

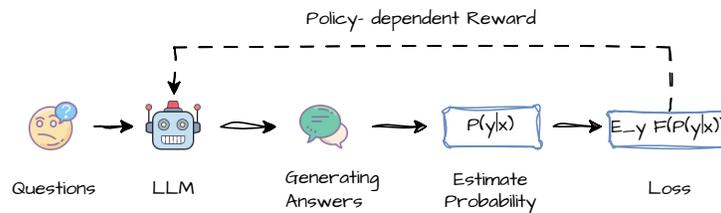
# CONFIDENCE IS ALL YOU NEED: FEW-SHOT RL FINE-TUNING OF LANGUAGE MODELS

Anonymous authors

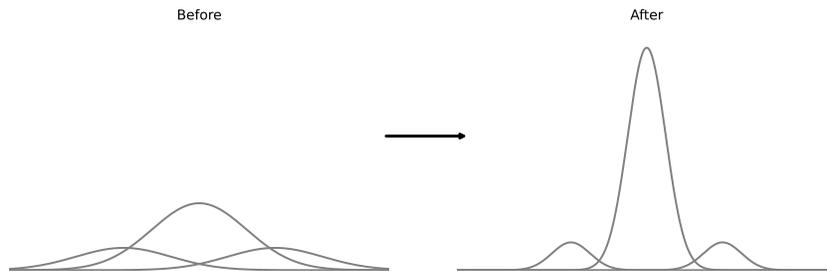
Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) have demonstrated strong performance on reasoning tasks, but post-training optimization remains essential for aligning their behavior with specific task objectives. Existing reinforcement learning (RL) approaches often rely on costly human annotations or external reward models, limiting their scalability in real-world applications. To address this, we propose Reinforcement Learning via Self-Confidence (RLSC)—a method that uses the model’s own confidence in its outputs as the reward signal, without requiring human labels, preference models, or manually crafted reward functions. RLSC is also highly sample-efficient: it only needs 1 to 8 samples per problem, and typically converges within 15 to 30 training steps. Under the Pass@1 evaluation metric, Qwen-Math-7B achieves significant performance improvements across several mathematical benchmarks: AIME2024 +6.7%, AMC23 +33.1%, Math500 +32.3%, Minerva +29.8%. On average, RLSC delivers a 23.68% improvement across these benchmarks. Notably, the effectiveness of RLSC is not limited to the Qwen series; it also leads to substantial performance gains on other mainstream models, including Olmo-7B, DeepSeek-R1-Distill-Qwen-1.5B, DeepSeek-R1-Distill-Llama-8B, Gemma-4B, and LLaMA-8B, etc. In summary, RLSC offers a simple, efficient, and scalable post-training method for pretrained language models, enabling significant performance gains with few training steps.



(a) Overview of the Reinforcement Learning via Self Confidence (RLSC) approach.



(b) Response probability distribution.

Figure 1: Combined visualization: (a) RL via Self Confidence workflow schema; (b) Probability distribution before and after training.

## 1 INTRODUCTION

Large language models [Achiam et al. \(2023\)](#) [Bai et al. \(2023\)](#) [Yang et al. \(2024b\)](#) [Liu et al. \(2024\)](#) [Guo et al. \(2025\)](#) [Guo et al. \(2024\)](#) [Grattafiori et al. \(2024\)](#) have demonstrated remarkable capabilities across various domains. However, to further align model behavior with specific task objectives, post-training optimization remains indispensable. Compared to supervised fine-tuning, reinforcement learning (RL) offers stronger exploration capabilities and is thus widely employed to enhance the reasoning performance of large models. Currently, reinforcement learning algorithms such as DPO [Rafailov et al. \(2023\)](#), PPO [Schulman et al. \(2017\)](#), RLHF [Ouyang et al. \(2022\)](#), and GRPO [Guo et al. \(2025\)](#) are extensively applied in the post-training of large language models. Among these, Reinforcement Learning with Verifiable Rewards (RLVR) [Zeng et al. \(2025\)](#) [Liu et al. \(2025\)](#) has become the mainstream approach for defining reward signals in RL training for large language models, significantly simplifying the complexity of reward design and improving training controllability and efficiency.

Although large language models have achieved remarkable progress in mathematical reasoning tasks, reinforcement learning (RL)-based fine-tuning methods often rely on costly human-labeled datasets or carefully designed reward functions. Performing RL training without labels could significantly reduce the cost of post-training optimization. Test-Time Reinforcement Learning (TTRL) [Zuo et al. \(2025\)](#) addresses this by generating up to 64 candidate answers per question and using majority voting to create pseudo-labels for training. However, this approach incurs substantial computational overhead, which limits its practicality for large-scale models.

To address the above challenges, we propose a new paradigm: *Reinforcement Learning via Self-Confidence (RLSC)*. Unlike methods that rely on external supervision—such as verifiable rewards or pseudo-label generation in TTRL, RLSC uses the model’s own confidence in its outputs as a reward signal, eliminating the need for any external supervision or manually constructed labels. By combining confidence analysis, the internal knowledge embedded within pre-trained language models can be effectively leveraged, enabling improved performance on downstream tasks—even without any external feedback.

