
|--------------------------------------------------------------|--------------|----------------|---------------|-----------------|--------------------|--------------|--------------|
| <b>Qwen2.5-Math-1.5B</b><br>+Entropy Loss, Train 20 step     | 36.0<br>63.4 | 6.7<br>8.8     | 28.1<br>33.8  | 8.1<br>14.3     | 22.2<br>26.5       | 4.6<br>3.3   | 17.6<br>25.0 |
| <b>Llama-3.2-3B-Instruct</b><br>+Entropy Loss, Train 10 step | 40.8<br>47.8 | 8.3<br>8.8     | 25.3<br>26.9  | 15.8<br>18.0    | 13.2<br>15.1       | 1.7<br>0.4   | 17.5<br>19.5 |
| <b>Qwen2.5-Math-7B</b><br>+Entropy Loss, Train 4 step        | 51.0<br>57.2 | 12.1<br>13.3   | 35.3<br>39.7  | 11.0<br>14.3    | 18.2<br>21.5       | 6.7<br>3.8   | 22.4<br>25.0 |

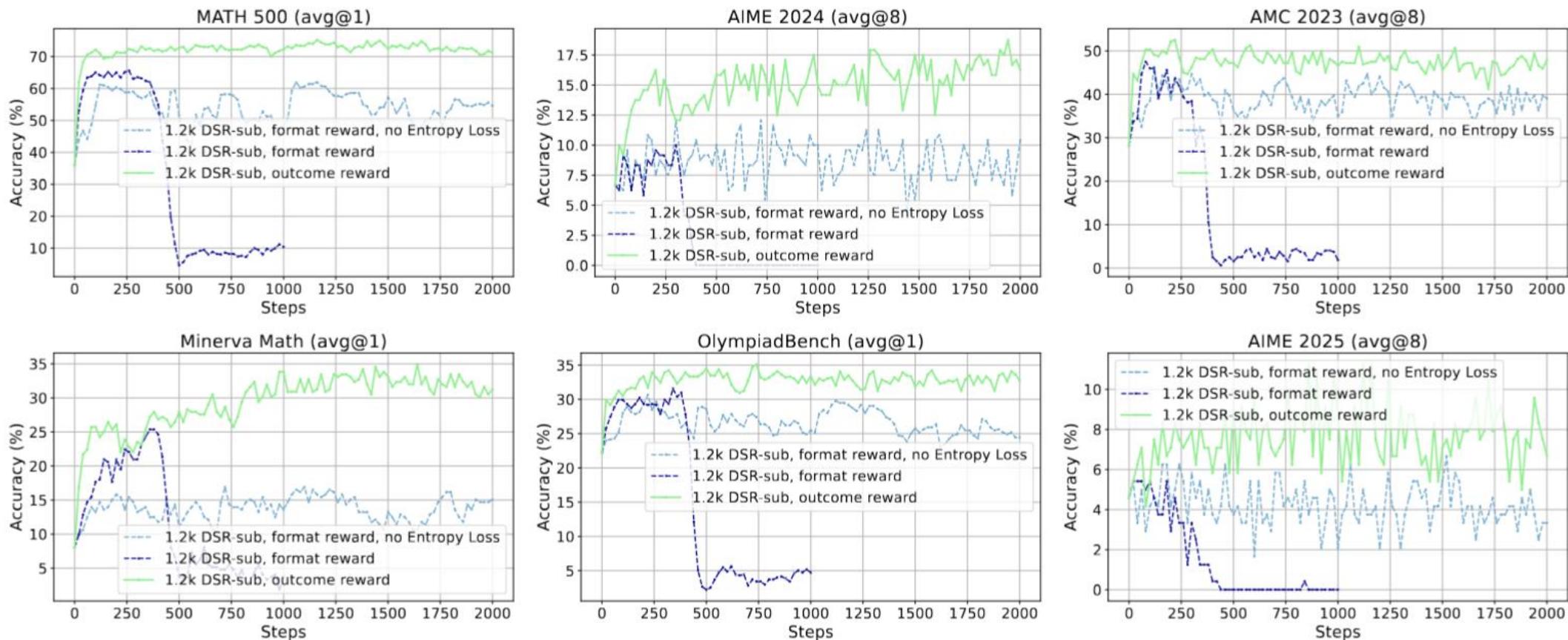
# Ablation Study

(4) Simply adding entropy loss alone can still improve model performance. So when the label is wrong, model still has some improvement from 1-shot RLVR.

| Row | Policy Loss | Weight Decay | KL Loss | Entropy Loss | Label  | Training Convergence | MATH 500    | AIME 2024   |
|-----|-------------|--------------|---------|--------------|--------|----------------------|-------------|-------------|
| 1   |             |              |         |              | 12.8   | NO                   | 39.8        | 7.5         |
| 2   | +           |              |         |              | 12.8   | YES                  | 71.8        | 15.4        |
| 3   | +           | +            |         |              | 12.8   | YES                  | 71.4        | 16.3        |
| 4   | +           | +            | +       |              | 12.8   | YES                  | 70.8        | 15.0        |
| 5   | +           | +            | +       | +            | 12.8   | YES                  | 74.8        | <b>17.5</b> |
| 6   | +           | +            | +       | , -0.003     | 12.8   | YES                  | 73.6        | 15.4        |
| 7   | +           |              |         | +            | 12.8   | YES                  | <b>75.6</b> | <u>17.1</u> |
| 8   |             | +            | +       |              | 12.8   | NO                   | 39.0        | 10.0        |
| 9   |             | +            | +       | +            | 12.8   | NO                   | 65.4        | 7.1         |
| 10  |             |              |         | +            | 12.8   | NO                   | 63.4        | 8.8         |
| 11  | +           | +            | +       | +            | 12.7   | YES                  | 73.4        | 17.9        |
| 12  | +           | +            | +       | +            | 4      | YES                  | 57.0        | 9.2         |
| 13  | +           | +            | +       | +            | 929725 | NO                   | 64.4        | 9.6         |

