# THE REASONING BOUNDARY PARADOX: HOW RE-INFORCEMENT LEARNING CONSTRAINS LANGUAGE MODELS

Phuc Minh Nguyen<sup>1,2</sup>, Chinh D. La<sup>1</sup>, Duy M. H. Nguyen<sup>3,5,6</sup>, Nitesh V. Chawla<sup>2,4</sup>
Binh T. Nguyen<sup>1,2\*</sup>, Khoa D. Doan<sup>1,2\*</sup>

<sup>1</sup>College of Engineering and Computer Science, VinUniversity

<sup>2</sup>Center for Environmental Intelligence, VinUniversity

<sup>3</sup>University of Stuttgart, <sup>4</sup>University of Notre Dame

<sup>5</sup>German Research Center for Artificial Intelligence (DFKI)

<sup>6</sup>Max Planck Research School for Intelligent Systems (IMPRS-IS)

#### **ABSTRACT**

Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as a key method for improving Large Language Models' reasoning capabilities, yet recent evidence suggests it may paradoxically shrink the reasoning boundary rather than expand it. This paper investigates the shrinkage issue of RLVR by analyzing its learning dynamics and reveals two critical phenomena that explain this failure. First, we expose negative interference in RLVR, where learning to solve certain training problems actively reduces the likelihood of correct solutions for others, leading to the decline of Pass@k performance, or the probability of generating a correct solution within k attempts. Second, we uncover the winnertake-all phenomenon: RLVR disproportionately reinforces problems with high likelihood, correct solutions under the base model while suppressing other initially low-likelihood ones. Through extensive theoretical and empirical analysis on multiple mathematical reasoning benchmarks, we show that this effect arises from the inherent on-policy sampling in standard RL objectives, causing the model to converge toward narrow solution strategies. Based on these insights, we propose a simple yet effective data curation algorithm that focuses RLVR learning on low-likelihood problems, achieving notable improvement in Pass@k performance. Our code is available at https://github.com/mail-research/ SELF-llm-interference.

#### 1 Introduction

Large Language Models (LLMs) have recently shown remarkable capabilities in complex logical tasks such as mathematical reasoning (Cobbe et al., 2021; Hendrycks et al., 2021) and programming (Jimenez et al., 2024; Yang et al., 2024b). A key factor behind this success is *Reinforcement Learning with Verifiable Rewards* (RLVR, Lambert et al. 2025; DeepSeek-AI et al. 2025). RLVR optimizes the LLMs against a binary signal based on objective correctness, eliminating the need for human annotations (OpenAI et al., 2024; MetaAI, 2024b). It is believed that RLVR encourages the emergence of novel reasoning strategies, such as self-reflection and iterative refinement (DeepSeek-AI et al., 2025; Zeng et al., 2025; Luo et al., 2025), enabling the LLMs to surpass the capabilities of their base models.

However, recent studies suggest that RLVR training does not expand the reasoning boundary beyond what the base model already possesses (Yue et al., 2025; Zhao et al., 2025; Zhu et al., 2025; Liu et al., 2025b). Notably, Liu et al. (2025c); Zhao et al. (2025) reveal that the base model already exhibits complex reasoning behaviors even before RLVR, while Yue et al. (2025), using the pass@k metric (i.e., the probability of generating a correct solution within k attempts), demonstrate that RLVR can even shrink the reasoning boundary, reducing the set of problems the model can solve within k

<sup>\*</sup>Co-Corresponding Authors: binh.nt2@vinuni.edu.vn & khoa.dd@vinuni.edu.vn

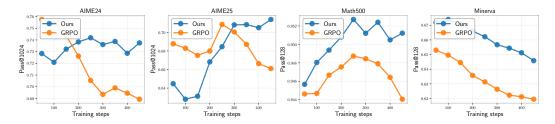


Figure 1: Pass@k evolution (smoothed) during RLVR training with Qwen2.5-Math-1.5B, comparing our proposed finetuning objective SELF to GRPO under a large sampling budget k. While GRPO shows a progressive decline, SELF exhibits consistent improvements in Pass@k throughout the training process.

trials. Identifying why this coverage shrinkage occurs is crucial for understanding and effectively leveraging RLVR to solve novel problems.

To investigate this issue, this paper analyzes the learning dynamics of RLVR training. Unlike standard RL, which is effective in learning novel strategies within a single and well-defined Markov Decision Process (MDP), in LLMs reasoning, each problem x induces its own MDP with a distinct and unknown reward function  $r(x, \cdot)$  (Setlur et al., 2025; Qu et al., 2025). Consequently, learning to solve a problem x, defined by one MDP, can affect the ability of the LM to solve another problem x', defined by another MDP. Indeed, our analysis reveals that RLVR is prone to a negative interference effect where learning to solve a subset of training problems reduces the likelihood of generating correct solutions for others. This dynamic can lead to a winner-take-all scenario (Schaul et al., 2019), where the LLMs disproportionately improve on a limited subset of problems that the base model can already solve with high likelihood while neglecting or even regressing on others, ultimately leading to a decline in Pass@k performance or problem coverage. Furthermore, we highlight the failure of current regularization techniques, such as clipping (Schulman et al., 2017) and KL regularization in preventing these effects in RLVR. Finally, these findings motivate us to propose SELF, a simple yet effective data curation algorithm that focuses model learning on a subset of problems with low likelihood, effectively improving the coverage performance across mathematical benchmarks. In summary, the main contributions of this work are as follows:

- 1. We study the learning dynamics of RLVR and reveal *negative interference* in LLM reasoning, where learning to solve a subset of training problems results in a negative effect on others, leading to *winner-take-all* in which the model only concentrates on solving a smaller set of problems.
- 2. We reveal the *winner-take-all* learning phenomenon, where RLVR tends to reinforce problems highly solvable under the base model while neglecting problems with initially low likelihood of correct solutions. Due to the existence of *negative interference*, these problems with lower success rates will progressively have lower likelihoods of generating the correct solutions, and on-policy learning cannot provide any explicit learning signal to improve them. This explains the reduction in the diversity of previously learned behaviors during RLVR training or reasoning boundary shrinkage in RLVR.
- 3. We propose SELF (Selective Examples with Low-likelihood and Forward-KL), a data curation algorithm that selectively solves only problems with low likelihood of arriving at correct answers while preserving previously learned behavior during RLVR. Our extensive empirical evaluation demonstrates that SELF not only improves sample efficiency but also effectively mitigates the coverage shrinkage problem in RLVR.

## 2 RELATED WORKS

Reinforcement Learning with Verifiable Rewards (RLVR). RLVR has emerged as a powerful method for improving LLMs in mathematical reasoning and programming (DeepSeek-AI et al., 2025). However, recent evidence suggests that reasoning abilities already exist in base models, with RLVR only amplifying rather than creating such capabilities (Shao et al., 2025; Yue et al., 2025; Liu et al., 2025c). Moreover, RLVR can degrade performance and reduce coverage of correct solutions, evidenced by decreased pass@k metrics (Yue et al., 2025; Dang et al., 2025; Zhao et al., 2025;

Wu et al., 2025). Liu et al. (2025a) attributes coverage reduction to limited RL training and over-specialization during pretraining, while Zhu et al. (2025) shows that increasing the likelihood of correct solutions improves accuracy but reduces diversity.

Loss of Plasticity in Neural Networks. Neural networks gradually lose adaptability, a phenomenon called *plasticity loss* (Klein et al., 2024; Juliani & Ash, 2024), particularly pronounced in RL due to its non-stationary nature (Klein et al., 2024; Juliani & Ash, 2024; Moalla et al., 2024). Proposed remedies include resetting parameters (Ash & Adams, 2020; Kielo & Lukin, 2024; D'Oro et al., 2023) and self-distillation (Igl et al., 2021). In particular, Tang & Berseth (2024); Tang et al. (2025) show plasticity loss correlates with drastic model confidence changes (churn), which undermines regularization techniques like clipping (Moalla et al., 2024). Our work explains why RLVR reduces behavioral diversity in LLM reasoning, where non-stationarity arises from both shifting distributions and varying prompt-induced objectives.

**Learning Dynamics.** Ren & Sutherland (2025) shows negative gradients in off-policy LLM fine-tuning create a *squeezing effect*. In RL, *interference* occurs when learning degrades performance on unseen or previously learned states (Liu et al., 2020; 2023). Schaul et al. (2019) demonstrates a *winner-take-all* effect (Schaul et al., 2019; Guo et al., 2018) in contextual bandits, where policies excel in some contexts while regressing in others. Our findings reveal *negative interference* as the key mechanism through which RLVR reduces solvable problems, with on-policy sampling causing disproportionate reinforcement of initially successful problems.

#### 3 PROBLEM SETTING AND PER-STEP INFLUENCE IN RLVR

**Language Model.** We first consider a language model (LM) policy  $\pi$ . For a given prompt x, the LM policy  $\pi$  will generate a response y in an auto-regressive manner:  $\pi(y|x) = \prod_t \pi(y_t|x, y_{< t})$ , where  $y_t$  is the t-th indexed token in y and  $y_{< t}$  is the partial completion before  $y_t$ .

Reinforcement Learning with Verifiable Rewards (RLVR). Given a verifiable reward function  $r(\boldsymbol{x}, \boldsymbol{y}) \in \{0, 1\}$  indicating whether a response  $\boldsymbol{y}$  is correct  $(\boldsymbol{y}^+)$  or incorrect  $(\boldsymbol{y}^-)$ , the goal of RLVR is to learn a new policy  $\pi_{\theta}$ , initialized from a base model (usually the pre-trained model)  $\pi_b$ , to generate responses maximizing the reward function with the following reward regularized maximization objective:

$$\max_{\pi_{\theta}} \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}, \boldsymbol{y} \sim \pi_{\theta}(\cdot | \boldsymbol{x})} \left[ r(\boldsymbol{x}, \boldsymbol{y}) - \beta \text{KL} \left( \pi_{\theta}(\cdot | \boldsymbol{x}) \parallel \pi_{b}(\cdot | \boldsymbol{x}) \right) \right], \tag{1}$$

where  $\beta$  is the KL regularization parameter. Eq. (1) represents a simplified version of common RL algorithms (PPO (Schulman et al., 2017), GRPO (DeepSeek-AI et al., 2025)). We provide a detailed derivation of the GRPO and PPO objectives in Appendix B.