We validate the proposed RLSC method on multiple models, including Qwen2.5-Math [Yang et al. \(2024b\)](#), Gemma [Team et al. \(2024\)](#), LLaMA [Grattafiori et al. \(2024\)](#), the distilled version of DeepSeek R1 [Guo et al. \(2025\)](#), and Olmo [Groeneveld et al. \(2024\)](#). Training was conducted on datasets such as GSM8K [Cobbe et al. \(2021\)](#), AIME24 [Li et al. \(2024b\)](#), Math500 [Hendrycks et al. \(2021c\)](#), Minerva Math [Hendrycks et al. \(2021b\)](#), and AMC23 [Li et al. \(2024a\)](#), using only 15–30 training steps. With a sampling temperature of 1, the model generated just 1 to 8 samples per problem and still achieved rapid performance gains.

All evaluations were performed under non-zero temperature sampling settings. RLSC showed particularly strong results on the Qwen2.5-Math-7B model: under the Pass@1 [Guo et al. \(2025\)](#) metric, it achieved a 6.67% improvement on AIME2024, 33.1% on AMC23, 32.3% on Math500, and 29.8% on Minerva Math. Across multiple math benchmarks, the average improvement reached 23.68%.

These results demonstrate that, when combined with the RLSC framework, a strong pre-trained language model can achieve significant gains in confidence and generalization—without the need for auxiliary datasets, human feedback, or handcrafted reward functions, and with only a few steps of training.

### Our main contributions are:

- A novel training paradigm:** We propose RLSC for the first time—an externally-supervision-free reinforcement learning method that uses the model’s own output confidence as a reward signal for self-supervised training.
- No reliance on labels or external rewards:** Unlike existing methods such as RLHF, DPO, PPO, GRPO, or TTRL, RLSC completely eliminates the need for costly human labels, reward models, or complex reward shaping, significantly reducing the cost of post-training.
- High training efficiency:** RLSC achieves convergence with only a few training steps (15–30) and a small number of samples (just 1 to 8 per problem), and can be fine-tuned in 30 minutes on a single 8×A100 GPU cluster using major LLMs such as Qwen2.5-Math and Gemma, etc.

108 **4. Outstanding generalization:** RLSC demonstrates significant performance improvements  
 109 on several mathematical benchmarks (AIME24, AMC23, Math500, Minerva, etc.), with an  
 110 average gain of 23.68%, maintaining strong performance on Pass@1 settings.  
 111

## 112 2 METHOD

### 113 2.1 FROM TTRL TO MODE SHARPENING

114 Test-Time Reinforcement Learning (TTRL) [Zuo et al. \(2025\)](#) improves large language models  
 115 (LLMs) by generating many outputs per input (typically 64) and applying majority voting to se-  
 116 lect the most frequent completion. This pseudo-label is then used to fine-tune the model. Although  
 117 this method has proven effective by experiments, it incurs substantial computational overhead during  
 118 training.  
 119

120 Inspired by the idea of majority voting, we asked the following question:

121 *what is the key underlying principle behind this voting process?*

122 Intuitively, majority voting selects the mode of the output distribution. Optimizing for agreement  
 123 between sampled completions sharpens the distribution: it increases the probability mass concen-  
 124 trated on the most likely answer. This also increases the chance of generating the same answer while  
 125 taking two independent samples.  
 126

127 Let  $p_\theta(y | x)$  denote the model probability of generating response  $y$  given input  $x$ , parameterized  
 128 by  $\theta$ . The probability that two independent samples from this distribution are identical is:

$$129 F(p_\theta) = \mathbb{E}_{y_1, y_2 \sim p_\theta(y|x)} [\mathbb{I}(y_1 = y_2)] = \sum_y p_\theta(y | x)^2 \quad (1)$$

130 This expression is maximized when the distribution collapses to a delta function centered on a single  
 131 most probable response - i.e., when the model is confident.  
 132

133 Therefore, we propose to directly maximize the following self-confidence objective:

$$134 F(p_\theta) = \mathbb{E}_{y \sim p_\theta(y|x)} [p_\theta(y | x)] \quad (2)$$

135 This objective acts as a continuation of the TTRL method. We suggest that it has similar effect on  
 136 probability distribution of model’s answers, while eliminating the need of generating large samples  
 137 that are needed for accurate majority voting. The key advantage is that it can work with any sample  
 138 size, even with a single output, and allows straightforward calculation based on model’s output with  
 139 no need for extra computation.  
 140

### 141 2.2 SELF-CONFIDENCE LOSS AND GRADIENT

142 To optimize the self-confidence objective 2, we compute its gradient with respect to the model’s  
 143 parameters  $\theta$ . Applying the log-trick, we obtain:  
 144

$$145 \begin{aligned} \nabla_\theta F(p_\theta) &= \sum_y \nabla_\theta p_\theta(y | x) p_\theta(y | x) \\ &= \mathbb{E}_{y \sim p_\theta} [\nabla_\theta p_\theta(y | x)] \\ &= \mathbb{E}_{y \sim p_\theta} [p_{old}(y | x) \nabla_\theta \log p_\theta(y | x)] \end{aligned} \quad (3)$$

146 Here,  $p_{old}(y | x)$  represents the probability of completion as it was at the time of model sampling.  
 147 As the rollouts are samples from the distribution by  $y \sim p_\theta$ , we can define the loss function as:

$$148 \mathcal{L}_1 = - \sum_y p_{old}(y | x) \log p_\theta(y | x) \quad (4)$$

This loss promotes higher log-probabilities for responses to which the old model assigned higher confidence. It does not require an external reward model or labels, relying solely on the model’s own confidence as feedback. This enables the use of cheaper, unlabeled data.