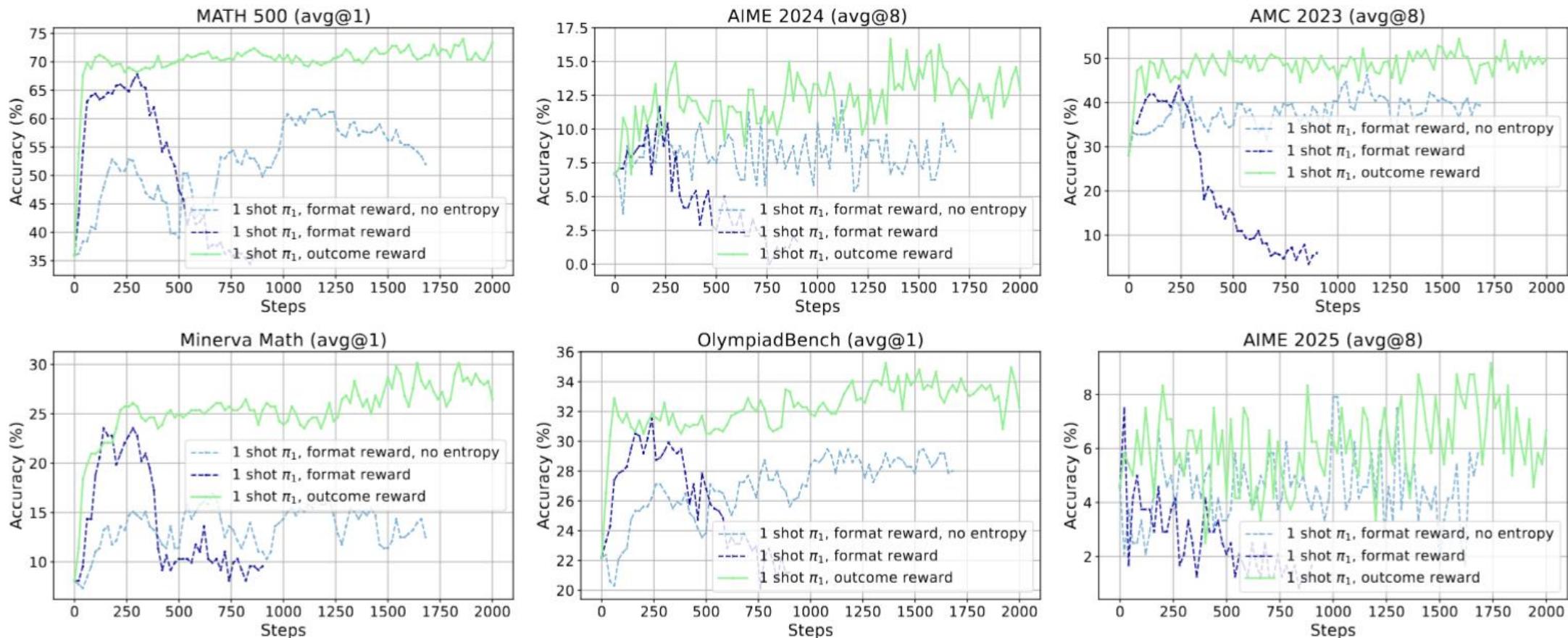
# (Only) Format Fixing?

- Only format reward can improve a lot on Qwen2.5-Math-1.5B
- Still has a gap with outcome reward
- (These two holds for both **full-set RLVR** and **1-shot RLVR!**)



# (Only) Format Fixing?

- Only format reward can improve a lot on Qwen2.5-Math-1.5B
- Still has a gap with outcome reward
- (These two holds for both **full-set RLVR** and **1-shot RLVR!**)



# (Only) Format Fixing?

- Only format reward can improve **a lot** on Qwen2.5-Math-1.5B
- Still has a gap with outcome reward
- (These two holds for both **full-set RLVR** and **1-shot RLVR!**)

Table 12: **RLVR with format reward can still improve model performance significantly, while still having a gap compared with that using outcome reward.** Here we consider adding entropy loss or not for format reward. Detailed results are also in Fig. 12 and Fig. 13.

| Dataset     | Reward Type | Entropy Loss | MATH 500    | AIME 2024   | AMC 2023    | Minerva Math | Olympiad-Bench | AIME 2025  | Avg.        |
|-------------|-------------|--------------|-------------|-------------|-------------|--------------|----------------|------------|-------------|
| NA          | NA          | NA           | 36.0        | 6.7         | 28.1        | 8.1          | 22.2           | 4.6        | 17.6        |
| DSR-sub     | Outcome     | +            | <b>73.6</b> | <b>17.1</b> | <b>50.6</b> | <b>32.4</b>  | <b>33.6</b>    | <b>8.3</b> | <b>35.9</b> |
| DSR-sub     | Format      | +            | 65.0        | 8.3         | 45.9        | 17.6         | 29.9           | 5.4        | 28.7        |
| DSR-sub     | Format      |              | 61.4        | 9.6         | 44.7        | 16.5         | 29.5           | 3.8        | 27.6        |
| $\{\pi_1\}$ | Outcome     | +            | <b>72.8</b> | <b>15.4</b> | <b>51.6</b> | <b>29.8</b>  | <b>33.5</b>    | <b>7.1</b> | <b>35.0</b> |
| $\{\pi_1\}$ | Format      | +            | 65.4        | 8.8         | 43.8        | 22.1         | 31.6           | 3.8        | 29.2        |
| $\{\pi_1\}$ | Format      |              | 61.6        | 8.3         | 46.2        | 15.4         | 29.3           | 4.6        | 27.6        |