**Per-step Influence.** Learning dynamics describes how the change in the model parameters  $\theta$ , induced by problem x, influences or changes the log-likelihood  $\log \pi_{\theta}$  of other problem-solution pairs (x', y'), which can be either training or test examples. Given an update step in LM parameters  $\theta^t \to \theta^{t+1}$ , the per-step influence measures the change in model confidence  $\log \pi_{\theta}(y'|x')$  on (x', y'):

$$\Delta \log \pi_{\theta}^{t}(\boldsymbol{y}'|\boldsymbol{x}') = \log \pi_{\theta^{t+1}}(\boldsymbol{y}'|\boldsymbol{x}') - \log \pi_{\theta^{t}}(\boldsymbol{y}'|\boldsymbol{x}')$$
 (2)

Using Taylor expansion, assuming with a sufficiently small learning rate  $\eta$ , we can decompose the per-step influence function as:

**Proposition 3.1.** Consider a problem-solutions pair (x, y) for updating using policy gradient objective, a sufficiently small learning rate  $\eta$ , the per-step influence on a problem-solution pair (x', y') is defined as:

$$\Delta \log \pi_{\theta}^{t}(\mathbf{y}'|\mathbf{x}') = \eta \mathcal{K}^{t}(\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}') A(\mathbf{x}, \mathbf{y}) + \mathcal{O}(\eta^{2})$$
(3)

where  $\mathcal{K}^t(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}') = \nabla_{\theta} \log \pi_{\theta^t}(\boldsymbol{y}|\boldsymbol{x})^{\top} \nabla_{\theta} \log \pi_{\theta^t}(\boldsymbol{y}'|\boldsymbol{x}')$  measures the influence between  $(\boldsymbol{x}, \boldsymbol{y})$  and  $(\boldsymbol{x}', \boldsymbol{y}')$ , and the advantage function  $A(\boldsymbol{x}, \boldsymbol{y})$  determines the direction of the model's update.

The dot product between two gradients,  $\mathcal{K}^t(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}')$ , can be interpreted as a model-specific similarity measure between problem-solution pairs, and  $A(\boldsymbol{x}, \boldsymbol{y})$  determines the *magnitude* and *direction* of the model update, indicating whether the gradient contributes positively or negatively.

#### 4 LEARNING DYNAMICS OF RLVR

Calculating  $\mathcal{K}^t(x, x', y, y')$  for LLMs is extremely expensive due to their large parameter counts. However, from Proposition 3.1, we have

$$\frac{|\Delta \log \pi_{\theta}^{t}(\boldsymbol{y}'|\boldsymbol{x}')|}{\eta C} \leq |\mathcal{K}^{t}(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}')|, \tag{4}$$

where C is a constant bounding the advantage function that holds due to finite binary reward values. This bound indicates that we can use the per-step influence  $\Delta \log \pi_{\theta}^t(y'|x')$  to approximate  $\mathcal{K}^t(x,x',y,y')$ . Given this insight, we introduce two quantitative metrics designed to better characterize the properties of  $\mathcal{K}(x,x',y,y')$  as follows.

**Definition 4.1** (Interference in Language Model). Given a behavior policy  $\mu$ , we define the interference of a language model policy,  $\pi_{\theta}^t$ , at iteration t as:

$$\Delta^{+}(\pi_{\theta^{t}}, \mu) = \mathbb{E}_{\substack{\boldsymbol{x}' \sim \mathcal{D}, \\ \boldsymbol{y}^{+} \sim \mu(\cdot|\boldsymbol{x})}} \left[ \Delta \log \pi_{\theta}^{t}(\boldsymbol{y}^{+}|\boldsymbol{x}') \right], \tag{5}$$

and the magnitude of the relative changes in model confidence on examples that lie outside of the training batch as:

$$||\Delta(\pi_{\theta^t}, \mu)|| = \mathbb{E}_{\substack{\boldsymbol{x}' \sim \mathcal{D}, \ \boldsymbol{y} \sim \mu(\cdot|\boldsymbol{x})}} \left(\Delta \log \pi_{\theta}^t(\boldsymbol{y}|\boldsymbol{x}')\right)^2.$$

Intuitively,  $\Delta^+(\pi_{\theta^t}, \mu)$  measures the *direction* of the relative influence  $\mathcal{K}(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}^+)$  on other problem–correct solution pairs under the behavior policy  $\mu$ . A positive  $\Delta \log \pi_{\theta}^t(\boldsymbol{y}^+|\boldsymbol{x})$ , indicates improvement, while a negative  $\Delta \log \pi_{\theta}^t(\boldsymbol{y}^+|\boldsymbol{x})$  corresponds to *negative interference*, i.e., learning to solve a training problem  $\boldsymbol{x}$  reduces the likelihood of producing other correct solutions. In our work, we consider  $\mu$  to be the base model  $\pi_b$  to study how RLVR affects the retention (or overwriting) of previously correct problem–solution pairs.

On the other hand,  $||\Delta(\pi_{\theta^t}, \mu)||$  measures the *squared magnitude* of the relative influence  $\mathcal{K}(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y})$  – i.e., how drastic the model confidence on examples generated from  $\mu$  changes after an update. If  $||\Delta(\pi_{\theta^t}, \pi_b)|| = 0$ , we can expect no coverage shrinkage to happen during RLVR training, as updating on training examples does not affect other non-training examples. An increase in the magnitude of influence  $||\Delta(\pi_{\theta^t}, \pi_b)||$  indicates that the language model struggles to distinguish between different problems, leading to highly correlated gradients across different examples. We provide a detailed explanation in Appendix C.

#### 4.1 EXPERIMENTAL SETUP

We focus on mathematical reasoning tasks to investigate the learning dynamics of RLVR. Specifically, we experiment on 2 series of models: Qwen2.5-Math-1.5B/7B (Yang et al., 2024a) and Llama-3.2-3B-Instruct (MetaAI, 2024b). We consider DeepScaleR (Luo et al., 2025) as the training dataset  $\mathcal{D}$ , which consists of approximately 40k mathematics problems collected from different sources. The detailed experimental setup is provided in Appendix A. We conduct experiments on 4 mathematical reasoning benchmarks: AIME24 & AIME25 (Art of Problem Solving), Math500 (Hendrycks et al., 2021; Lightman et al., 2024), and Minerva (Lewkowycz et al., 2022), following existing works on LLMs reasoning (Yue et al., 2025; Zhu et al., 2025; Li et al., 2025a). Unless otherwise specified, we use  $\mathcal{D}_{\text{test}}$  to denote the combined test set from all four benchmarks.

To analyze the model's learning dynamics during training, we create a probing dataset  $\mathcal{D}_{\text{prob}}$  (based on the training prompts  $x \in \mathcal{D}$ ) using the base model. For each prompt x, we generate 4 responses  $\{y_i\}_{i=1}^4$ , to balance between computation and solution diversity. We then rely on  $\mathcal{D}_{\text{prob}}$  to calculate the effect of interference  $\Delta^+$  and the relative change in model confidence on each update  $||\Delta(\pi_{\theta}, \mu)||$  throughout the training process. We present our findings in the following section.

## 4.2 NEGATIVE INTERFERENCE IN RLVR TRAINING OF LANGUAGE MODELS

**Decreasing Pass**@k as training progresses. As shown in Fig. 2(A), RLVR training steadily improves the accuracy (Pass@1) of  $\pi_{\theta}$ , averaged across four benchmarks, as training progresses.

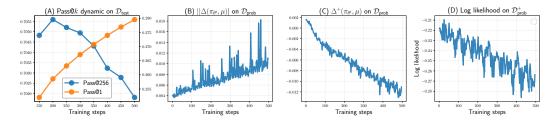


Figure 2: Learning dynamics of RLVR with key trends: (A) RLVR tends to improve average accuracy, but reduce the coverage of solvable solutions, as measured by Pass@256 (averaged across 4 test benchmarks); (B) the increasing effect of influence strength as training progresses; (C) the increase in *negative interference*; and (D) the decline of model confidence on previously correct solutions.

However, when evaluated with a larger sampling budget (e.g., Pass@256), RLVR displays a clear decline, as the performance drops below that of the base model  $\pi_b$  after about 300 steps and continues to deteriorate towards the end of training. The results indicate that RLVR provides little to no improvement in Pass@k at large k's throughout training, a phenomenon consistent with recent findings in (Yue et al., 2025).

The increase of influence on other examples in one training batch update. We also observe an increase in  $||\Delta(\pi_{\theta^t},\mu)||$  during training, as shown in Fig. 2(B). The growing strength of influence at each update step on outside examples suggests stronger similarity among gradients across samples (Tang & Berseth, 2024; Tang et al., 2025), indicating that LLMs struggle to differentiate between distinct problems. As we demonstrate in Sec. 5, this effect can reduce the diversity of previously learned behaviors, where LLMs tend to concentrate on a single reasoning strategy across different problems, and also weaken the regularization effect in RLVR.

Increasing effect of negative interferences. We continue to investigate the interference on the other problem-correct solution pairs in Fig. 2(C), and observe a decreasing trend in  $\Delta^+(\pi_{\theta^t},\mu)$ ; we also observe a similar decreasing trend in the log-likelihood of correct solutions  $\log \pi_{\theta}(y^+|x)$  (Fig. 2(D)). These observations indicate that the learned policy increasingly reduces the likelihood of other problem-correct solution pairs generated by the base model. We call this behavior the *negative interference*, where learning to solve certain problems in  $\mathcal{D}$  inadvertently hinders the performance on other problem-correct solution pairs.

Later in Sec. 5, we will show that, due to RLVR's on-policy sampling, highly solvable problems in the base model tend to capture most of the learning signal. With *negative inference*, this subset of problems will increasingly dominate the learning process while "weaker" ones are suppressed, ultimately constraining the policy to excel only in this subset of contexts while regressing in others. This also explains why RLVR improves average accuracy but reduces coverage. Under a small sampling budget k (e.g., k=1), RLVR incentivizes these "winner" problems to reach higher average accuracy, but this comes at the cost of *negative interference*, where the model's confidence in correct solutions to other problems is degraded. Consequently, the set of problems solvable under a larger budget k becomes narrower.

Interference as an indicator of decreasing Pass@k. The previous observations suggest that the decline in Pass@k performance is likely the result of negative interference in RLVR learning dynamics. To confirm this hypothesis, we uniformly sample two checkpoints T and T' with T' > T during RLVR training. For each pair, we compute the interference  $\Delta^+(\pi_{\theta^{T'}}, \pi_{\theta^T})$  under correct problem-solution pairs of  $\mathcal{D}_{\text{test}}$  and the corresponding change in Pass@k, i.e.,  $\Delta$ Pass@k = Pass@k( $\pi_{\theta^{T'}}$ ) - Pass@k( $\pi_{\theta^T}$ ) with k = 128 (the value of k is selected based on the related prior

works (Yue et al., 2025; Liu et al., 2025a)). We then measure the correlation between these two quantities across four benchmarks.