We also generalize this to a broader class of differentiable functions  $\mathcal{L}(p_{old}, p_{\theta})$ . An effective variant smooths the weighting by adding a constant  $\alpha > 0$ :

$$\mathcal{L}_2 = - \sum_y (p_{old}(y | x) + \alpha) \log p_{\theta}(y | x) \quad (5)$$

We empirically find that even small values of  $\alpha$  (for example, 0.1) improve both convergence and generalization.

Table 1: Loss functions and corresponding optimized functionals

Name	Loss function	Functional
RLSC loss	$p_{old} \log p$	$\mathbb{E}_{p_{\theta}}[p_{\theta}]$
Shannon Entropy	$(1 + \log p_{old}) \log p$	$\mathbb{E}_{p_{\theta}}[\log p_{\theta}]$
Advantage Score	$A(y) \log p$	$\mathbb{E}_{p_{\theta}}[A(y)]$

Some alternative estimation methods and their corresponding loss functions are also proposed, as shown in Table 1.

These statements capture the core idea and motivation behind the RLSC method, with its algorithm presented in pseudocode in Algorithm 1.

---

**Algorithm 1** The pseudo-code of the RLSC for LLM

---

```

# model.generate(prompt): generates multiple completions
# model.forward(input): returns token logits

for question in dataset:
    # generate completions
    completions = model.generate(question, temperature, num_samples)

    # get gradable probabilities
    logits = model.forward(question.repeat() + completions)[question.
        length:-1]
    all_log_probs = log_softmax(logits / temperature)
    log_p = all_log_probs.gather(token_ids).sum

    # compute loss
    loss = - (exp(log_p).detach() + alpha) * log_p

    loss.backward()
    optimizer.step()

```

---

### 3 EXPERIMENTS

#### 3.1 TRAINING SETUP

We leverage the self-confidence signal as a training supervision mechanism to fine-tune a range of models, including Qwen Yang et al. (2024b), LLaMA Grattafiori et al. (2024), Gemma Team et al. (2024), Olmo Groeneveld et al. (2024), and distilled versions of Deepseek-R1 Guo et al. (2025), among others. The training procedure is as follows:

1. For each question, call generate(temperature = 0.6, num of samples = n) to produce n candidate responses (with n = 1 to 8);
2. For each (prompt + response) pair, tokenize the sequence and compute token-level log-probabilities. Then, resample the tokens with a temperature of 1 to ensure the original distribution remains unchanged;
3. Apply an assistant mask to keep only the response tokens, and sum the log-probabilities across the masked region to obtain the sequence-level response likelihood;
4. Compute the loss based on the sequence-level log-likelihood and update model parameters via backpropagation using the Adam optimizer.

**Training Dataset.** We performed reinforcement learning-based training on the GSM8K Cobbe et al. (2021), AIME24 Li et al. (2024b), AMC23 Li et al. (2024a) and MATH500 Hendrycks et al. (2021c) datasets.

### 3.2 EVALUATION AND RESULTS

**Benchmarks.** We evaluated our method on several challenging benchmark datasets, including mathematical reasoning tasks such as AIME24 Li et al. (2024b), MATH500 Hendrycks et al. (2021c), AMC23 Li et al. (2024a), MMLU STEM Hendrycks et al. (2021a), Minerva Math Hendrycks et al. (2021b), Olympiadbench He et al. (2024).

**Evaluation Metrics.** For model evaluation, we adopt the Pass@1 metric as proposed in DeepSeek-R1 Guo et al. (2025), which measures the average accuracy under high-temperature and multi-sample generation settings, as defined in Equation 6.

**Evaluation Metrics.** For model evaluation, we adopt the **Pass@1** metric as proposed in DeepSeek-R1 Guo et al. (2025), which measures the average accuracy under high-temperature and multi-sample generation settings, as defined in Equation 6. All evaluations are conducted using a generation setting of 3072 new tokens, with **temperature** = 0.6 and **top-p** = 0.95.

$$\text{Pass@1} = \frac{1}{k} \sum_{i=1}^k p_i \tag{6}$$

To ensure a fair comparison, we evaluated both our trained model checkpoints and the baseline using the same evaluation template, with all experimental settings kept identical, detailed evaluation procedures can be found in Appendix B.2.

Table 2: Performance comparison across different benchmarks (Pass@1). All compared models have the same size (7B). For each benchmark, the **best** result is highlighted in bold, the second best is underlined, and the *third best* is italicized, \* indicates the best-performing result among models trained without labeled data. Eurur-2-7B-SFT is developed based on Qwen2.5-Math.

