

Interpretable Intrinsic Cues for Efficient Reinforcement Learning with Large Language Models

Anonymous ACL submission

Abstract

Reinforcement learning with verifiable rewards (RLVR) has improved the reasoning ability of large language models, yet training remains costly because many rollouts contribute little to optimization relative to their heavy computational demands. This study investigates how simply leveraging interpretable and intrinsic data properties, which come at almost no additional computational cost during training, can markedly improve data efficiency for RLVR. We propose PREPO¹, an RLVR model with two complementary components. First, we use prompt perplexity as a proxy for model adaptability in learning, and adopt a schedule to guide the model from well-understood prompts to progressively challenging ones. Second, we amplify the diversity among rollouts by differentiating their relative entropy and prioritizing sequences with greater exploratory behavior. Together, these mechanisms reduce rollout demand while preserving competitive performance. On Qwen and Llama models, PREPO achieves effective results on mathematical reasoning benchmarks with up to $3\times$ fewer rollouts than baselines. Beyond empirical gains, we provide theoretical and in-depth analyses that explain how our method improves the data efficiency of RLVR.

1 Introduction

Reinforcement learning (RL) has become central in improving the reasoning capabilities of large language models (LLMs) by optimizing self-generated rollouts [Guo et al., 2025, Team et al., 2025, Chen et al., 2025a]. Recent advances in reinforcement learning with verifiable reward (RLVR) demonstrate that it is a simple yet effective method for scaling reasoning performance [Shao et al., 2024, Yu et al., 2025]. However, applying RLVR to models that generate long reasoning traces will incur substantial computational overhead in the

rollout stage, significantly hampering RL training efficiency and becoming the primary training bottleneck [Zhong et al., 2024].

The exploration of effective strategies for leveraging data in RL training remains relatively underdeveloped. Prior studies [Zhang et al., 2025, Albalak et al., 2025] have suggested using pass rate as an indicator of data difficulty to strengthen the training signal. Nonetheless, this approach often requires multiple rounds of sampling to attain a sufficiently stable measurement. In fact, it converts online rollouts into an offline context, which ultimately does not reduce the overall computational burden. Other metrics, such as human-defined criteria [Chen et al., 2025b, Parashar et al., 2025], *e.g.*, specific domains or topics, have the disadvantage of being influenced by the tagging process and by biases in human perceptions and experiences. In addition, alternative approaches often depend on auxiliary trained models or embedding techniques to reflect data semantics. We assert that these techniques not only impose significant computational and memory costs but also introduce mismatches and limited applicability across diverse policy models and training paradigms. Furthermore, they overlook the inherent dynamism of the RL training process, lagging far behind the pace of training updates. As a remedy, other methods incorporate historical information from training to reflect the dynamic nature of data. Nevertheless, it increases memory requirements for the RL framework and introduces extraneous noise during training.

A natural way to improve training efficiency is to select data, *i.e.*, prune uninformative prompts or rollouts while preserving those that drive learning. There are emerging approaches based on parameterized modeling [Qu et al., 2025], replay buffers [Liu et al., 2025], or selective rollout execution [Zheng et al., 2025]. Instead of being aided by inductive biases from both humans and external models, we address this long-standing research question from

¹We will release data and codes soon upon acceptance.

083 a new perspective:

084 *Can the intrinsic data properties deriving from*
085 *the training process improve the efficiency of*
086 *RLVR?*

087 In this study, we propose a simple method
088 with almost negligible computation cost,
089 Perplexity-schedule with Relative-Entropy Policy
090 Optimization (PREPO), a method that combines
091 a perplexity-based schedule with sequence-level
092 entropy weighting to realize *intrinsic exploration*.
093 Specifically, PREPO traces perplexity *before*
094 rollout generation to prune the prompts, and
095 applies entropy weighting *after* rollout generation
096 to emphasize uncertain responses. It is worth
097 noting that metrics collected during standard
098 RL training can be reused to compute both
099 components, thereby ensuring computational
100 efficiency of our method. Moreover, our method is
101 coherently integrated with the policy model and
102 training process, offering a favorable trade-off
103 between suitability and flexibility. Beyond that,
104 our approach uncovers the intrinsic nature of
105 data during RL training, providing fine-grained
106 interpretability of training dynamics. Across Qwen
107 and Llama models, PREPO surpasses existing
108 data-pruning baselines and remains competitive
109 with the baseline, while reducing rollout usage by
110 more than 40% (see Fig. 1 for Qwen2.5-Math-7B).
111 These results show that RLVR can be made
112 substantially more efficient by leveraging the
113 intrinsic properties of prompt and rollout data.