The correlation results are presented in Fig. 3. As expected, we observe a strong correlation between  $\Delta^+(\pi_{\theta^{T'}}, \pi_{\theta^T})$  and the decrease in Pass@k, suggesting that interference is a key factor driving the degradation of Pass@k performance during training.

# 5 THE EFFECT OF *on-policy* IN RLVR

In this section, we examine which problems RLVR likely prioritizes. First, we show that RLVR primarily incentivizes the language model  $\pi_{\theta}$  to generate responses that already have high likelihood under the base model, regardless of their correctness. This arises from the nature of *on-policy sampling*: for problems where the base model can produce multiple correct solutions, RLVR tends to reinforce the most probable one; but for those where the correct solutions lie in the "valley" regions of the distribution, RLVR, instead, reinforces the base model's highest-confidence – yet potentially incorrect – responses, offering no meaningful learning signal to escape these incorrect modes.

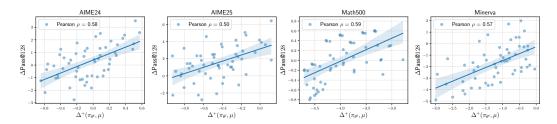


Figure 3: Interference as an indicator of Pass@k decrease across various benchmarks.

Then, we highlight a *winner-take-all* effect that leads to a reduction of previously learned behaviors, as models tend to propagate reasoning strategies from highly solvable problems to others. Finally, we show that existing regularization techniques fail to prevent this issue.

**Experimental Setup.** We consider a subset of problems from the Math500 dataset, denoted as  $\mathcal{D}$ ; the experiments for the other datasets are provided in Appendix H. We utilize the *perplexity* metric (Jurafsky & Martin, 2025):

$$PPL_{\mu}(\pi, \mathcal{D}) = \mathbb{E}_{\substack{\boldsymbol{x} \sim \mathcal{D}, \\ \boldsymbol{y} \sim \mu(\cdot | \boldsymbol{x})}} \left[ \exp \left( -\frac{1}{|\boldsymbol{y}|} \log \pi(\boldsymbol{y} | \boldsymbol{x}) \right) \right]$$
(6)

where  $\mathcal{D} = \{x_i\}_{i=1}^n$  is the set of prompts, the responses are generated by the reference distribution  $y \sim \mu(\cdot|x)$  conditioned on x (|y| denotes the sequence length), and  $\pi$  is the model whose perplexity is being evaluated. We also define the change in average accuracy between the base policy and the learned policy for each problem x:

$$\Delta r(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\cdot|\boldsymbol{x})} \left[ r(\boldsymbol{x}, \boldsymbol{y}) \right] - \mathbb{E}_{\boldsymbol{y} \sim \pi_{b}(\cdot|\boldsymbol{x})} \left[ r(\boldsymbol{x}, \boldsymbol{y}) \right]$$
(7)

and then partition  $\mathcal{D}$  into two subsets:  $\mathcal{D}^{\downarrow}$ , containing problems where RLVR reduces the average accuracy compared to the base model (i.e.,  $\Delta r(\boldsymbol{x}) < 0$ ), and  $\mathcal{D}^{\uparrow}$ , containing problems where RLVR improves upon the base model (i.e.,  $\Delta r(\boldsymbol{x}) > 0$ ). We provide additional experimental details in Appendix H.

**RLVR tends to increase high likelihood regions under the base model.** Fig. 4 shows that during training, the LM  $\pi_{\theta}$  increasingly generates responses that already have high likelihood under the base model  $\pi_b$  (**Leftmost**), while the initial problem-solution pairs sampled from  $\pi_b$  increasingly have lower likelihood (or higher  $PPL_{\pi_b}(\pi_{\theta_t}, \mathcal{D})$ ) in the RLVR policy  $\pi_{\theta^t}$  (**Middle**). Together, these results indicate that the sampled responses from  $\pi_{\theta^t}$  concentrate on the highest likelihood regions in the base model. We also observe in Fig. 4 (**Middle**) that the problem-correct solution pairs sampled from  $\pi_b$  have even lower model confidence (higher  $PPL_{\pi_b}(\pi_{\theta_t}, \mathcal{D}^+)$ ) and a similar decreasing trend, suggesting that for each problem, RLVR collapses into the solution with the highest likelihood

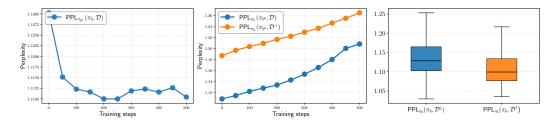


Figure 4: Perplexity during RLVR training. **Leftmost**: the data sampled from each intermediate checkpoint  $\pi_{\theta^t}$  exhibits an increasingly high model confidence under  $\pi_b$ . **Middle**: RLVR models  $\pi_{\theta^t}$  exhibit reduced confidence in data previously generated by the base model, regardless of their correctness. **Rightmost**: problems that RLVR improves already have a high likelihood of generating correct solutions under  $\pi_b$ , while coverage-reduced problems initially have a low likelihood of producing correct answers.

under  $\pi_b$  regardless of whether these high likelihood responses is correct or not. Finally, Fig. 4 (**Rightmost**) demonstrates that, for problems where the base model already assigns high likelihood to correct solutions (lower PPL( $\pi_{\theta_t}, \mathcal{D}^{\uparrow}$ )), RLVR reinforces these responses, leading to performance gains (or higher average accuracies). In contrast, when the correct responses initially lie in low-likelihood regions (higher PPL( $\pi_{\theta_t}, \mathcal{D}^{\downarrow}$ )), RLVR tends to reduce its probability, reinforcing the base model's initial biases, resulting in degraded performance.

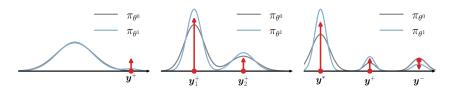


Figure 5: An illustration describing the dynamic of on-policy learning in Eq. (8). **Leftmost**: correct response  $y^+$  in low-likelihood regions induce minimal effect. **Middle**: with multiple correct responses  $y_1^+$  and  $y_2^+$ , updates favor the one with higher initial likelihood. **Rightmost**: negative gradients on incorrect response  $y^-$  can raise correct ones  $y^+$ , but greedy responses  $y^*$  increase the most.

Why does RLVR reinforce problems with high-likelihood correct solutions in the base model? To explain why this phenomenon happens, we first confirm that it occurs when *on-policy sampling* is used for learning (the details are in Appendix D). More specifically, this phenomenon also appears even in a simple bandit problem with a Softmax policy with V actions. We can look into the on-policy objective REINFORCE for an action y:

$$\nabla_{\theta} \mathcal{L}^{\text{REINFORCE}}(\pi_{\theta}) = A(\boldsymbol{y}) \nabla_{\theta} \pi_{\theta}(\boldsymbol{y}) = A(\boldsymbol{y}) \pi_{\theta}(\boldsymbol{y}) \nabla_{\theta} \log \pi_{\theta}(\boldsymbol{y})$$
(8)

and arrive at the following insights (additionally, we utilize Fig. 5 for illustration):

- Low-likelihood tokens provide less meaningful updates. Eq. (8) suggests that low likelihood, correct responses receive infrequent and small updates because of the explicit  $\pi_{\theta}(y)$  scaling factor and the fact that they are rarely sampled (Fig. 5 Leftmost), while high-likelihood tokens dominate the learning process. Moreover, if  $\pi_{\theta}(y) = 0$ , the gradient vanishes, creating an additional saddle point where zero-probability tokens cannot be updated.
- When there are multiple optimal actions  $y_1^+$  and  $y_2^+$ , the one with with the higher likelihood  $\pi_{\theta}(\cdot)$  will dominate. Due to the scaling of the model probability in Eq. 8, the update will tilt probability toward the already dominant correct mode, as depicted in Fig. 5 (**Middle**). This effect will exacerbate when the gap between the two probabilities increases.
- Negative gradients tend to reinforce high-likelihood tokens, regardless of their correctness. Onpolicy learning with negative gradient exhibits similar behavior to off-policy with negative gradient. When the correct solution resides in a "valley" of the distribution, negative gradient tends

to increase high-likelihood tokens regardless of correctness (Fig. 5 **Rightmost**), further diminishing the low-likelihood (but correct) ones (especially when the distribution is peaky). This is also observed in Ren & Sutherland (2025).

• If  $y^+$  is unlikely to be sampled under the k-sampling budget  $\pi_{\theta}$ , RLVR provides no learning signal. This arises from the advantage calculation, where A(y) = 0 if all sampled actions are incorrect, resulting in no update.

We provide additional theoretical and empirical details of this effect by analyzing the ratio of model predictions  $\pi_{\theta^{t+1}}/\pi_{\theta^t}$  in Appendix D.

**On-policy learning and negative interference lead to Winner-take-all.** The above analysis suggests that highly solvable problems dominate the learning signal, while problems with lower success rates contribute less due to gradient vanishing. Learning primarily on these highly solvable problems also increases the chance of *negative interference* on problems with lower success rates, and eventually exacerbates the effect of *winner-take-all*; as training progresses, these problems with lower success rates will progressively receive less explicit learning signal and have lower likelihoods of sampling correct solutions.

Winner-take-all explains the coverage shrinkage and performance degradation in Minerva Benchmark. Shao et al. (2025) report that the Qwen2.5-Math model family employs two distinct reasoning modes: code-based and natural-language reasoning. However, during RLVR training, we observe a progressive collapse to natural language reasoning. As shown in Fig. 6 (Left), the frequency of the generated solutions using code reasoning steadily decreases over training, while the frequency of natural language reasoning dominates at the end of training.

We observe that natural language reasoning has a better initial accuracy in the training dataset  $\mathcal{D}$ 

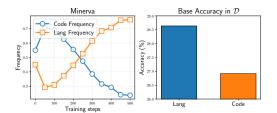


Figure 6: The LM switches from *code reasoning* to *language reasoning* in the Minerva benchmark, due to the reduction of diversity of previously learned behaviors.

(Fig. 6 (**Right**)). As RLVR training progresses, the LM tends to concentrate on this subset of problems that are already highly solvable under the base model. Winner-take-all suggests that RLVR amplifies the highest-likelihood solutions while suppressing other correct ones, leading to reduced diversity among these solvable problems and a collapse to natural language reasoning. As the influence strength  $||\Delta(\pi_{\theta^t}, \mu)||$  grows (as shown in Fig. 2), this loss of diversity and the resulting bias towards high-probability responses under base model  $\pi_b$  becomes amplified propagate, affecting other problem–solution pairs, including the test problems. Consequently, in cases where the base model's initial bias is incorrect, this self-bias amplification will result in coverage shrinkage and negative interference. This explains that on the Minerva benchmark (the test set), while code reasoning achieves higher initial accuracy, as RLVR training gradually switches to language reasoning, the performance on code reasoning problems drops, as observed in Fig. 1.