Model	Base	Reward	AIME24	MATH500	AMC23	MMLU	Minerva	Olympiad	Avg
<i>With using extra labeled data.</i>									
Qwen2.5-Math-7B-Instruct Yang et al. (2024b)	*	SFT+RL	8.8	<b>81.6</b>	58.8	<b>70.3</b>	38.8	37.6	49.3
Open-Reasoner-Zero Hu et al. (2025)	Qwen2.5	Rule-Based	13.3	79.6	51.3	<u>68.1</u>	38.2	<b>41.2</b>	48.6
SimpleRL-Zoo Zeng et al. (2025)	Qwen2.5-Math	RLVR	<u>28.3</u>	74.9	59.4	56.1	23.0	34.2	46.0
Oat-Zero Liu et al. (2025)	Qwen2.5-Math	RLVR	<b>31.7</b>	79.1	<b>68.8</b>	40.6	<u>40.4</u>	<u>39.9</u>	<i>50.1</i>
Qwen2.5-Math-7B-GRPO*	Qwen2.5-Math	RLVR	23.3	78.3	<u>67.5</u>	56.9	<b>40.9</b>	39.3	<u>51.0</u>
<i>Without using any extra labeled data.</i>									
Qwen2.5-Math-7B (Base Line) Yang et al. (2024b)	*	*	10.8	46.7	34.4	47.6	10.4	14.7	27.4
Spurious Reward Shao et al. (2025)	Qwen2.5-Math	Spurious Reward	10.8	50.9	41.9	48.4	9.7	17.8	29.9
TTRL Zuo et al. (2025)	Qwen2.5-Math	Majority Voting	18.3*	74.4	58.8	54.8	36.9	37.0	46.7
AZR Zhao et al. (2025)	Qwen2.5	Self Reward	11.7	64.1	39.4	62.9	18.3	29.5	37.6
PRIME-Zero Cui et al. (2025)	Qwen2.5	Implicit Reward	16.7	55.1	45.0	48.9	10.5	19.5	32.6
PRIME Cui et al. (2025)	Eurus-2-7B-SFT	Implicit Reward	13.3	72.6	61.3	41.5	39.5	34.4	43.8
<b>RLSC (Ours)</b>	Qwen2.5-Math	Self-Confidence	17.5	79.0*	<u>67.5*</u>	<u>63.1*</u>	<u>40.2*</u>	<u>39.4*</u>	<b>51.1*</b>

We evaluated the above models using the official Qwen Math evaluation scripts Yang et al. (2024a). Qwen2.5-Math-7B-GRPO\* is our model, trained for 1 epoch on the Deepscaler Luo et al. (2025) dataset with the GRPO and RLVR mechanisms. Except for the Spurious-Reward Shao et al. (2025) and TTRL Zuo et al. (2025) models, which we re-implemented based on the scripts released with the Spurious-Reward paper (see training details in Appendix B.1), all other models were evaluated

using trained weights from HuggingFace. The evaluation parameters were kept consistent across all models.

As shown in Table 2, our method achieves the highest average performance across five mathematical benchmarks and one benchmark encompassing Science, Technology, Engineering, and Mathematics (STEM), including some models trained with labeled data. Notably, among models trained without labels, our approach outperforms existing methods on the MATH500, AMC23, MMLU Stem, Minerva, and Olympiad benchmarks.

Regarding training efficiency, most other models typically require at least five hours to train, with TTRL being particularly time-consuming—training on an 8xA100 setup took 15 hours due to its Majority Voting procedure (64 responses per prompt). In contrast, our method requires only 15–30 training steps, completing in approximately 30 minutes on the same 8xA100 GPU setup.

Table 3: Pass@1 performance comparison across different models. Evaluation was conducted with non-zero temperature sampling, averaging 4 responses per question.

Model	AIME24	MATH500	AMC23	MMLU	Minerva	Olympiad	Avg
<i>LLaMA-8B</i>							
LLaMA-8B Grattafiori et al. (2024)	0.0	12.0	1.9	45.6	<b>6.4</b>	3.1	11.5
Ours	<b>0.0</b>	<b>12.3</b>	<b>7.5</b>	<b>46.9</b>	6.0	<b>3.5</b>	<b>12.7</b>
$\Delta$	<b>+0.0</b>	<b>+0.3</b>	<b>+5.6</b>	<b>+1.3</b>	<b>-0.4</b>	<b>+0.4</b>	<b>+1.2</b>
<i>LLaMA-8B-It</i>							
LLaMA-8B-It Grattafiori et al. (2024)	3.3	49.5	19.4	49.4	19.1	<b>17.6</b>	26.38
Ours	<b>8.3</b>	<b>49.6</b>	<b>23.1</b>	<b>50.4</b>	<b>24.0</b>	17.1	<b>28.75</b>
$\Delta$	<b>+5.0</b>	<b>+0.1</b>	<b>+3.7</b>	<b>+1.0</b>	<b>+4.9</b>	<b>-0.5</b>	<b>+2.37</b>
<i>OLMo-2-7B</i>							
OLMo-2-7B Groeneveld et al. (2024)	3.3	9.7	5.0	40.9	2.2	2.3	10.57
Ours	<b>3.3</b>	<b>19.5</b>	<b>17.5</b>	<b>45.9</b>	<b>4.2</b>	<b>4.6</b>	<b>15.83</b>
$\Delta$	<b>+0.0</b>	<b>+9.8</b>	<b>+12.5</b>	<b>+5.0</b>	<b>+2.0</b>	<b>+2.3</b>	<b>+5.26</b>
<i>Gemma-4B-Pt</i>							
Gemma-4B-pt Team et al. (2024)	0.0	7.4	3.8	36.4	2.8	2.2	8.77
Ours	<b>0.8</b>	<b>14.1</b>	<b>8.6</b>	<b>40.0</b>	<b>5.1</b>	<b>3.7</b>	<b>12.05</b>
$\Delta$	<b>+0.8</b>	<b>+6.7</b>	<b>+4.8</b>	<b>+3.6</b>	<b>+2.3</b>	<b>+1.5</b>	<b>+3.28</b>
<i>Qwen-Math-1.5B</i>							
Qwen-Math-1.5B Yang et al. (2024b)	5.0	34.8	30.0	33.3	8.7	19.9	21.95
Ours	<b>6.7</b>	<b>66.5</b>	<b>43.1</b>	<b>48.1</b>	<b>26.2</b>	<b>28.7</b>	<b>36.55</b>
$\Delta$	<b>+1.7</b>	<b>+31.7</b>	<b>+13.1</b>	<b>+14.8</b>	<b>+17.5</b>	<b>+8.8</b>	<b>+14.6</b>
<i>Qwen-Math-7B</i>							
Qwen-Math-7B Yang et al. (2024b)	10.8	46.7	34.4	47.6	10.4	14.7	27.44
Ours	<b>17.5</b>	<b>79.0</b>	<b>67.5</b>	<b>63.1</b>	<b>40.2</b>	<b>39.4</b>	<b>51.12</b>
$\Delta$	<b>+6.7</b>	<b>+32.3</b>	<b>+33.1</b>	<b>+15.5</b>	<b>+29.8</b>	<b>+24.7</b>	<b>+23.68</b>