# (Only) Format Fixing?

- Fixing format and improving general reasoning happen at the same time

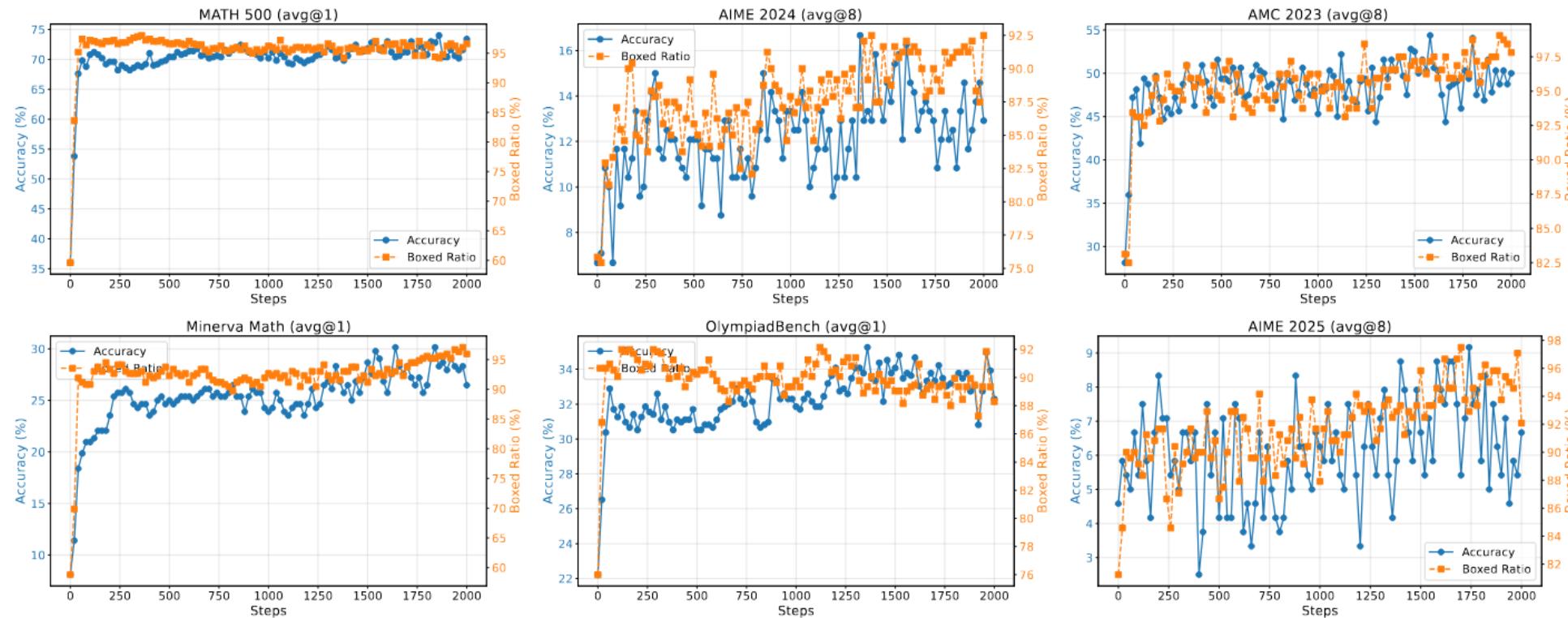


Figure 14: Relation between the number of `\boxed{}` and test accuracy. We can see that they have a strong positive correlation. However, after the number of `\boxed{}` enters a plateau, the evaluation results on some evaluation tasks continue improving (like Minerva Math, OlympiadBench and MATH500).

# (Only) Format Fixing?

- Fixing format and improving general reasoning happen at the same time

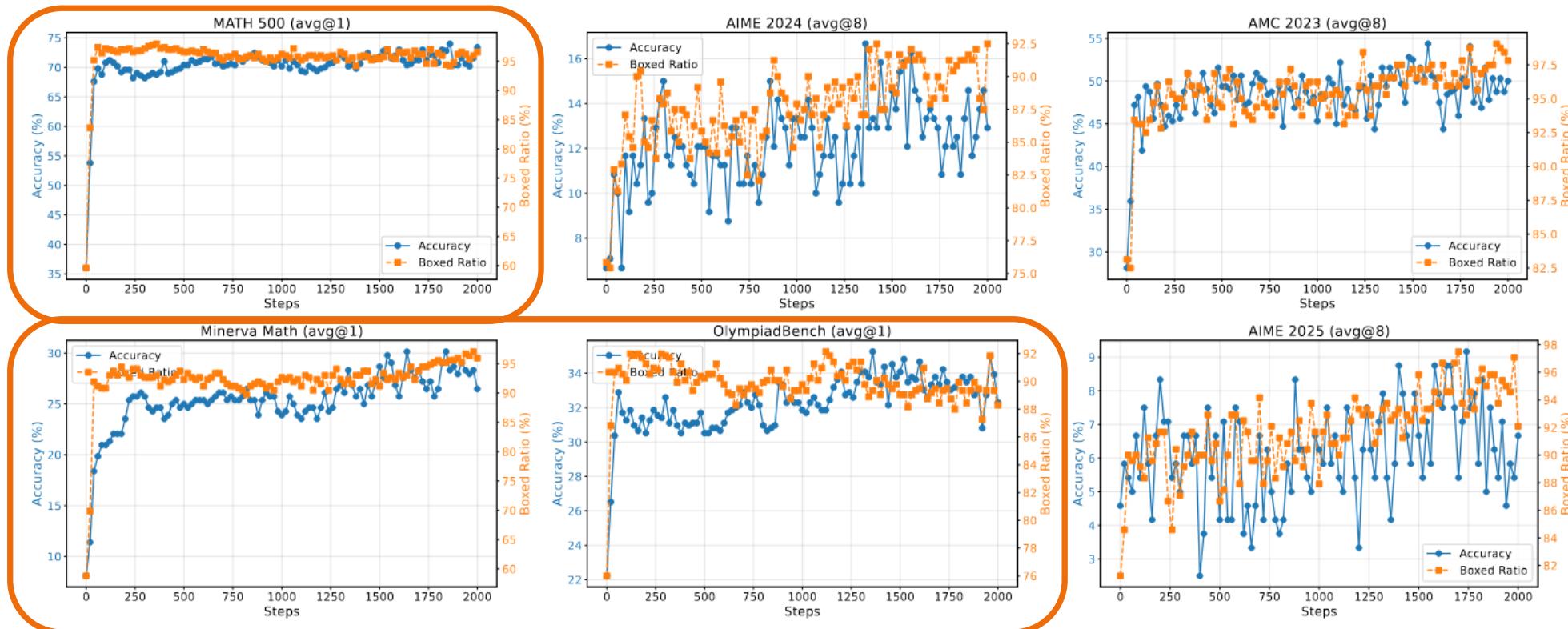


Figure 14: **Relation between the number of \boxed{} and test accuracy.** We can see that they have a strong positive correlation. However, after the number of `\boxed{}` enters a plateau, the evaluation results on some evaluation tasks continue improving (like Minerva Math, OlympiadBench and MATH500).