Existing regularizations fail to mitigate diversity shrinkage. In RLVR training, trust region constraints such as clipping aim to prevent the learned policy from deviating too far from the previously updated policy. As shown in Fig. 7, the average probability ratio  $\rho = \pi_{\theta}/\pi_{\text{old}}$  that violates the trust region constraint (lies outside the clipping range  $(1 - \epsilon, 1 + \epsilon)$  ( $\epsilon = 0.2$ )) remains stable in the training batch. However, under the probing dataset  $\mathcal{D}_{\text{prob}}$ , this violation steadily grows as training progresses. This suggests that the clipping mechanism has little influence on data outside the training batch, thus failing to prevent the *increasing negative interference* throughout the training process. Another regularizer, the Reverse

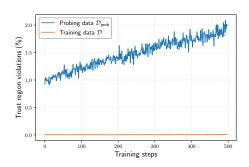


Figure 7: Clipping violations on training data and probing data.

KL, is also ineffective, as Reverse KL is known to exhibit mode-seeking behavior, where it has little

to no penalty on regions where  $\pi_{\theta}$  puts low probability mass. Furthermore, Reverse KL often favors collapsing probability mass onto high-likelihood mode under the base model  $\pi_b$ .

### 6 A DATA CURATION TECHNIQUE TO MITIGATE WINNER-TAKE-ALL EFFECT.

Building on the above analysis, we propose a novel and effective algorithm that focuses *learning* only on problems where the greedy response fails. This acts as a proxy metric to exclude highly solvable problems from learning, thus preventing them from dominating the learning signal. To further preserve the diversity of learned behaviors, we replace the Reverse KL regularization with a Forward KL (SFT loss) objective, which penalizes the model when it begins to forget previously learned behaviors. We refer to this method as SELF (Selective Examples with Low-likelihood and Forward-KL):

$$\mathcal{J}(\pi_{\theta}) = \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}, (\boldsymbol{y}^*, \boldsymbol{y}) \sim \pi_{\theta}(\cdot | \boldsymbol{x})} \left[ 1\{r(\boldsymbol{x}, \boldsymbol{y}^*) \in \mathcal{C}(\boldsymbol{x})\} r(\boldsymbol{x}, \boldsymbol{y}) - \beta \text{KL}(\pi_b || \pi_{\theta}) \right]$$
(9)

where  $y^*$  is the greedy response and  $C(x) = \{y | r(x, y) = 1\}$  is the set of correct completions given a problem x.

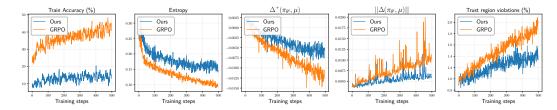


Figure 8: Learning dynamics of GRPO and the proposed method, with key observations: 1) SELF focus on learning problems with low success rate, 2) SELF exhibits better diversity, 3) SELF mitigates the effect of *negative interference* and reduce the influence strength to other examples, result in a strong regularization effect.

We compare SELF learning dynamics against the standard GRPO objective in Fig. 8. Since SELF focuses on problems with a low likelihood of being solved correctly, it achieves a lower average training reward than GRPO. However, it preserves better diversity throughout training, as measured by higher token-level entropy across the training process. In addition, SELF effectively mitigates the *negative interference effect* as we observe a significant improvement in  $\Delta^+(\pi_\theta, \pi_b)$  and a smaller magnitude of influence  $||\Delta(\pi_\theta, \mu)||$  on other examples. Furthermore, SELF reduces the extent of trust region violations for examples outside the training batch.

Table 1: Evaluation results on various mathematical benchmarks (%). The best performance under each dataset is marked with **boldface**.

Model	Method		ME <b>2024</b> Pass@102		ME <b>2025</b> Pass@102		th500 Pass@250		inerva Pass@256
Qwen2.5-Math-1.5B	Base model GRPO	9.73 <b>15.16</b>	75.86 67.77	4.45 <b>8.50</b>	67.80 64.05	58.83 <b>72.67</b>	95.58 95.0	16.82 23.96	<b>70.67</b> 65.07
	SELF	14.33	80.94	8.05	69.45	71.0	96.0	23.14	68.75
Qwen2.5-Math-7B	Base model GRPO	16.94 <b>27.95</b>	81.83 84.77	6.97 14.13	67.57 58.60	59.71 <b>80.36</b>	96.0 92.8	11.60 29.40	66.91 60.66
	SELF	25.81	89.62	14.62	72.0	79.82	97.80	30.42	71.70
Llama-3.2-3B-Instruct	Base model GRPO	3.53 11.76	<b>53.7</b> 49.25	0.32 0.38	48.22 25.28	26.61 <b>54.06</b>	<b>91.4</b> 88.2	8.25 15.81	49.26 50.73
	SELF	13.30	51.69	0.65	51.83	52.12	91.0	17.0	57.35

Table 1 shows the evaluation of SELF and the baselines on various mathematical benchmarks. SELF achieves comparable Pass@1 performance compared to GRPO, but consistently better performance for larger values of k. While GRPO achieves good Pass@1 performance, it also reinforces

problems that already have high success rates, leading to narrower coverage at larger k. Interestingly, larger models tend to suffer from even more severe coverage shrinkage than smaller models, despite achieving higher Pass@1. In contrast, SELF scales more effectively, outperforming GRPO on Pass@1 in some experiments while also frequently surpassing the base model under larger k budgets. We also provide details of the computational cost of SELF in Appendix G, where SELF demonstrates a negligible computation difference to GRPO.

#### 7 CONCLUSION

We investigate why RLVR exhibits coverage shrinkage of solvable problems, indicated by reduced Pass@k performance, compared to the base model, by analyzing its learning dynamics. We identify a detrimental effect, which we call *negative interference*, where learning on a subset of training problems can hinder performance on others. This phenomenon serves as a key indicator of the performance decline in larger k budgets. Furthermore, we find that RLVR tends to concentrate learning on highly solvable problems, leading to reduced diversity and amplification of self-bias toward both previously solvable and unseen problems. Motivated by these findings, we propose SELF, a simple yet effective algorithm that only learn on problems with low success rates, effectively mitigating the coverage shrinkage problem.

#### REFERENCES

Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting REINFORCE-style optimization for learning from human feedback in LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12248–12267, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.662. URL https://aclanthology.org/2024.acl-long.662/.

Art of Problem Solving. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME\_Problems\_and\_Solutions. Accessed: 2025-04-20.

Jordan T. Ash and Ryan P. Adams. On warm-starting neural network training. In NeurIPS, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 288cd2567953f06e460a33951f55daaf-Abstract.html.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.

Xingyu Dang, Christina Baek, J Zico Kolter, and Aditi Raghunathan. Assessing diversity collapse in reasoning. In *Scaling Self-Improving Foundation Models without Human Supervision*, 2025. URL https://openreview.net/forum?id=AMiKsHLjQh.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Runxin Xu, Qihao Zhu, and Yiliang ... Xiong. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, January 22 2025. doi: 10.48550/arXiv.2501.12948. URL https://arxiv.org/abs/2501.12948.

Pierluca D'Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G Bellemare, and Aaron Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=OpC-9aBBVJe.

Michelle Guo, Albert Haque, De-An Huang, Serena Yeung, and Li Fei-Fei. Dynamic task prioritization for multitask learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL https://openreview.net/forum?id=7Bywt2mQsCe.
- Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*, 2024.
- Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson. Transient non-stationarity and generalisation in deep reinforcement learning. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=Qun8fv4qSby.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
- Arthur Juliani and Jordan T. Ash. A study of plasticity loss in on-policy deep reinforcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=MsUf8kpKTF.
- Daniel Jurafsky and James H. Martin. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, with Language Models. 3rd edition, 2025. URL https://web.stanford.edu/~jurafsky/slp3/. Online manuscript released August 24, 2025.
- Matthew Kielo and Vladimir Lukin. It's time to move on: Primacy bias and why it helps to forget. In *The Third Blogpost Track at ICLR 2024*, 2024. URL https://openreview.net/forum?id=7ZEwqUyKMC.
- Timo Klein, Lukas Miklautz, Kevin Sidak, Claudia Plant, and Sebastian Tschiatschek. Plasticity loss in deep reinforcement learning: A survey. *CoRR*, abs/2411.04832, 2024. URL https://doi.org/10.48550/arXiv.2411.04832.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxi Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training. In Second Conference on Language Modeling, 2025. URL https://openreview.net/forum?id=iluGbfHHpH.
- Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=IFXTZERXdM7.
- Yuetai Li, Zhangchen Xu, Fengqing Jiang, Bhaskar Ramasubramanian, Luyao Niu, Bill Yuchen Lin, Xiang Yue, and Radha Poovendran. Temporal sampling for forgotten reasoning in LLMs. In 2nd AI for Math Workshop @ ICML 2025, 2025a. URL https://openreview.net/forum?id=1medQAgr1g.
- Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Zhi-Quan Luo, and Ruoyu Sun. Preserving diversity in supervised fine-tuning of large language models. In *The Thirteenth International Conference on Learning Representations*, 2025b. URL https://openreview.net/forum?id=NQEe7B7bSw.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=v8L0pN6EOi.

- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models. *CoRR*, abs/2505.24864, May 2025a. URL https://doi.org/10.48550/arXiv.2505.24864.
- Vincent Liu, Adam White, Hengshuai Yao, and Martha White. Towards a practical measure of interference for reinforcement learning, 2020. URL https://arxiv.org/abs/2007.03807.
- Vincent Liu, Han Wang, Ruo Yu Tao, Khurram Javed, Adam White, and Martha White. Measuring and mitigating interference in reinforcement learning, 2023. URL https://arxiv.org/ abs/2307.04887.
- Zichen Liu, Changyu Chen, Wenjun Li, Tianyu Pang, Chao Du, and Min Lin. There may not be aha moment in r1-zero-like training a pilot study. https://oatllm.notion.site/oat-zero, 2025b. Notion Blog.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. In 2nd AI for Math Workshop @ ICML 2025, 2025c. URL https://openreview.net/forum?id=jlpC1zavzn.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing olpreview with a 1.5b model by scaling rl. https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-Ol-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005k 2025. Notion Blog.
- MetaAI. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. 2024b. URL https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/.
- Skander Moalla, Andrea Miele, Razvan Pascanu, and Caglar Gulcehre. No representation, no trust: Connecting representation, collapse, and trust issues in PPO. In *Seventeenth European Workshop on Reinforcement Learning*, 2024. URL https://openreview.net/forum?id=8esP9hxdsM.
- OpenAI, Josh Achiam, and et al. Steven Adler. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.
- Yuxiao Qu, Matthew Y. R. Yang, Amrith Setlur, Lewis Tunstall, Edward Emanuel Beeching, Ruslan Salakhutdinov, and Aviral Kumar. Optimizing test-time compute via meta reinforcement finetuning. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=TqODUDsU4u.
- Yi Ren and Danica J. Sutherland. Learning dynamics of LLM finetuning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tPNHOoZF19.
- Tom Schaul, Diana Borsa, Joseph Modayil, and Razvan Pascanu. Ray interference: a source of plateaus in deep reinforcement learning. *arXiv preprint arXiv:1904.11455*, April 25 2019. doi: 10.48550/arXiv.1904.11455. URL https://arxiv.org/abs/1904.11455.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. URL https://arxiv.org/abs/1707.06347.
- Amrith Setlur, Matthew Y. R. Yang, Charlie Snell, Jeremy Greer, Ian Wu, Virginia Smith, Max Simchowitz, and Aviral Kumar. e3: Learning to explore enables extrapolation of test-time compute for llms, 2025. URL https://arxiv.org/abs/2506.09026.
- Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei Du, Nathan Lambert, Sewon Min, Ranjay Krishna, Yulia Tsvetkov, Hannaneh Hajishirzi, Pang Wei Koh, and Luke Zettlemoyer. Spurious rewards: Rethinking training signals in rlvr, 2025. URL https://arxiv.org/abs/2506.10947.

- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:* 2409.19256, 2024.
- Hongyao Tang and Glen Berseth. Improving deep reinforcement learning by reducing the chain effect of value and policy churn. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=cQoAgPBARc.
- Hongyao Tang, Johan Obando-Ceron, Pablo Samuel Castro, Aaron Courville, and Glen Berseth. Mitigating plasticity loss in continual reinforcement learning by reducing churn. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=EkoFXfSauv.
- Fang Wu, Weihao Xuan, Ximing Lu, Zaid Harchaoui, and Yejin Choi. The invisible leash: Why rlvr may not escape its origin, 2025. URL https://arxiv.org/abs/2507.14843.
- An Yang, Beichen Zhang, and et al. Binyuan Hui. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024a. URL https://arxiv.org/abs/2409.12122.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL https://openreview.net/forum?id=mXpq6ut8J3.
- Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale. *CoRR*, abs/2503.14476, March 2025. URL https://doi.org/10.48550/arXiv.2503.14476.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*, April 18 2025. URL https://arxiv.org/abs/2504.13837.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerlzoo: Investigating and taming zero reinforcement learning for open base models in the wild, 2025. URL https://arxiv.org/abs/2503.18892.
- Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: Rl post-training amplifies behaviors learned in pretraining. *arXiv* preprint *arXiv*:2504.07912, April 10 2025. URL https://arxiv.org/abs/2504.07912.
- Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. The surprising effectiveness of negative reinforcement in llm reasoning. *arXiv preprint arXiv:2506.01347*, June 2025. URL https://arxiv.org/abs/2506.01347.

#### A EXPERIMENTAL DETAILS

#### A.1 TRAINING DETAILS

**Models.** We utilize two family of models for RLVR training: Qwen2.5-Math-1.5B and Qwen2.5-Math-7B (Yang et al., 2024a); and Llama-3.2-3B-Instruct (MetaAI, 2024b).

**Training configurations.** We utilize the VERL framework (Sheng et al., 2024) for RLVR training. We train each model on 4 GPUs with a constant learning rate of  $10^{-6}$ , the maximum response length is 3072 for Qwen2.5-Math models, and 8192 maximum generation tokens for Llama-3.2-3B-Instruct. We use a mini batch size and a rollout batch size of 64. For each prompt, we collect 8 rollouts to compute advantages for the GRPO update. We use a sampling temperature  $\tau=1$ . We do not apply Reverse KL divergence or entropy loss in our training, and the total training steps are 500, similar to prior works (Shao et al., 2025; Yu et al., 2025; Luo et al., 2025). We set the hyperparameter controlling the trade-off between reward maximization and Forward KL regularization to  $\beta=10^{-4}$ .

**Reward function.** For the reward function r(x, y) calculation, we follow the implementation of (Zeng et al., 2025; Hu et al., 2024) with the following rule for faster convergence:

- If the response provides a final answer in the boxed format and is correct, it receives a reward of +1.
- If the response provides a final answer in the boxed format, but it is incorrect, we set the reward as -0.5.
- Otherwise, the reward is set to -1 for incorrect answers.

#### A.2 PROBING DATASET CONSTRUCTION

One major problem when analyzing learning dynamics is the huge response space  $\mathcal{Y}$ : the number of possible responses  $\boldsymbol{y} \in \mathcal{V}^L$ , not only includes natural language but also non-sensical texts. In this paper, we are interested in the response region from the base model  $\pi_b$ , due to the following reason: the concentration coefficient  $C(\pi_\theta, \pi_b)$  is small, since RLVR tends to reduce the coverage compared to the base model, we can expect that base model distribution can provide coverage of the learned policy  $\pi_\theta$ .

To create the probing dataset  $\mathcal{D}_{\text{prob}}$ , for each problem in the training dataset  $x \in \mathcal{D}$ , we sample k = 4 responses from the base model with temperature  $\tau = 0.9$  to ensure diversity and accuracy.

## A.3 PROMPT TEMPLATE

We use the prompt templates for Qwen2.5-Math models and Llama-3.2-3B-Instruct following (Yang et al., 2024a; Luo et al., 2025), as shown in Table 2 and Table 3.

Table 2: Prompt template for Qwen2.5-Math.

```
<|im_start|>system
Please reason step by step, and put your final answer within \boxed{}.<|im_end|>
<|im_start|>user
{input}
Let's think step by step and output the final answer within \boxed{}.<|im_end|>
<|im_start|>assistant
```

Table 3: Prompt template for Llama-3.2-3B-Instruct.