The results are shown in Table 3. In multiple benchmark tests, all models achieved significant improvements, with particularly outstanding performance in the Qwen2.5-Math-7B (an average improvement of 23.68%). Even when the initial model performed poorly, effective performance gains were still realized.

### 3.3 ANALYSIS OF EMERGENT BEHAVIOR

**Impact of Training Datasets.** We found that different training datasets lead to significant variations in model performance. Through analysis, we believe that simpler datasets are more effective in enhancing the balance between the model’s confidence and accuracy. As shown in Figure 2, we used greedy decoding and report only one generated response per prompt.

**More Confident Responses.** We observed that fine-tuning based on RLSC enables the model to produce answers that are more concise and confident. Unlike traditional fine-tuning methods, our model learns to identify the correct answer early, thereby avoiding lengthy and redundant reasoning.

**Analysis of Accuracy Improvements Across Different Models.** Experimental results demonstrate that the proposed method consistently improves performance across multiple models, yet a notable discrepancy arises between Qwen and LLaMA. Specifically, Qwen exhibits relatively low policy entropy at the early stage of training, indicating higher decision determinism, which enables the

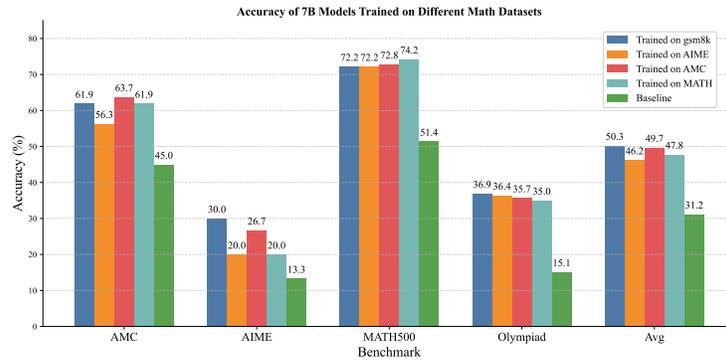


Figure 2: Comparison of the performance of models trained on various datasets, alongside the baseline.

sharpened distribution to enhance performance more rapidly. In contrast, LLaMA starts with significantly higher policy entropy, resulting in lower certainty. By further tracking the training process, we observe that *LLaMA tends to improve its confidence by shortening the generated sequence length*, thereby gradually achieving convergence. This phenomenon suggests that the distinct policy entropy characteristics at the beginning of training may be a key factor underlying the differences in enhancement effectiveness across models.

**Visualization of Logits Distribution.** We visualized the logits distribution during the sampling process before and after training as shown in Figure 3. The results show that the variance of the distribution is relatively small before training, while it becomes significantly larger after training. Meanwhile, the red dashed line in the figure indicates the position of the maximum value. It can be clearly observed that the model becomes more confident in its predictions after training.

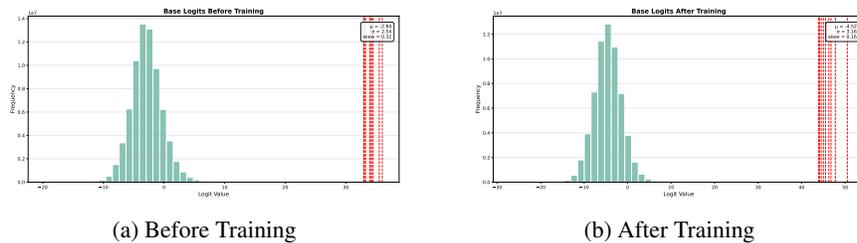


Figure 3: Logits distribution before and after training. Red dashed lines indicate the top 10 maximum values.

**Training Efficiency: Model Average Accuracy Across All Steps.** Since our model is trained using an on-policy approach and gradients are computed based on entire sequences level rather than tokens level, it converges relatively quickly. Moreover, with 4 rollouts per prompt, the model achieves rapid convergence. As shown in Figure 4, Qwen2.5-Math-1.5B converges in 20 steps, while in the same configuration, other models are also able to converge in 15 to 30 steps.

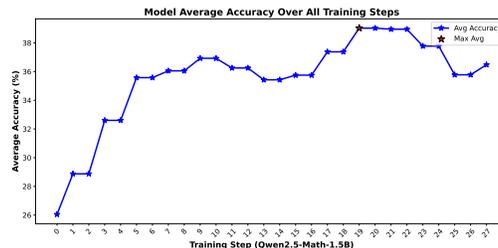


Figure 4: Average performance during training over 4 benchmarks.

**Confidence and Entropy.** As shown in Figure 5, the model’s average confidence steadily increases during training, while the entropy decreases accordingly. This visualization indicates that the objective of sharpening the distribution has been successfully achieved.