114 2 Related Work

115 **Data efficiency in RLVR.** A growing body of
116 work has explored data efficiency for RLVR, with
117 particular attention to online data selection. Unlike
118 offline methods that require pretraining or costly
119 rollouts to estimate sample quality [Qu et al., 2025,
120 Zhang et al., 2025, Chen et al., 2025b, Kamaloo
121 et al., 2025], online approaches aim to reduce over-
122 head by dynamically filtering or prioritizing sam-
123 ples during training. Online difficulty filtering [Bae
124 et al., 2025] removes prompts that contribute little
125 to reasoning improvement, while predictive prompt
126 allocation [Qu et al., 2025] directs rollouts toward
127 more promising inputs. Curriculum-based strate-
128 gies further adapt training to the model’s evolving
129 competence [Zhang et al., 2025, Chen et al., 2025b].
130 Other online selection approaches prioritize sam-
131 ples using gradient-informed signals [Kamaloo
132 et al., 2025, Chen et al., 2025c] or policy-advantage

estimates [Wang and Guofeng, 2025]. Parallel ef-
083 orts focus on reducing rollout redundancy during
084 training. Down-sampling strategies [Li et al., 2025]
085 and efficient replay buffer designs [Liu et al., 2025]
086 lessen the burden of repeatedly training on uninfor-
087 mative samples. Collectively, these online methods
088 emphasize the importance of allocating computa-
089 tional resources to samples that drive progress in
090 reasoning, a goal that aligns with the direction of
091 our work. 142

Entropy Mechanism in RLVR. Entropy has
143 long been studied in reinforcement learning
144 through entropy-regularized objectives that pro-
145 mote exploration in control settings. Recent stud-
146 ies extend this idea to RLVR for reasoning LLMs.
147 Cui et al. [2025] identify rapid entropy collapse
148 as a major failure mode and propose covariance-
149 based updates to slow its decay. Wang et al. [2025]
150 show that high-entropy “forking tokens,” though
151 rare, account for most reasoning gains, highlight-
152 ing entropy as a token-level signal of informativ-
153 ness. Cheng et al. [2025] incorporates a clipped,
154 gradient-detached entropy term into the advantage
155 function, encouraging more exploratory responses
156 but introducing additional hyperparameters. Build-
157 ing on these insights, we propose a parameter-free
158 approach to reinforce entropy-driven exploration
159 in RLVR. A related line of work explores weight-
160 ing in policy-gradient updates through truncated
161 importance sampling, such as CISPO [Chen et al.,
162 2025a]. These methods address a different dimen-
163 sion of RLVR optimization, and can be viewed as
164 complementary to our focus on online data selec-
165 tion and entropy-based weighting. 166

167 3 Preliminary Analysis

168 3.1 Low-PPL Prompts Tend to Yield Higher 169 Pass Rate

170 We begin by examining the relationship between
171 prompt perplexity (PPL) and task difficulty us-
172 ing the DAPO-Math-17K dataset [Yu et al., 2025]
173 (17,917 samples). For both Qwen and Llama mod-
174 els, Figure 2 shows a clear negative correlation be-
175 tween PPL and passrate@16, where passrate@16
176 measures the fraction of prompts solved by at least
177 one of 16 generations. Lower-PPL prompts gener-
178 ally yield higher success rates. Table 1 correlation
179 is statistically significant across models, suggest-
180 ing that PPL can serve as a lightweight signal to
181 identify more informative prompts for training.