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>
Cutting Knowledge Date: December 2023
<|start_header_id|>user<|end_header_id|>
{input}
Let's think step by step and output the final answer within \boxed{}.<|eot_id|>
<|eot_id|><|start_header_id|>assistant<|end_header_id|>
```

#### A.4 EVALUATION DETAILS

Across all checkpoints, we evaluate using a temperature  $\tau=0.6$  and a top-p value of 0.95, following prior works (Yue et al., 2025; Zeng et al., 2025; Lewkowycz et al., 2022), with a maximum generation of 3072 tokens for Qwen2.5-Math due to limited context length. For Llama-3.2-3B-Instruct, we allow the model to generate a maximum of 16384 tokens.

#### B DETAILS DERIVATION OF GRPO

Given a problem x and the generated solutions  $\{y\}_{i=1}^G$ , where G>2 is the number of generated solutions per prompt. We define each token index t in the response y as  $y_t$  with  $1 \le t \le |y|$ , where |y| is the sequence length of the reasoning traces. Specifically, for each problem x, GRPO objective DeepSeek-AI et al. (2025) will first calculate group-wise normalized advantage function:

$$A(\boldsymbol{x},\boldsymbol{y}_i) = \frac{r(\boldsymbol{x},\boldsymbol{y}_i) - \hat{\mu}}{\hat{\sigma}}$$
 (10) where  $\hat{\mu} = \frac{1}{G} \sum_{j=1}^{G} r(\boldsymbol{x},\boldsymbol{y}_j)$ , and  $\hat{\sigma} = \sqrt{\frac{1}{G-1} \sum_{j=1}^{G} (r(\boldsymbol{x},\boldsymbol{y}_j) - \hat{\mu})^2}$ .

Let  $\pi_{\text{old}}$  denote the policy from the previous update step that generates the reasoning traces. The probability ratio  $\rho_t = \frac{\pi_{\theta}(y_t|\boldsymbol{x},\boldsymbol{y}_{< t})}{\pi_{\text{old}}(y_t|\boldsymbol{x},\boldsymbol{y}_{< t})}$  at each token index t, along with PPO clipping mechanism to serves as a proxy of trust region contraint with clipping threshold  $\epsilon$ , the GRPO objective is defined as:

$$J(\theta) = \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}, \, \boldsymbol{y} \sim \pi_{\text{old}}(\cdot | \boldsymbol{x})} \left[ \sum_{t=1}^{|\boldsymbol{y}|} \min \left( \rho_t(\boldsymbol{y}; \theta) \, A(\boldsymbol{x}, \boldsymbol{y}_i), \, \operatorname{clip}(\rho_t(\boldsymbol{y}; \theta), 1 - \epsilon, 1 + \epsilon) A(\boldsymbol{x}, \boldsymbol{y}_i) \right) \right] -\beta \, \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}} \left[ \operatorname{KL}(\pi_{\theta}(\cdot | \boldsymbol{x}) \parallel \pi_b(\cdot | \boldsymbol{x})) \right].$$

$$(11)$$

The KL regularization helps the learned policy not deviate too far from the base model. Recently, Yu et al. (2025) suggests that removing KL regularization can provide better performance. Following this, we set  $\beta = 0$  to eliminate KL regularization in our work.

#### C PER-STEP INFLUENCE PROOFS AND FURTHER ANALYSIS

# C.1 Proof of Proposition 3.1

**Proposition C.1.** Consider a problem-solutions pair (x, y) for updating using policy gradient objective, a sufficiently small learning rate  $\eta$ , the per-step influence on a problem-solution pair (x', y') is defined as:

$$\Delta \log \pi_{\theta}^{t}(\mathbf{y}'|\mathbf{x}') = \eta \mathcal{K}^{t}(\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}') A(\mathbf{x}, \mathbf{y}) + \mathcal{O}(\eta^{2})$$
(12)

where  $\mathcal{K}^t(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}) = \nabla_{\theta} \log \pi_{\theta^t}(\boldsymbol{y}|\boldsymbol{x})^{\top} \nabla_{\theta} \log \pi_{\theta^t}(\boldsymbol{y}'|\boldsymbol{x}')$  measures the influence between  $(\boldsymbol{x}, \boldsymbol{y})$  and  $(\boldsymbol{x}', \boldsymbol{y}')$ , and the advantage function  $A(\boldsymbol{x}, \boldsymbol{y})$  determines the magnitude and direction of the model's update.

*Proof.* We first approximate  $\log \pi_{\theta^{t+1}}(y|x)$  using Taylor expansion. For simplicity, we denote  $\pi_{\theta^t} = \pi_{\theta}^t$ :

$$\log \pi_{\theta}^{t+1}(\boldsymbol{y}|\boldsymbol{x}) = \log \pi_{\theta}^{t}(\boldsymbol{y}|\boldsymbol{x}) + \nabla_{\theta} \log \pi_{\theta}^{t}(\boldsymbol{y}|\boldsymbol{x})^{\top} (\theta^{t+1} - \theta^{t}) + \mathcal{O}(\eta^{2})$$
(13)

Using the definition of stochastic gradient ascent with a sufficiently small learning rate  $\eta$ , we have that:

$$\theta^{t+1} = \theta^t + \eta \nabla_{\theta} \mathcal{J}(\pi_{\theta}) \tag{14}$$

where  $\nabla_{\theta} \mathcal{J}(\pi_{\theta}) = A(\boldsymbol{x}, \boldsymbol{y}) \nabla_{\theta} \log \pi_{\theta}(\boldsymbol{y}|\boldsymbol{x})$  is the standard policy gradient objective (Qu et al., 2025; Ahmadian et al., 2024), we plug in the above objective in Eq. 13:

$$\Delta \log \pi_{\theta}^{t}(\mathbf{y}'|\mathbf{x}') = \eta \mathcal{K}(\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}') A(\mathbf{x}, \mathbf{y}) + \mathcal{O}(\eta^{2})$$
(15)

which concludes the proof.

#### C.2 Interpretation of the relative strength of influence $||\Delta(\pi_{\theta^t}\mu)||$

If we view the language model  $\pi_{\theta}$  as a neural network with a problem and solution pair o = (x, y) as an input, where the output is the log-probability

$$\log \pi_{\theta}(\boldsymbol{y}|\boldsymbol{x}) = \frac{1}{|\boldsymbol{y}|} \sum_{t=1}^{|\boldsymbol{y}|} \log \pi_{\theta}(y_t|\boldsymbol{x}, \boldsymbol{y}_{< t})$$
(16)

where  $y_{< t}$  is the partial completion before the token index t. We consider the gradient of  $\pi_{\theta}$  with respect to parameter  $\theta$  at (x, y) as  $\nabla_{\theta} \log \pi_{\theta}(y|x) \in \mathbb{R}^d$  where d is the parameter dimension. We can define a matrix of gradient dot products  $\mathcal{K}(x, x', y, y') = \mathcal{K}(o, o')$  across all different promptresponse pair:

$$\mathcal{K} = \begin{bmatrix}
\mathcal{K}(\boldsymbol{o}_1, \boldsymbol{o}_1) & \mathcal{K}(\boldsymbol{o}_1, \boldsymbol{o}_2) & \cdots \\
\mathcal{K}(\boldsymbol{o}_2, \boldsymbol{o}_1) & \mathcal{K}(\boldsymbol{o}_2, \boldsymbol{o}_2) & \cdots \\
\vdots & \vdots & \ddots
\end{bmatrix} = \nabla_{\theta} \log \pi_{\theta} (\mathcal{Y}|\mathcal{X})^{\top} \nabla_{\theta} \log \pi_{\theta} (\mathcal{Y}|\mathcal{X}) \tag{17}$$

where  $\nabla_{\theta} \log \pi_{\theta}(\mathcal{Y}|\mathcal{X}) = [\nabla_{\theta} \log \pi_{\theta}(y_1|x_1), \nabla_{\theta} \log \pi_{\theta}(y_2|x_2), \cdots]$  denotes the matrix of the gradient of different prompt-solution pairs (x, y).

In an idealized scenario, the kernel matrix  $\mathcal{K}$  is diagonal, meaning that all off-diagonal terms are zeros, i.e.,  $\mathcal{K}(o_i,o_j)=0$  for  $i\neq j$ . Under this assumption, improving performance on a given training problem-solution pair would not influence any other unseen problems. Consequently, there would be neither *negative interference*, where training on one example degrades performance on others, nor *coverage shrinkage*, where the set of solvable problems becomes narrower.

In practice, however, this assumption rarely holds for language models. Because parameters are shared across all examples, updates to one example inevitably affect others. This manifests as non-zero off-diagonal entries in  $\mathcal{K}$ , which capture the degree of cross-influence between examples.

The growing magnitude of the relative influence strength,  $||\Delta(\pi_{\theta^t}, \mu)||$ , can thus be interpreted as evidence that the number and magnitude of these off-diagonal interactions are increasing. Specifically, this signals that the gradients of different examples are becoming more linearly dependent, since the rank of the kernel  $\mathcal{K}$  is equivalent to the rank of the gradient matrix  $\nabla_{\theta} \log \pi_{\theta}(\mathcal{Y}|\mathcal{X})$ . This growing entanglement explains why optimization on a subset of problems may inadvertently interfere with generalization to others.

# D ANALYZING LEARNING DYNAMIC OF ON-POLICY LEARNING IN TOY BANDIT PROBLEM.

We first consider a contextual bandit problem with a parametrized softmax policy with V actions. where the input data x is a d dimensional feature vector  $x \in \mathbb{R}^d$ . The policy is defined as a linear

layer with Softmax output:

$$\pi_{\theta}(y_i|\boldsymbol{x}) = \frac{\exp(z_i(\boldsymbol{x}))}{\sum\limits_{j=1}^{V} \exp(z_j(\boldsymbol{x}))}$$
(18)

We consider the objective of REINFORCE and Off-policy maximum likelihood objectives for 2 actions  $y_1$  and  $y_2$  that receive a positive reward (r(x, y) = 1).

$$\mathcal{J}^{\text{REINFORCE}}(\pi_{\theta}) = \pi_{\theta}(\boldsymbol{a}_{1}|\boldsymbol{x}) + \pi_{\theta}(\boldsymbol{a}_{2}|\boldsymbol{x}), \quad \mathcal{J}^{\text{MLE}}(\pi_{\theta}) = \log \pi_{\theta}(\boldsymbol{a}_{1}|\boldsymbol{x}) + \log \pi_{\theta}(\boldsymbol{a}_{2}|\boldsymbol{x}) \quad (19)$$

Using gradient descent, we can perform parameter update at each update step t with a learning rate  $\eta$ :

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^{t} + \eta \nabla_{\boldsymbol{\theta}} \mathcal{J}^{\text{REINFORCE}} = \boldsymbol{\theta}^{t} - \eta \boldsymbol{x} \left[ (\pi_{\boldsymbol{\theta}^{t}}(y_{1})(\pi_{\boldsymbol{\theta}^{t}}(\boldsymbol{y}) - \boldsymbol{e}(y_{1}))^{\top} + \pi_{\boldsymbol{\theta}^{t}}(y_{2})(\pi_{\boldsymbol{\theta}^{t}}(\boldsymbol{y}) - \boldsymbol{e}(y_{2}))^{\top} \right]$$