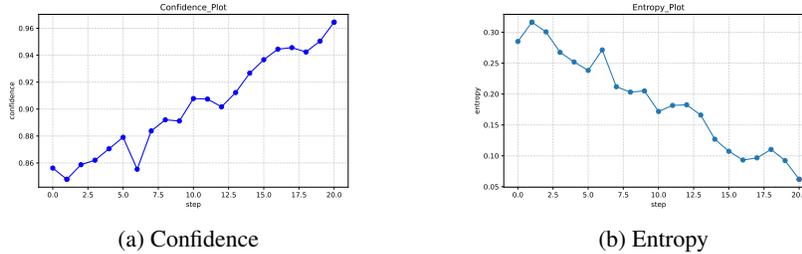


Figure 5: Variation of Average Confidence and Entropy Across Training Steps (Train on gsm8k).

**Cherry Pick.** Additionally, we present the outcomes of the two cases on AMC 2023 and AIME 2024, providing a clear, step-by-step exposition of the solution strategies.

### Case 1: Meeting Point of Two Bikers

**Problem.** Cities  $A$  and  $B$  are 45 miles apart. Alicia lives in  $A$  and Beth lives in  $B$ . Alicia bikes towards  $B$  at 18 miles per hour. Leaving at the same time, Beth bikes toward  $A$  at 12 miles per hour. How many miles from City  $A$  will they be when they meet? Mark your solution with `\boxed`.

**Solution.** To determine how many miles from City  $A$  Alicia and Beth will be when they meet, follow these steps:

1. **Combined speed:** Alicia bikes at 18 mph, Beth at 12 mph.

$$\text{Combined speed} = 18 + 12 = 30 \text{ mph}$$

2. **Time to meet:** The distance between the cities is 45 miles.

$$t = \frac{\text{Distance}}{\text{Combined speed}} = \frac{45}{30} = 1.5 \text{ hours}$$

3. **Distance Alicia travels:**

$$\text{Distance} = \text{Speed} \times \text{Time} = 18 \times 1.5 = 27 \text{ miles}$$

**Answer:** They will meet 27 miles from City  $A$ .

`\boxed{27}`

### Case 2: Finding $xy$ from Logarithmic Equations

**Problem.** There exist real numbers  $x$  and  $y$ , both greater than 1, such that

$$\log_x(y^x) = \log_y(x^{4y}) = 10.$$

Find  $xy$ . Mark your solution with `\boxed`.

**Solution.** Solve step by step.

1. Rewrite the logarithmic equations using logarithm properties:

$$\log_x(y^x) = x \log_x(y) = 10, \quad \log_y(x^{4y}) = 4y \log_y(x) = 10.$$

2. Let  $\log_x(y) = a$  and  $\log_y(x) = b$ :

$$x \cdot a = 10 \implies x = \frac{10}{a}, \quad 4y \cdot b = 10 \implies y = \frac{5}{2b}.$$

3. Using the change-of-base formula:

$$a = \frac{\log(y)}{\log(x)}, \quad b = \frac{\log(x)}{\log(y)}.$$

4. Solving numerically gives

$$xy = 25.$$

**Answer:**

`\boxed{25}`

## 4 RELATED WORK

**Reinforcement Learning in Reasoning Tasks.** In recent years, **Reinforcement Learning (RL)** has played a pivotal role in enhancing the reasoning capabilities of Large Language Models (LLMs). Models such as DeepSeek-R1 [Guo et al. \(2025\)](#), ChatGPT o1 [Achiam et al. \(2023\)](#), QwQ [Team \(2025\)](#), and Qwen have demonstrated impressive reasoning skills by decomposing complex problems into intermediate steps and engaging in deep deliberation prior to producing final responses capabilities often acquired and refined through reward-driven learning mechanisms.

A widely adopted paradigm is Reinforcement Learning from Human Feedback (RLHF) [Ouyang et al. \(2022\)](#), which aligns model behavior with human preferences using human annotations or learned preference models as reward signals. RLHF typically employs Proximal Policy Optimization (PPO) [Schulman et al. \(2017\)](#) for policy updates. However, RLHF is heavily reliant on costly human annotations, which limits scalability.

To alleviate this dependency, Reinforcement Learning with Verifiable Rewards (RLVR) [Luong et al. \(2024\)](#) [Lambert et al. \(2024\)](#) [Zeng et al. \(2025\)](#) [Liu et al. \(2025\)](#) adopts a paradigm where rewards are computed directly from question-answer pairs  $(x, y^*)$ , by comparing model outputs with reference answers. By leveraging such a scheme, RLVR greatly reduces the resource requirements of training, though it still depends on human-labeled questions and answers, which constrains scalability.

**Test-Time Training.** More recently, **Test-Time Training (TTT)** [Zuo et al. \(2025\)](#) has emerged as a promising direction for further optimizing model behavior during inference. Notable examples include *SelfPlay Critic (SPC)* [Chen et al. \(2025\)](#) and *Absolute Zero Reasoner (AZR)* [Zhao et al. \(2025\)](#), which employ adversarial dual-model frameworks inspired by game-theoretic learning. In these approaches, one model acts as a "sneaky generator" crafting challenging or misleading reasoning steps, while the other serves as a "critic" that learns to detect errors. These methods eliminate the need for human supervision but depend on external tools (e.g., Python executors or code verifiers) to supply feedback signals.

Another prominent TTT approach is TTRL [Zuo et al. \(2025\)](#). It constructs pseudo-labels by sampling multiple candidate responses for each question and applying a majority-vote mechanism. The resulting consensus serves as a proxy label to compute rewards for model updates. Although TTRL avoids explicit human supervision, it requires a large number of samples (e.g., 64 per question), which leads to significant computational overhead.