# (Only) Format Fixing?

- Fixing format and improving general reasoning happen at the same time

Table 14: 1-shot RLVR does not do something like put the answer into the `\boxed{}`. “Ratio of disagreement” means the ratio of questions that has different judgement between Qwen-Eval and QwQ-32B judge. Here we let QwQ-32B judged based on if the output contain correct answer, without considering if the answer is put in the `\boxed{}`.

|                                | Step 0 | Step 20 | Step 60 | Step 500 | Step 1300 | Step 1860 |
|--------------------------------|--------|---------|---------|----------|-----------|-----------|
| Ratio of <code>\boxed{}</code> | 59.6%  | 83.6%   | 97.4%   | 96.6%    | 96.6%     | 94.2%     |
| Acc. judge by Qwen-Eval        | 36.0   | 53.8    | 69.8    | 70.4     | 72.2      | 74.0      |
| Acc. judge by QwQ-32B          | 35.8   | 57.2    | 70.6    | 71.8     | 73.6      | 74.6      |
| Ratio of disagreement          | 4.2%   | 5%      | 1.2%    | 1.4%     | 1.8%      | 1.8%      |



# (Only) Format Fixing?

- Analysis: Qwen2.5-Math families have a lot of repetitive outputs, which results that fixing format itself brings lots of improvement
    - Qwen2.5-Math-1.5B: **~40% outputs** contain infinite loop output in MATH500!
    - Qwen2.5-Math-7B: **~20% outputs** contain infinite loop output in MATH500!

根据题意，设  $x$  为总人数，则有  $\frac{1}{3}x = 14$ ，解得  $x = 42$ 。所以，总人数是 42 人。

# (Only) Format Fixing?

- Analysis: Qwen2.5-Math families have a lot of repetitive outputs, which results that fixing format itself brings lots of improvement
    - Qwen2.5-Math-1.5B: **~40% outputs** contain infinite loop output in MATH500!
    - Qwen2.5-Math-7B: **~20% outputs** contain infinite loop output in MATH500!

NO", "gt": "9", "score": [false], "code": ["What is the largest 3-digit number that is a multiple of 12?"]}

# (Only) Format Fixing?

- Analysis: Qwen2.5-Math families have a lot of repetitive outputs, which results that fixing format itself brings lots of improvement
    - Qwen2.5-Math-1.5B: **~40% outputs** contain infinite loop output in MATH500!
    - Qwen2.5-Math-7B: **~20% outputs** contain infinite loop output in MATH500!

# (Only) Format Fixing?

---

- Analysis: Qwen2.5-Math families have a lot of repetitive outputs, which results that fixing format itself brings lots of improvement
  - Qwen2.5-Math-1.5B: **~40% outputs** contain infinite loop output in MATH500!
  - Qwen2.5-Math-7B: **~20% outputs** contain infinite loop output in MATH500!

*Maybe in the future, a necessary **baseline** will be **RLVR** with format reward  
(or strongest prompt)*

# Pi1 for in-context learning

---

- In-context learning: add a “**Question-Answer example**” (here is pi1) before evaluating downstream question

<|im\_start|>system

Please reason step by step, and put your final answer within \boxed{}.<|im\_end|>

<|im\_start|>user

Question: The pressure  $\langle P \rangle$  exerted by wind on a sail varies jointly as the area  $\langle A \rangle$  of ...

Answer: Given:

-  $\langle P \rangle \propto A \cdot V^3$

-  $\langle P = k \cdot A \cdot V^3 \rangle$  where  $\langle k \rangle$  is the constant of proportionality. Using the given data:

...

Therefore, the wind velocity when the pressure on  $\langle 4 \rangle$  square feet of sail is  $\langle 32 \rangle$  pounds is approximately  $\langle 12.7 \rangle$  miles per hour.

Question: Find the sum of all integer bases  $b > 9$  for which  $17_b$  is a divisor of  $97_b$ .

Answer:<|im\_end|>

<|im\_start|>assistant

# Pi1 for in-context learning

- Pi1 can even improve Qwen2.5-Math-7B's MATH500 from **51.0 -> 75.4**, and OlympiadBench from **18.2 -> 41.3** with in-context learning!!
- Perform much better than Qwen's official 4 examples on these two models

Table 13:  $\pi_1$  even performs well for in-context learning on Qwen2.5-Math-7B.