$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^{t} + \eta \nabla_{\boldsymbol{\theta}} \mathcal{J}^{\text{MLE}} = \boldsymbol{\theta}^{t} - \eta \boldsymbol{x} \left[ ((\pi_{\boldsymbol{\theta}^{t}}(\boldsymbol{y}) - \boldsymbol{e}(\boldsymbol{y}_{1})^{\top} + (\pi_{\boldsymbol{\theta}^{t}}(\boldsymbol{y}) - \boldsymbol{e}(\boldsymbol{y}_{2}))^{\top} \right]$$
(20)

where  $\pi_{\theta}(y) \in \mathbb{R}^V$  is the vector represents probability distribution over all actions, e(y) is a one-hot vector indicates the target action. A key difference between the REINFORCE and MLE objectives is that, in REINFORCE, the probability of the sampled actions directly scales the update step. As a result, actions with higher likelihood receive larger updates, while low-likelihood actions lead to much smaller changes. To analyze how the model's probabilities in each action change, we define a ratio  $\alpha_i = \frac{\pi_i^{t+1}}{\pi^{t_i}}$  and use The following lemma describes its behavior:

**Lemma D.1.** The ratio of confidence change for each i can be represented as:

$$\alpha_{i}^{\textit{MLE}} = \frac{\pi_{i}^{t+1}}{\pi_{i}^{t}} = \frac{\sum_{j=1}^{V} e^{z_{j}^{t}}}{\sum_{j=1}^{V} \beta_{j} e^{z_{j}^{t}}}; \quad \alpha_{i}^{\textit{REINFORCE}} = \frac{\sum_{j=1}^{V} e^{z_{j}^{t}}}{\sum_{j=1}^{V} \gamma_{j} e^{z_{j}^{t}}}$$
(21)

Note that the values of  $\beta_j$  also depends on whether  $i \in \{y_1, y_2\}$ , hence for Case 1 ( $i \in \{y_1, y_2\}$ ), and Case 2 ( $i \notin \{y_1, y_2\}$ ):

Case 1

$$\beta_{j} = \begin{cases} e^{-\eta' \left( 2(\pi_{j}^{t} - \pi_{i}^{t} + 1 \right)} & \text{if } j \notin \{y_{1}, y_{2}\} \\ 1 & \text{if } j = y_{1} ; \\ e^{-2\eta' \left( \pi_{j}^{t} - \pi_{i}^{t} \right)} & \text{if } j \neq \{y_{1}, y_{2}\} \end{cases}$$

$$if j = y_{1} ; \quad \gamma_{j} = \begin{cases} e^{-\eta' \left( \Delta \pi^{t} (\pi_{j}^{t} - \pi_{i}^{t}) + \pi_{i}^{t} \right)} & \text{if } j \notin \{y_{1}, y_{2}\} \\ 1 & \text{if } j = y_{1} \end{cases} (22)$$

$$e^{-\eta' \left( (\Delta \pi^{t} - 1)(\pi_{j}^{t} - \pi_{i}^{t}) + \pi_{i}^{t} \right)} & \text{if } j = y_{2}$$

Case 2.

$$\beta_{j} = \begin{cases} e^{-2\eta' \left(\pi_{j}^{t} - \pi_{i}^{t}\right)} & \text{if } j \notin \{y_{1}, y_{2}\} \\ e^{-\eta' \left(2\left(\pi_{j}^{t} - \pi_{i}^{t}\right) - 1\right)} & \text{if } j \in \{y_{1}, y_{2}\} \end{cases}; \quad \gamma_{j} = \begin{cases} e^{-\eta' \Delta \pi^{t} \left(\left(\pi_{j}^{t} - \pi_{i}^{t}\right)\right)} & \text{if } j \notin \{y_{1}, y_{2}\} \\ e^{-\eta' \left(\Delta \pi^{t} \left(\pi_{j}^{t} - \pi_{i}^{t}\right) - \pi_{j}^{t}\right)} & \text{if } j \in \{y_{1}, y_{2}\} \end{cases}$$
(23)

where  $\Delta \pi^t = \pi^t_{y_1} + \pi^t_{y_2}, \eta' = \eta ||x||_2^2$ .

*Proof.* We first need to link the logits vector  $z^{t+1}$  and  $z^t$ . From Eq.20,  $z^{t+1}$  can be recursively written as:

$$z^{t+1} = (\boldsymbol{\theta}^{t+1})^{\top} \boldsymbol{x}$$

$$= (\boldsymbol{\theta}^{t} + \eta \boldsymbol{x} \nabla_{z} \mathcal{J}))^{\top} \boldsymbol{x}$$

$$= (\boldsymbol{\theta}^{t})^{\top} \boldsymbol{x} + \eta (\boldsymbol{x} \nabla_{z} \mathcal{J})^{\top} \boldsymbol{x}$$

$$= z^{t} + \eta' \nabla_{z} \mathcal{J}$$
(24)

where  $\eta' = \eta ||x||_2^2$ . We can write down for each value in vector  $z^{t+1}$  for MLE and REINFORCE objectives as:

$$\begin{split} \text{MLE: } z_i^{t+1} &= \begin{cases} z_i^t - 2\eta' \pi_i^t & \text{if } i \notin \{y_1, y_2\} \\ z_i^t - 2\eta' \pi_i^t + \eta' & \text{if } i \in \{y_1, y_2\} \end{cases} \\ \text{REINFORCE: } z_i^{t+1} &= \begin{cases} z^t - \eta' (\Delta \pi^t \cdot \pi_i^t) & \text{if } i \notin \{y_1, y_2\} \\ z^t - \eta' (\Delta \pi^t - 1) \pi_i^t & \text{if } i \in \{y_1, y_2\} \end{cases} \end{split}$$

Then, we can write down the probability change  $\pi_i^{t+1}$  for each action i. For case  $i \in \{y_1, y_2\}$ , we have the following derivation for MLE objective:

$$\pi_{y_1}^t = \frac{e^{z_i^{t+1}}}{\sum_{j=1}^V e^{z_j^{t+1}}} = \frac{e^{z_i^t - 2\eta' \pi_i^t + \eta'}}{\sum_{j \notin \{y_1, y_2\}} e^{z_j^t - 2\eta' \pi_j^t} + e^{z_{y_1}^t - 2\eta' \pi_{y_1}^t + \eta'} + e^{z_{y_2}^t - 2\eta' \pi_{y_2}^t + \eta'}}$$
(25)

$$= \frac{e^{z_i^t}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^t - \eta' \left(2(\pi_j^t - \pi_i^t) + 1\right)} + e^{z_{y_1}^t - 0} + e^{z_{y_2}^t - 2\eta' \left(\pi_{y_2}^t - \pi_i^t\right)}}$$
(26)

Similarly, for REINFORCE objective,

$$\pi_{y_1}^t = \frac{e^{z_i^{t+1}}}{\sum_{j=1}^V e^{z_j^{t+1}}} = \frac{e^{z_i^t - \eta'(\Delta \pi^t - 1)\pi_i^t}}{\sum_{j \notin \{y_1, y_2\}} e^{z_j^t - \eta'(\Delta \pi^t \cdot \pi_j^t)} + e^{z_{y_1}^t - \eta'(\Delta \pi^t - 1)\pi_{y_1}^t} + e^{z_{y_2}^t - \eta'(\Delta \pi^t - 1)\pi_{y_2}^t}}$$
(27)

$$= \frac{e^{z_i^t}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^t - \eta' \left(\Delta \pi^t (\pi_j^t - \pi_i^t) + \pi_i^t\right)} + e^{z_{y_1}^t - 0} + e^{z_{y_2}^t - \eta' \left((\Delta \pi^t - 1)(\pi_{y_2}^t - \pi_i^t)\right)}}$$
(28)

For case  $i \notin \{y_1, y_2\}$ , we have the following derivation for MLE objective:

$$\pi_i^t = \frac{e^{z_i^{t+1}}}{\sum_{j=1}^V e^{z_j^{t+1}}} = \frac{e^{z_j^{t} - \eta'(2\pi_i^t)}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^{t} - 2\eta'\pi_j^t} + e^{z_{y_1}^{t} - 2\eta'\pi_{y_1}^t + \eta'} + e^{z_{y_2}^{t} - 2\eta'\pi_{y_2}^t + \eta'}}$$
(29)

$$= \frac{e^{z_i^t}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^t - 2\eta' \left(\pi_j^t - \pi_i^t\right)} + e^{z_{y_1}^t - \eta'(2(\pi_{y_1}^t - \pi_i^t) - 1)} + e^{z_{y_2}^t - \eta'(2(\pi_{y_2}^t - \pi_i^t) - 1)}}$$
(30)

Similarly for REINFORCE objective:

$$\pi_i^t = \frac{e^{z_i^{t+1}}}{\sum_{j=1}^V e^{z_j^{t+1}}} = \frac{e^{z_j^{t} - \eta'(\Delta \pi^t \cdot \pi_j^t)}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^{t} - \eta'(\Delta \pi^t \cdot \pi_j^t)} + e^{z_{y_1}^t - \eta'(\Delta \pi^t - 1)\pi_{y_1}^t} + e^{z_{y_2}^t - \eta'(\Delta \pi^t - 1)\pi_{y_2}^t}}$$
(31)

$$= \frac{e^{z_i^t}}{\sum\limits_{j \notin \{y_1, y_2\}} e^{z_j^t - \eta' \left(\Delta \pi^t (\pi_j^t - \pi_i^t)\right)} + e^{z_{y_1}^t - \eta' \left(\Delta \pi^t (\pi_{y_1}^t - \pi_i^t) - \pi_{y_1}^t\right)} + e^{z_{y_2}^t - \eta' \left(\Delta \pi^t (\pi_{y_2}^t - \pi_i^t) - \pi_{y_2}^t\right))}}$$
(32)

We can now better understand how each  $\pi_i$  changes after one update. Specifically, if  $\alpha_i > 1$ , the corresponding  $\pi_i$  increases, and vice versa. To determine the value of  $\alpha_i$ , we can treat any  $\beta_j > 1$  as contributing to the conclusion that  $\alpha_i < 1$ , while any  $\beta_j < 1$  against it. The value of the corresponding  $e^{z_j^t}$  and  $|\beta_j - 1|$  controls how strong the contribution is.

Low-likelihood tokens tend to receives less meaningful updates. Using the log-derivative trick,  $\nabla_{\theta}\pi_{y}=\pi_{y}\nabla_{\theta}\log\pi_{y}$ , we see that REINFORCE inherently applies a conservative update rule: tokens with higher probability receive proportionally larger gradient updates, while low-probability tokens are updated much more weakly. This disproportionate updating can have a detrimental effect in cases where optimal actions may reside in low-probability regions of the distribution, yet RE-INFORCE fails to provide sufficient reinforcement for these tokens to be effectively explored and learned.

Assuming  $\pi^t_{y_1} > \pi^t_{y_2}$ , then  $\pi^t_{y_1}$  is guaranteed is increase, i.e.,  $\alpha_{y_1} > 1$ . To see this, consider the value of  $\gamma$  in Case 1 of the REINFORCE objective. For any  $j \notin y_1, y_2$ , we have  $\gamma_j < 1$  because

$$\Delta \pi^{t} (\pi_{j}^{t} - \pi_{i}^{t}) + \pi_{i}^{t} = \Delta \pi^{t} \pi_{j}^{t} + (1 - \Delta \pi^{t}) \pi_{i}^{t} > 0.$$
(33)

Moreover, since  $\gamma_{y_1}=1$  while  $\gamma_{y_2}<1$ , the probability of the highest-likelihood action  $y_1$  will increase. In addition, because the update step size depends on the probability of the sampled action, the smaller  $\pi^t_{y_1}$  is, the more strongly  $\pi^t_{y_2}$  will decrease, leading the model to gradually concentrate its mass on the dominant mode  $\pi^t_{y_1}$ . However, this is not the case for the maximum likelihood objective; later, as we show in section 2, the maximum likelihood objective will tend to increase the lower probability  $\pi^t_{y_2}$ , avoiding model collapse.

#### D.1 EMPIRICAL VERIFICATION WITH A SIMPLE BANDIT PROBLEM.

To empirically verify the analysis above, we analyze using a simple logistic regression task following Ren & Sutherland (2025). We set the vocabulary size V=50, with parameter dimension d=5, a learning rate of  $\eta=0.5$ , and a randomly generated input vector  $\boldsymbol{x}$ . We compared the update effect of MLE and REINFORCE objectives in 2 scenarios: when the initial distribution  $\pi_{\theta^0}$  is relatively uniform and when it's highly peaked around certain dimensions.

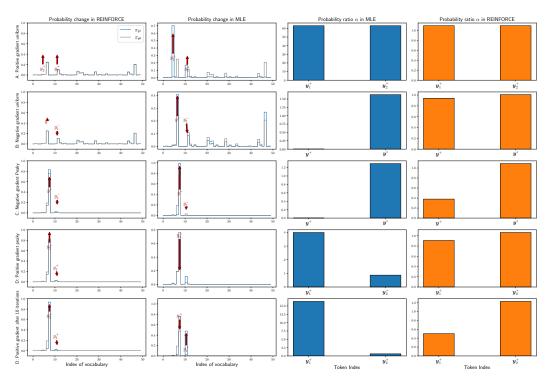


Figure 9: Experimental verification and comparison between REINFORCE and MLE objectives using a simple contextual bandit problem.