**Summary and Motivation.** The reliance on annotations, external models, or hand-crafted rewards underscores the need for a simple and efficient label-free reinforcement signal that can be derived directly from the model itself.

## 5 CONCLUSION

We introduced **Reinforcement Learning via Self-Confidence (RLSC)**, a lightweight fine-tuning method that eliminates the need for labels, preference models, or handcrafted rewards. Unlike prior approaches such as TTRL, which rely on large-scale majority voting, RLSC formalizes the underlying principle mathematically.

Our key contribution is a derivation: we show that majority voting implicitly optimizes for agreement within the model’s output distribution - and we transform that into a differentiable, self-supervised objective. This “mode sharpening” functional enables efficient reinforcement learning directly from the model’s own confidence.

We find that for nearly all pretrained models, RLSC converges reliably within 4–8 rollouts per question and 15–30 training steps, while still delivering substantial accuracy improvements without external supervision.

This work demonstrates that high-quality post-training can emerge not from external labels, but from a model’s internal signal - when that signal is derived with care. We believe RLSC provides both a practical tool and a conceptual bridge between majority-voting-based pseudo-labeling and principled self-supervision.

## 486 ETHICS STATEMENT

487

488 This work strictly adheres to the ICLR Code of Ethics. The research does not involve human sub-  
 489 jects, nor does it pose potential harm to society, individuals, or the environment. There are no  
 490 conflicts of interest, nor any issues related to discrimination, bias, or fairness. All experiments and  
 491 analyses were conducted following academic norms and ethical guidelines.

492

## 493 REPRODUCIBILITY STATEMENT

494 All content reported in this paper is accurate. Detailed information on the evaluation and training  
 495 methods can be found in the Supplementary Materials to facilitate reproducibility. The Supple-  
 496 mentary Materials include data processing steps, model, training hyperparameters, and evaluation  
 497 metrics.

498

## 499 REFERENCES

500

501 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
 502 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
 503 report. *arXiv preprint arXiv:2303.08774*, 2023.

504 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
 505 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

506

507 Jiaqi Chen, Bang Zhang, Ruotian Ma, Peisong Wang, Xiaodan Liang, Zhaopeng Tu, Xiaolong Li,  
 508 and Kwan-Yee K Wong. Spc: Evolving self-play critic via adversarial games for llm reasoning.  
 509 *arXiv preprint arXiv:2504.19162*, 2025.

510 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
 511 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
 512 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

513 Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu  
 514 Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint*  
 515 *arXiv:2502.01456*, 2025.

516

517 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
 518 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd  
 519 of models. *arXiv preprint arXiv:2407.21783*, 2024.

520 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
 521 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the  
 522 science of language models. *arXiv preprint arXiv:2402.00838*, 2024.

523

524 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
 525 Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—  
 526 the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.

527 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
 528 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
 529 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

530 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
 531 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for  
 532 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*  
 533 *arXiv:2402.14008*, 2024.

534

535 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
 536 Steinhardt. Measuring massive multitask language understanding. *Proceedings of the Interna-*  
 537 *tional Conference on Learning Representations (ICLR)*, 2021a.

538 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
 539 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,  
 2021b.

- 540 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
541 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
542 *preprint arXiv:2103.03874*, 2021c.
- 543
- 544 Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum.  
545 Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base  
546 model. *arXiv preprint arXiv:2503.24290*, 2025.
- 547 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brah-  
548 man, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. T\”ulu 3: Pushing  
549 frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- 550
- 551 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
552 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in  
553 ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*,  
554 13:9, 2024a.
- 555 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
556 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in  
557 ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*,  
558 13:9, 2024b.
- 559 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
560 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
561 *arXiv:2412.19437*, 2024.
- 562
- 563 Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee,  
564 and Min Lin. Understanding rl-zero-like training: A critical perspective. *arXiv preprint*  
565 *arXiv:2503.20783*, 2025.
- 566 Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y Tang, Manan Roongta, Colin Cai,  
567 Jeffrey Luo, Tianjun Zhang, Li Erran Li, et al. Deepscaler: Surpassing o1-preview with a 1.5 b  
568 model by scaling rl. *Notion Blog*, 2025.
- 569
- 570 Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Rea-  
571 soning with reinforced fine-tuning. *arXiv preprint arXiv:2401.08967*, 3, 2024.
- 572 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
573 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
574 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
575 27730–27744, 2022.
- 576
- 577 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
578 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
579 *in Neural Information Processing Systems*, 36:53728–53741, 2023.
- 580 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
581 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 582
- 583 Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei  
584 Du, Nathan Lambert, Sewon Min, Ranjay Krishna, et al. Spurious rewards: Rethinking training  
585 signals in rlvr. *arXiv preprint arXiv:2506.10947*, 2025.
- 586 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-  
587 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma  
588 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- 589
- 590 Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL  
591 <https://qwenlm.github.io/blog/qwq-32b/>.
- 592 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
593 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*  
*arXiv:2407.10671*, 2024a.

594 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-  
595 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical  
596 expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024b.

597  
598 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-  
599 zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv*  
600 *preprint arXiv:2503.18892*, 2025.

601 Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin Xu, Matthieu Lin, Shenzhi Wang, Qingyun  
602 Wu, Zilong Zheng, and Gao Huang. Absolute zero: Reinforced self-play reasoning with zero  
603 data. *arXiv preprint arXiv:2505.03335*, 2025.