| Dataset                              | Method     | MATH<br>500 | AIME<br>2024 | AMC<br>2023 | Minerva<br>Math | Olympiad-<br>Bench | AIME<br>2025 | Avg.        |
|--------------------------------------|------------|-------------|--------------|-------------|-----------------|--------------------|--------------|-------------|
| <b>Qwen2.5-Math-1.5B</b>             |            |             |              |             |                 |                    |              |             |
| NA                                   | NA         | 36.0        | 6.7          | 28.1        | 8.1             | 22.2               | 4.6          | 17.6        |
| { $\pi_1$ }                          | RLVR       | 72.8        | <b>15.4</b>  | <b>51.6</b> | <b>29.8</b>     | <b>33.5</b>        | <b>7.1</b>   | <b>35.0</b> |
| Qwen official 4 examples for MATH500 | In-Context | 59.0        | 8.3          | 34.7        | 19.9            | 25.6               | 5.4          | 25.5        |
| Qwen official Example 1 for MATH500  | In-Context | 49.8        | 1.7          | 16.9        | 19.9            | 19.9               | 0.0          | 18.0        |
|                                      | In-Context | 34.6        | 2.5          | 14.4        | 12.1            | 21.0               | 0.8          | 14.2        |
| <b>Qwen2.5-Math-7B</b>               |            |             |              |             |                 |                    |              |             |
| NA                                   | NA         | 51.0        | 12.1         | 35.3        | 11.0            | 18.2               | 6.7          | 22.4        |
| { $\pi_1$ }                          | RLVR       | <b>79.2</b> | <b>23.8</b>  | <b>60.3</b> | 27.9            | 39.1               | 10.8         | <b>40.2</b> |
| { $\pi_1$ }                          | In-Context | 75.4        | 15.8         | 48.4        | <b>30.1</b>     | <b>41.3</b>        | <b>13.3</b>  | 37.4        |
| Qwen official 4 examples for MATH500 | In-Context | 59.2        | 4.2          | 20.9        | 20.6            | 24.4               | 0.8          | 21.7        |
| Qwen official Example 1 for MATH500  | In-Context | 54.0        | 4.2          | 23.4        | 18.4            | 21.2               | 2.1          | 20.6        |

# Pi1 for in-context learning

---

**Still tricky:**

- **Not work for all models**, like fail on Qwen2.5-Math-72B and Llama3.2-3B-Instruct (slightly worse than Qwen's official 4 examples)
- **Highly example-dependent**. Pi13 works well on RLV, but fail on in-context learning (worse than original zero-shot learning)

# Application: Does RLVR has high label robustness?

- In RLVR training, 1100 data with wrong labels + 100 data with correct labels can performs worse than 1 data with correct label.

Table 15: **Influence of Random Wrong Labels.** Here “Error Rate” means the ratio of data that has the random wrong labels.

| Dataset                         | Error Rate | MATH 500 | AIME 2024 | AMC 2023 | Minerva Math | Olympiad-Bench | AIME 2025 | Avg. |
|---------------------------------|------------|----------|-----------|----------|--------------|----------------|-----------|------|
| NA                              | NA         | 36.0     | 6.7       | 28.1     | 8.1          | 22.2           | 4.6       | 17.6 |
| <b>Qwen2.5-Math-1.5B + GRPO</b> |            |          |           |          |              |                |           |      |
| DSR-sub                         | 0%         | 73.6     | 17.1      | 50.6     | 32.4         | 33.6           | 8.3       | 35.9 |
| DSR-sub                         | 60%        | 71.8     | 17.1      | 47.8     | 29.4         | 34.4           | 7.1       | 34.6 |
| DSR-sub                         | 90%        | 67.8     | 14.6      | 46.2     | 21.0         | 32.3           | 5.4       | 31.2 |
| $\{\pi_1\}$                     | 0%         | 72.8     | 15.4      | 51.6     | 29.8         | 33.5           | 7.1       | 35.0 |
| <b>Qwen2.5-Math-1.5B + PPO</b>  |            |          |           |          |              |                |           |      |
| DSR-sub                         | 0%         | 72.8     | 19.2      | 48.1     | 27.9         | 35.0           | 9.6       | 35.4 |
| DSR-sub                         | 60%        | 71.6     | 13.3      | 49.1     | 27.2         | 34.4           | 12.1      | 34.6 |
| DSR-sub                         | 90%        | 68.2     | 15.8      | 50.9     | 26.1         | 31.9           | 4.6       | 32.9 |
| $\{\pi_1\}$                     | 0%         | 72.4     | 11.7      | 51.6     | 26.8         | 33.3           | 7.1       | 33.8 |

# Discussion

---

- Base models already has strong reasoning capability, and its prior affects a lot for the RLVR stage.
- How to select/curate proper data for RLVR is critical
  - 1-shot RLVR works does not necessarily means that scaling RL dataset is useless
- How to understand 1-shot RLVR and post-saturation generalization?
  - policy loss has implicit generalization
- Better exploration (entropy loss is unstable).
- Other domain (code) and application (label robustness)

# Spurious Reward & Data Contamination

---

## Spurious Rewards: Rethinking Training Signals in RLVR

---

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

<sup>1</sup>University of Washington <sup>2</sup>Allen Institute for Artificial Intelligence

<sup>3</sup>University of California, Berkeley

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



GitHub Repo

## REASONING OR MEMORIZATION? UNRELIABLE RESULTS OF REINFORCEMENT LEARNING DUE TO DATA CONTAMINATION

Mingqi Wu<sup>1\*</sup>, Zhihao Zhang<sup>1,2\*</sup>, Qiaole Dong<sup>1\*</sup>,  
Zhiheng Xi<sup>1</sup>, Jun Zhao<sup>1</sup>, Senjie Jin<sup>1</sup>, Xiaoran Fan<sup>1</sup>, Yuhao Zhou<sup>1</sup>,  
Huijie Lv<sup>1,2</sup>, Ming Zhang<sup>1</sup>, Yanwei Fu<sup>1</sup>, Qin Liu<sup>3</sup>, Songyang Zhang<sup>2</sup>, Qi Zhang<sup>1,2†</sup>