Uniform case: we consider negative and positive gradient effect, in the negative gradient case, we consider negative label  $y^-=11$  and investigate the "indirect" push-up effect on  $y_2=5$ . In the positive gradient scenario, we use 2 positive actions  $y_1=11, y_2=5$  for the SGD update. As demonstrated in the first 2 rows in Fig. 9, both REINFORCE and MLE behave similarly in both negative and positive gradient scenarios, where both  $y_1, y_2$  are increased. However, MLE aggressively imposes a large update, where the probability is approximately 60 times larger on these two tokens after the update, as measured by the probability ratio  $\alpha$ . Similarly, a negative gradient tends to increase the highest-likelihood token  $y^*$ , while decreasing other actions, including the positive action  $y_2$ . This can be explained by REINFORCE by imposing a probability step-size, resulting in a more conservative update, where the models less deviate from their previous step compared to the MLE objective.

**Peaky case:** In LLM regime, the LMs usually initialize from a pre-trained checkpoint, where the model exhibits a non-uniform distribution. We consider the scenario where the models concentrate on a few dimensions in action distribution. Similarly, negative label  $y^- = 11$  and investigate the

"indirect" push-up effect on  $y_2 = 8$ . In the positive gradient scenario, we use 2 positive actions  $y_1 = 11, y_2 = 8$  for the SGD update. As shown in the last 3 rows in Fig. 9, when the distribution is peaky, the negative gradient exacerbates the effect of increasing the highest probability token while decreasing other actions, regardless of correctness Ren & Sutherland (2025). Interestingly, in a positive gradient scenario, between the 2 positive actions, MLE tends to push up the actions with the lowest likelihood (as measured by  $\alpha$ ). In the last row, we observe that after 10 iterations of training, MLE can converge to a new distribution with equal probabilities between 2 actions. However, REINFORCE shows favor for the actions with the highest probability; this effect will exacerbate after multiple updates, where it collapses the highest probability  $\pi(y_1 = 8)$ . This suggests that offline learning can still provide benefits compared to REINFORCE in avoiding model collapse.

# E ADDITIONAL RESULTS ON DIFFERENT MODELS

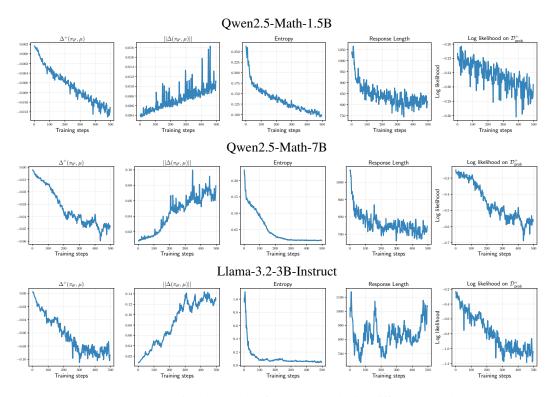


Figure 10: Learning dynamics of RLVR model with different models.

Consistent learning dynamics across different models. In this section, we provide additional results that our findings hold across different model families and sizes. Specifically, we conduct GRPO experiments on a larger model, Qwen2.5-Math-7B, as well as a different family, Llama-3.2-3B-Instruct. As shown in Fig. 10, our analysis consistently holds across these settings.

**Larger models tend to suffer more from** *negative interference*. Interestingly, we find that the negative interference phenomenon becomes more severe when scaling up to larger models such as Qwen2.5-Math-7B. This also helps explain why, under RLVR training, the Pass@k performance of Qwen2.5-Math-7B with large sampling budgets k can even underperform than the smaller Qwen2.5-Math-1.5B after GRPO.

**Detailed Pass**@k performance. We present detailed Pass@k curves for Qwen2.5-Math-1.5B and Qwen2.5-Math-7B of different methods below:

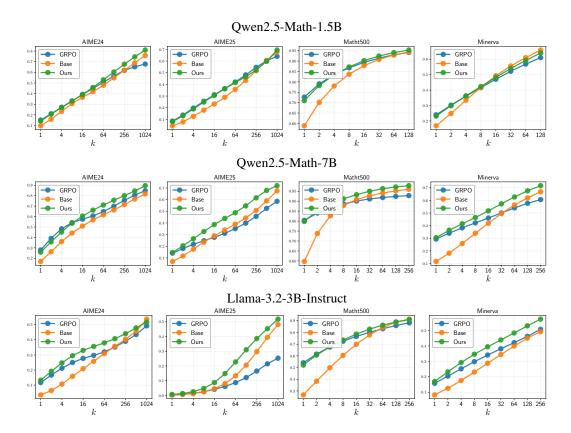


Figure 11: Pass@k curves for base model, RLVR model and SELF with different models.

#### F PERFORMANCE DEGRADATION IN W-REINFORCE

We also experiment with recent work **W-REINFORCE** Zhu et al. (2025), which upweights the negative gradient signal to encourage diversity. As shown in Fig. 2, we find that **W-REINFORCE** suffers significantly from the catastrophic forgetting issue, where the average accuracy decreases even before 300 steps. We hypothesize that it is due to the high variance of the REINFORCE objections.

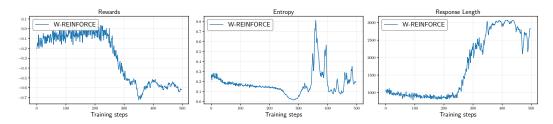


Figure 12: Learning dynamics of **W-REINFORCE**, we observe that **W-REINFORCE** suffer from training instability under prolong training.

tive. Furthermore, progressively learning with a negative gradient can result in improper increases in tail probabilities Li et al. (2025b), leading the models to generate non-meaningful tokens.

# G COMPUTATIONAL COST OF SELF

Our method require an additional greedy generation to filter highly solvable problems. To analyze the computational cost of SELF and GRPO, we utilize per-iteration time of each update step with 3 factors: total generation time  $T_{\rm gen}$ , backward time  $T_{\rm backward}$ , and forward time  $T_{\rm forward}$ . We provide

our results in Fig. 13 We find that SELF incurs an additional cost in generation time compared to

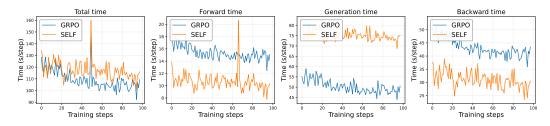


Figure 13: Computational results for training <code>Qwen2.5-Math-1.5B</code> for 100 steps.

GRPO. However, due to the filtering mechanism, both the forward and backward passes are faster, which compensates for the extra generation time and results in a negligible overall computational cost relative to GRPO.

# H PERPLEXITY DETAIL ANALYSIS

In this section, we provide a detailed perplexity analysis. Across all benchmarks, we randomly sample 16 reasoning traces for each problem in the test dataset to calculate perplexity.

To collect problems that increased after RLVR  $\mathcal{D}^{\uparrow}$ , we estimate the average accuracy for each problem and sort the problem indices according to average accuracy changes between the base policy and RLVR policy:

$$\Delta r(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\boldsymbol{y}|\boldsymbol{x})} \left[ r(\boldsymbol{x}, \boldsymbol{y}) \right] - \mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\boldsymbol{y}|\boldsymbol{x})} \left[ r(\boldsymbol{x}, \boldsymbol{y}) \right]$$
(34)

For AIME 2024/2025, we sample the top-3 problems that increased the most for  $\mathcal{D}^{\uparrow}$  and top-3 problems decreased for  $\mathcal{D}^{-}$ , as these benchmarks only consist of 30 problems. For Math500 and Minerva, we sample top-64 problems for both  $\mathcal{D}^{+}$  and  $\mathcal{D}^{\downarrow}$ . We present out results in Fig. 14.

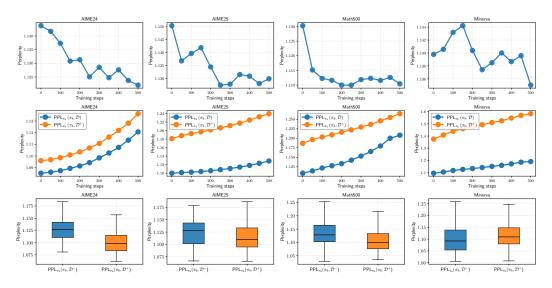


Figure 14: Dynamics of Perplexity during RLVR. The first row shows that RLVR increasingly generates low-perplexity responses under the base model. The second row demonstrates that the initial base solutions, regardless of correctness, exhibit increasing perplexity; the LMs are less likely to generate initial, diverse solutions. The final row shows that for problems where the correct solution already has a high probability to generate under base model, RLVR will incentivize, improving these problems, while problems with low success rate exhibits performance degradation.

# I ADDITIONAL DETAILS IN REDUCING PREVIOUSLY LEARNED BEHAVIORS IN RLVR.

Similar to Shao et al. (2025), we define *code reasoning* are responses that contain the string python. For each prompt x, the LMs policy can generate either *code reasoning traces* or *lang reasoning traces*. We can define a mixture of distribuion between these two reasoning traces:

$$\pi_b(\boldsymbol{y}|\boldsymbol{x}) = \alpha(\boldsymbol{x})\pi_{\text{code}}(\boldsymbol{y}|\boldsymbol{x}) + (1 - \alpha(\boldsymbol{x}))\pi_{\text{lang}}(\boldsymbol{y}|\boldsymbol{x})$$
(35)

where  $\alpha(x)$  is the mixture coefficient that depends on the prompt. For each prompt x, we sample k=144 responses to estimate the average accuracy of the mixture distributions  $\pi_{\rm code}$  and  $\pi_{\rm lang}$ , denoted as  $\rho_{\rm code}(x)$  and  $\rho_{\rm lang}(x)$ . We then define their average accuracies on the test set  $\mathcal{D}_{\rm test}$  as:

$$\mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}_{\text{test}}} \left[ \rho_{\text{code}}(\boldsymbol{x}) \right], \quad \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}_{\text{test}}} \left[ \rho_{\text{lang}}(\boldsymbol{x}) \right]. \tag{36}$$

We report results for each reasoning behavior in the <code>Qwen2.5-Math</code> models across all benchmarks below:

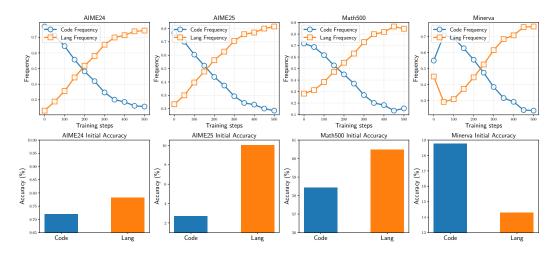


Figure 15: Reasoning strategy switching from code to lang throughout training process. While AIME24/25 and Math500 exhibits initial better accuracy in lang than code, Minerva shows better performance in *code reasoning traces*.

We observe that the models eventually switch almost entirely to *lang reasoning traces*. We hypothesize that this shift explains the improvements on AIME24/25 and Math500, where *lang reasoning traces* is more effective initially. However, the loss of *code reasoning traces* reduces performance on Minerva, where the base model originally benefited from producing *code reasoning traces*.

#### J USAGE OF LARGE LANGUAGE MODELS (LLMS)

We used ChatGPT solely for revising the writing of the paper. Note that revision here strictly means enhancing the clarity and readability of the text (e.g., fixing typos or constructing latex tables), and not for any other purposes.