604 Yuxin Zuo, Kaiyan Zhang, Shang Qu, Li Sheng, Xuekai Zhu, Biqing Qi, Youbang Sun, Ganqu  
605 Cui, Ning Ding, and Bowen Zhou. Ttrl: Test-time reinforcement learning. *arXiv preprint*  
606 *arXiv:2504.16084*, 2025.

## 609 A USE OF LARGE LANGUAGE MODELS (LLMs)

610  
611 We used a large language model (LLM) to assist with grammar correction and language polishing.  
612 All scientific content, analysis, and results were generated by the authors, and the LLM was only  
613 used for improving readability and clarity.

## 615 B CONFIGURATION

### 617 B.1 TRAINING CONFIGURATION

618  
619 All experiments were performed on a cluster equipped with 8xNVIDIA A100 GPUs, ensuring a  
620 consistent hardware setup across all training runs.

#### 622 B.1.1 RLSC

623 We train the RLSC model with the following settings:

- 625 • **Batch size:** 1
- 626 • **Generation per question:** 1–8 responses (Rollouts)
- 627 • **Top-p sampling:** 0.95
- 628 • **Temperature:** 0.95
- 629 • **Repetition penalty:** 1.05
- 630 • **Total training steps:** 30
- 631 • **Generated sequence length:** 512 or 768 or 1024 tokens
- 632 • **Gradient accumulation:** 1 or 2
- 633 • **Optimizer:** AdamW
- 634 • **Learning rate:**  $5 \times 10^{-6}$
- 635 • **Reward:** RLSC

636  
637  
638 All hyperparameters were kept consistent across experiments, unless otherwise specified.

#### 641 B.1.2 QWEN2.5-MATH-7B-GRPO\*

642 We train the Qwen2.5-Math-7B-GRPO\* model with the following settings:

- 643 • **Training dataset:** DeepScaleR
- 644 • **Train batch size:** 128
- 645 • **Max prompt length:** 1024 tokens

- 648 • **Max response length:** 3072 tokens
- 649 • **Generation per question:** 8 responses (rollouts)
- 650 • **Sampling strategy:** Temperature=1
- 651 • **Optimizer:** AdamW
- 652 • **Learning rate:**  $5 \times 10^{-6}$
- 653 • **Epochs:** 1
- 654 • **Reward:** RLVR

### 657 B.1.3 TTRL

658 We train the TTRL for  $\approx 13$  hours with the following settings:

- 661 • **Batch size:** 128
- 662 • **Micro batch size:** 4
- 663 • **Rollout batch size:** 64
- 664 • **Micro rollout batch size:** 4
- 665 • **Temperature:** 1.0
- 666 • **Total global steps:** 50
- 667 • **Generated sequence length:** 3072
- 668 • **Prompt max length:** 1024
- 669 • **N votes per prompt:** 64
- 670 • **Optimizer:** DeepSpeedCPUAdam
- 671 • **Actor learning rate:**  $5 \times 10^{-7}$
- 672 • **Critic learning rate:**  $9 \times 10^{-6}$
- 673 • **KL coefficient:** 0.00 (with KL loss)
- 674 • **Advantage estimator:** group\_norm
- 675 • **Discount factor  $\gamma$ :** 1.0
- 676 • **GAE parameter  $\lambda$ :** 1.0
- 677 • **Nodes/GPUs:** 1 node with 4 GPUs for actor, critic, ref (A100); vLLM with 4 engines,
- 678 tensor parallel size 1
- 679 • **Precision:** bfloat16
- 680 • **Reward:** Majority Voting

681 All hyperparameters were kept consistent across experiments, unless otherwise specified.

### 682 B.1.4 SPURIOUS REWARD

683 We train the Spurious Reward with the following settings:

- 691 • **Batch size:** 128
- 692 • **Micro batch size:** 4
- 693 • **Rollout batch size:** 64
- 694 • **Micro rollout batch size:** 4
- 695 • **Temperature:** 1.0
- 696 • **Total global steps:** 50
- 697 • **Generated sequence length:** 3072
- 698 • **Prompt max length:** 1024
- 699 • **Optimizer:** DeepSpeedCPUAdam

- 702 • **Actor learning rate:**  $5 \times 10^{-7}$
- 703 • **Critic learning rate:**  $9 \times 10^{-6}$
- 704 • **KL coefficient:** 0.00 (with KL loss)
- 705 • **Advantage estimator:** group\_norm
- 706 • **Discount factor  $\gamma$ :** 1.0
- 707 • **GAE parameter  $\lambda$ :** 1.0
- 708 • **Nodes/GPUs:** 1 node with 4 GPUs for actor, critic, ref (A100); vLLM with 4 engines,
- 709 tensor parallel size 1
- 710 • **Precision:** bfloat16
- 711 • **Reward:** Spurious Reward

712 All hyperparameters were kept consistent across experiments, unless otherwise specified.

## 713 B.2 EVALUATION CONFIGURATION

714 We evaluated the performance of all models using the official Qwen mathematical evaluation scripts.  
715 Except for the re-implemented models explicitly noted, all other models were tested directly with  
716 pre-trained weights from HuggingFace, with evaluation parameters kept consistent across all mod-  
717 els.

- 723 • **Generation per question:** 4
- 724 • **Temperature:** 0.6, unless otherwise specified.
- 725 • **Top-p sampling:** 0.95
- 726 • **Generated max new sequence length:** 3072, unless otherwise specified.
- 727 • **Prompt type:** Math CoT — Please reason step by step, and enclose your final answer  
728 within `\boxed{}`.

729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755