<sup>1</sup> Fudan University

<sup>2</sup> Shanghai Artificial Intelligence Laboratory

<sup>3</sup> University of California, Davis

{qz}@fudan.edu.cn {qinli}@ucdavis.edu

# Spurious Reward & Data Contamination

---

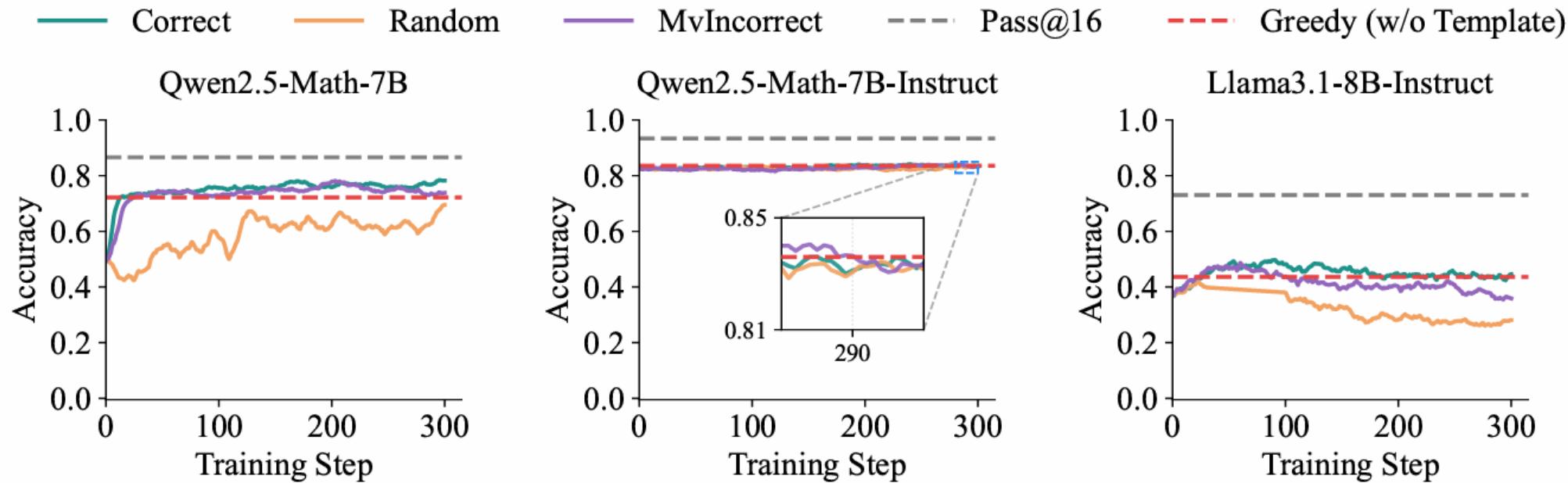


Figure 3: Accuracy on the **MATH-500** for Qwen2.5-Math-7B, Qwen2.5-Math-7B-Instruct, and Llama3.1-8B-Instruct trained with RLVR under various reward signals. Greedy and pass@16 scores are reported *without* template.

# Spurious Reward & Data Contamination

---

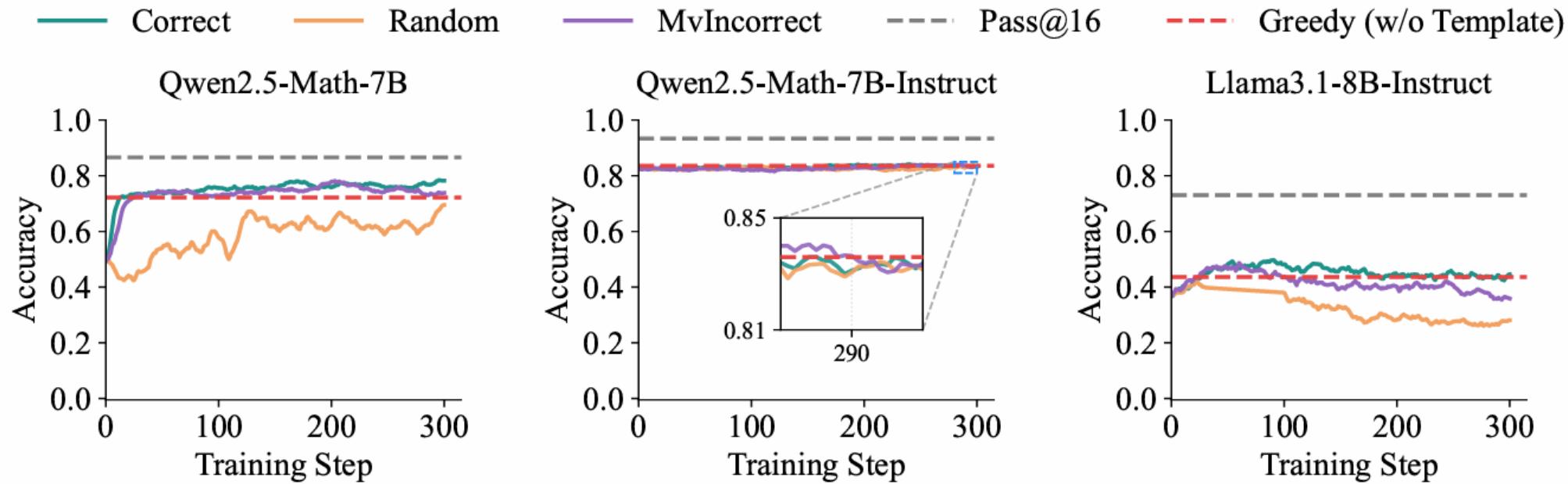


Figure 3: Accuracy on the **MATH-500** for Qwen2.5-Math-7B, Qwen2.5-Math-7B-Instruct, and Llama3.1-8B-Instruct trained with RLVR under various reward signals. Greedy and pass@16 scores are reported *without* template.

# Spurious Reward & Data Contamination

---

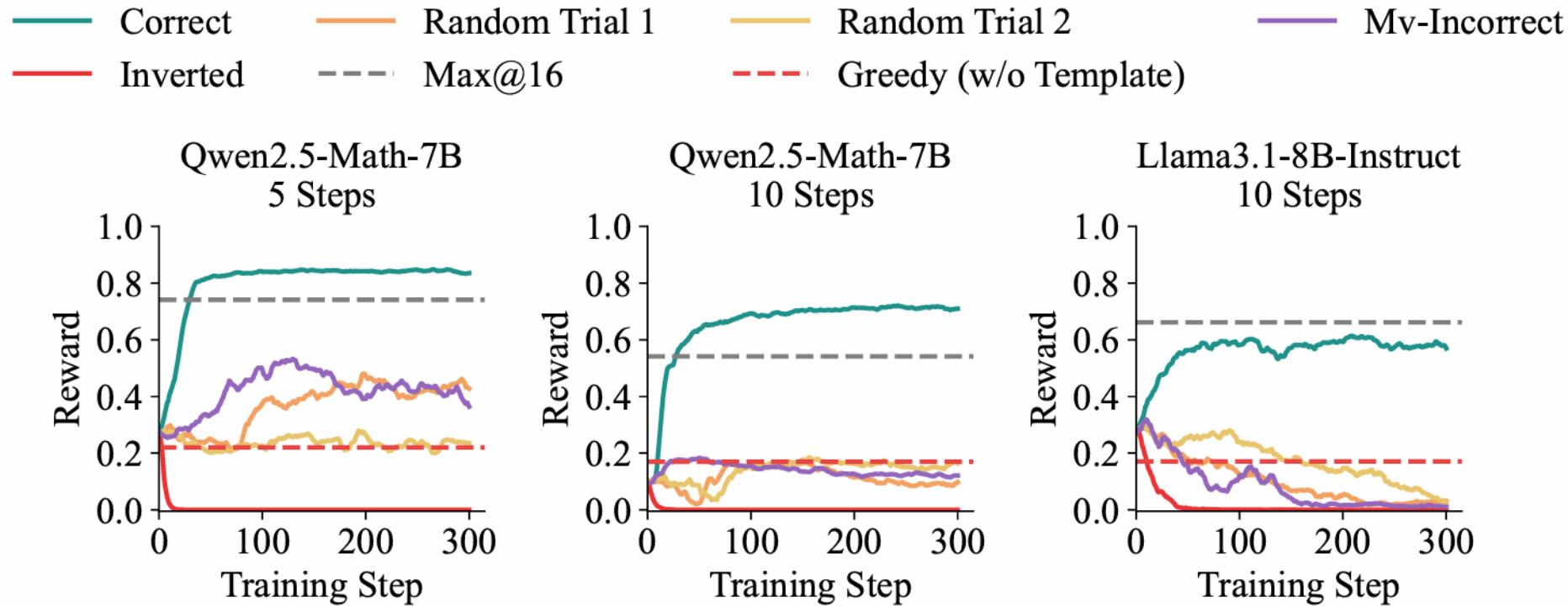


Figure 7: Reward of Qwen2.5-Math-7B and Llama3.1-8B-Instruct on *RandomCalculation*. Results are presented for datasets with 5-step and 10-step calculations.

# Spurious Reward & Data Contamination

---

Table 2: Accuracy (Exact Match, EM) and ROUGE-L scores on several datasets (lower scores in gray ) under different prompt prefix ratios in greedy decoding mode without applying chat template, namely *Greedy (w/o Template)* configuration.

| Model           | Dataset       | Size | 80%-Problem |              | 60%-Problem |              | 40%-Problem |              |
|-----------------|---------------|------|-------------|--------------|-------------|--------------|-------------|--------------|
|                 |               |      | RougeL      | EM           | RougeL      | EM           | RougeL      | EM           |
| Qwen2.5-Math-7B | MATH-500      | 500  | 81.25       | <b>65.80</b> | 78.06       | <b>54.60</b> | 69.01       | <b>39.20</b> |
|                 | AMC           | 83   | 77.38       | <b>55.42</b> | 70.25       | <b>42.17</b> | 75.17       | <b>36.14</b> |
|                 | AIME2024      | 30   | 74.04       | <b>56.67</b> | 55.31       | 20.00        | 57.72       | 16.67        |
|                 | AIME2025      | 30   | 54.71       | 16.67        | 34.88       | 0.00         | 27.43       | 0.00         |
|                 | MinervaMath   | 272  | 36.08       | 2.94         | 31.22       | 0.37         | 29.35       | 0.00         |
|                 | LiveMathBench | 100  | 42.76       | 5.00         | 32.78       | 0.00         | 29.97       | 0.00         |
| Qwen2.5-7B      | MATH-500      | 500  | 66.42       | <b>40.20</b> | 60.98       | 21.20        | 50.36       | 8.20         |
|                 | AMC           | 83   | 73.24       | <b>49.40</b> | 64.42       | 33.73        | 63.79       | 28.92        |
|                 | AIME2024      | 30   | 59.80       | <b>30.00</b> | 48.69       | 13.33        | 44.65       | 10.00        |
|                 | AIME2025      | 30   | 54.61       | 10.00        | 37.59       | 0.00         | 30.30       | 0.00         |
|                 | MinervaMath   | 272  | 35.24       | 2.94         | 32.35       | 0.37         | 27.89       | 0.00         |
|                 | LiveMathBench | 100  | 41.15       | 4.00         | 32.74       | 0.00         | 27.95       | 0.00         |
| Llama3.1-8B     | MATH-500      | 500  | 48.33       | 17.80        | 40.55       | 3.80         | 32.07       | 0.60         |
|                 | AMC           | 83   | 44.54       | 4.82         | 30.62       | 0.00         | 27.10       | 0.00         |
|                 | AIME2024      | 30   | 50.50       | 13.33        | 30.80       | 0.00         | 26.08       | 0.00         |
|                 | AIME2025      | 30   | 47.04       | 10.00        | 33.49       | 0.00         | 25.20       | 0.00         |
|                 | MinervaMath   | 272  | 36.24       | 2.21         | 29.52       | 0.00         | 27.11       | 0.00         |
|                 | LiveMathBench | 100  | 35.55       | 5.00         | 31.93       | 0.00         | 26.88       | 0.00         |

# Spurious Reward & Data Contamination

Table 2: Accuracy (Exact Match, EM) and ROUGE-L scores on several datasets (lower scores in gray ) under different prompt prefix ratios in greedy decoding mode without applying chat template, namely *Greedy (w/o Template)* configuration.

| Model           | Dataset       | Size | 80%-Problem |              | 60%-Problem |              | 40%-Problem |              |
|-----------------|---------------|------|-------------|--------------|-------------|--------------|-------------|--------------|
|                 |               |      | RougeL      | EM           | RougeL      | EM           | RougeL      | EM           |
| Qwen2.5-Math-7B | MATH-500      | 500  | 81.25       | <b>65.80</b> | 78.06       | <b>54.60</b> | 69.01       | <b>39.20</b> |
|                 | AMC           | 83   | 77.38       | <b>55.42</b> | 70.25       | <b>42.17</b> | 75.17       | <b>36.14</b> |
|                 | AIME2024      | 30   | 74.04       | <b>56.67</b> | 55.31       | 20.00        | 57.72       | 16.67        |
|                 | AIME2025      | 30   | 54.71       | 16.67        | 34.88       | 0.00         | 27.43       | 0.00         |
|                 | MinervaMath   | 272  | 36.08       | 2.94         | 31.22       | 0.37         | 29.35       | 0.00         |
|                 | LiveMathBench | 100  | 42.76       | 5.00         | 32.78       | 0.00         | 29.97       | 0.00         |

| RL Dataset                               | Dataset Size | MATH 500                                      | AIME 2024                                     | AMC 2023                                      | Minerva Math                                  | Olympiad-Bench                                | AIME 2025                                   | Avg.                                          |
|------------------------------------------|--------------|-----------------------------------------------|-----------------------------------------------|-----------------------------------------------|-----------------------------------------------|-----------------------------------------------|---------------------------------------------|-----------------------------------------------|
| <b>Qwen2.5-Math-7B [24] + GRPO</b>       |              |                                               |                                               |                                               |                                               |                                               |                                             |                                               |
| NA DSR-sub                               | NA 1209      | 51.0 <sub>+0.0</sub><br>78.6 <sub>+27.6</sub> | 12.1 <sub>+0.0</sub><br>25.8 <sub>+13.7</sub> | 35.3 <sub>+0.0</sub><br>62.5 <sub>+27.2</sub> | 11.0 <sub>+0.0</sub><br>33.8 <sub>+22.8</sub> | 18.2 <sub>+0.0</sub><br>41.6 <sub>+23.4</sub> | 6.7 <sub>+0.0</sub><br>14.6 <sub>+7.9</sub> | 22.4 <sub>+0.0</sub><br>42.8 <sub>+20.4</sub> |
| { $\pi_1$ }                              | 1            | <b>79.2</b> <sub>+28.2</sub>                  | 23.8 <sub>+11.7</sub>                         | 60.3 <sub>+25.0</sub>                         | 27.9 <sub>+16.9</sub>                         | 39.1 <sub>+20.9</sub>                         | 10.8 <sub>+4.1</sub>                        | 40.2 <sub>+17.8</sub>                         |
| { $\pi_1, \pi_{13}$ }                    | 2            | <b>79.2</b> <sub>+28.2</sub>                  | 21.7 <sub>+9.6</sub>                          | 58.8 <sub>+23.3</sub>                         | 35.3 <sub>+24.3</sub>                         | 40.9 <sub>+22.7</sub>                         | 12.1 <sub>+5.4</sub>                        | 41.3 <sub>+18.9</sub>                         |
| { $\pi_1, \pi_2, \pi_{13}, \pi_{1209}$ } | 4            | 78.6 <sub>+27.6</sub>                         | 22.5 <sub>+10.4</sub>                         | 61.9 <sub>+26.6</sub>                         | <b>36.0</b> <sub>+25.0</sub>                  | 43.7 <sub>+25.5</sub>                         | 12.1 <sub>+5.4</sub>                        | 42.5 <sub>+20.1</sub>                         |
| Random                                   | 16           | 76.0 <sub>+25.0</sub>                         | 22.1 <sub>+10.0</sub>                         | <b>63.1</b> <sub>+27.3</sub>                  | 31.6 <sub>+20.6</sub>                         | 35.6 <sub>+17.4</sub>                         | 12.9 <sub>+6.2</sub>                        | 40.2 <sub>+17.8</sub>                         |
| { $\pi_1, \dots, \pi_{16}$ }             | 16           | 77.8 <sub>+26.8</sub>                         | <b>30.4</b> <sub>+18.3</sub>                  | 62.2 <sub>+26.0</sub>                         | 35.3 <sub>+24.3</sub>                         | 39.9 <sub>+21.7</sub>                         | 9.6 <sub>+2.9</sub>                         | 42.5 <sub>+20.1</sub>                         |

Format reward baseline for Minerva and aime25: 24.3 & 6.7

# Spurious Reward & Data Contamination

---

- I think data contamination indeed happens, but it would not make 1-shot RLVR's conclusion fail.
- Mid-training or pretraining with open-source training data is really important for research.

# Authors

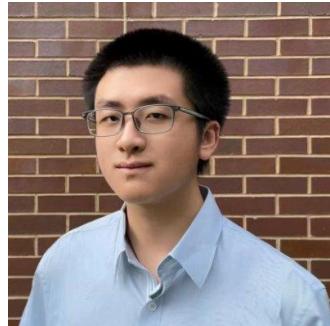
---



**Yiping  
Wang**



**Qing  
Yang**



**Zhiyuan  
Zeng**



**Liliang  
Ren**



**Lucas  
Liu**



**Baolin  
Peng**



**Hao  
Cheng**



**Xuehai  
He**



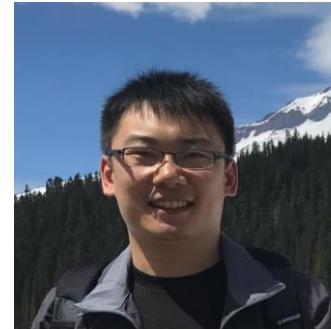
**Kuan  
Wang**



**Jianfeng  
Gao**



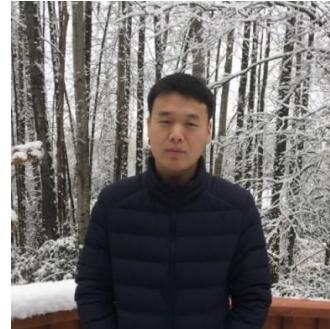
**Weizhu  
Chen**



**Shuohang  
Wang**



**Simon S.  
Du**



**Yelong  
Shen**

Thank You