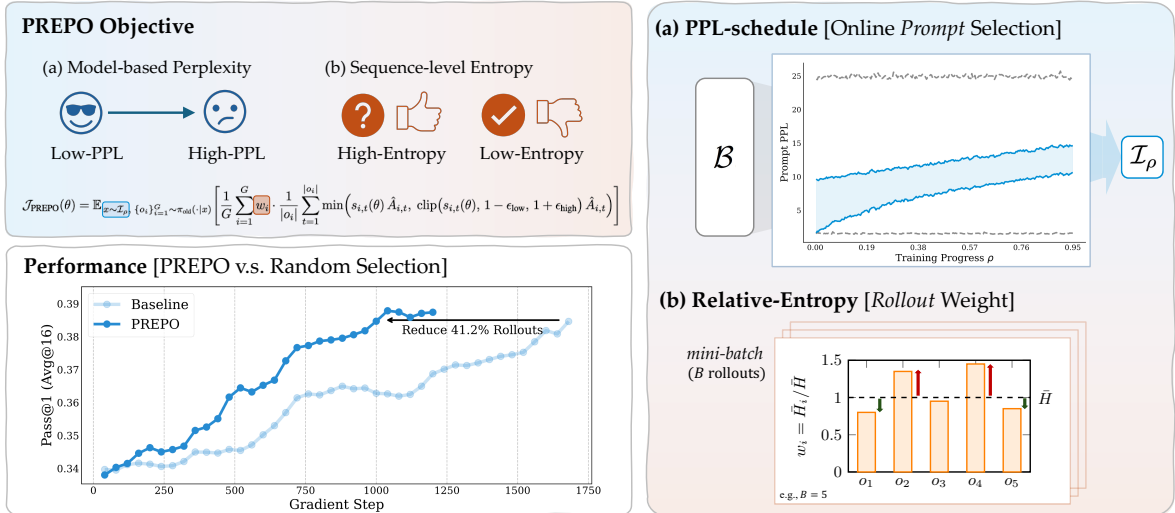


Figure 1: **Overview of PREPO.** The PREPO objective integrates perplexity-based schedule learning and sequence-level entropy weighting into a unified optimization scheme. On Qwen2.5-Math-7B, PREPO achieves higher performance while requiring only 41.2% of the rollouts used by random selection, showing improved efficiency. Specifically, PREPO has two complementary components: (a) **PPL-schedule**, which actively selects prompts according to model-based perplexity, starting with lower-PPL prompts (“adapted” to the model) and progressively introducing higher-PPL ones (“obscured” to the model) as training progresses. (b) **Relative-Entropy Weighting**, which adjusts rollout contributions by comparing each sequence’s entropy against the batch average, amplifying the high-entropy rollouts (“novel attempts”) and downweighting the low-entropy (“regular response”).

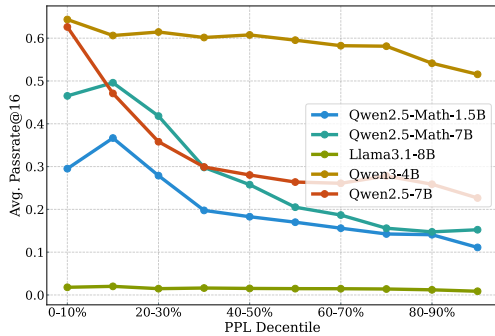


Figure 2: Prompt PPL versus average passrate@16.

Table 1: Correlation between prompt PPL and pass-rate@16. (***) $p < 0.001$, (**) $p < 0.05$)

Qwen2.5-7B	Qwen2.5-7B	Qwen2.5-1.5B	Qwen3-4B	LLama3.1-8B
-0.233***	-0.183**	-0.186***	-0.169***	-0.199***

3.2 Training Dynamics of LOW-PPL and HIGH-PPL Prompts

To better understand the role of PPL during training, we compare prompts from the lowest 20% (LOW-PPL) and highest 20% (HIGH-PPL) of the distribution (see Appendix B for examples of the two groups). Figure 3 illustrates their training dynamics on Qwen2.5-Math-7B. The two groups exhibit complementary behavior: LOW-PPL prompts

drive rapid improvements in reward and validation accuracy during early training, though at the cost of faster entropy collapse, whereas HIGH-PPL prompts preserve entropy and lower zero-advantage ratio, yielding stronger performance in later stages.

These trends are consistent across other Qwen models (Appendix C). Specifically, (a) HIGH-PPL prompts are associated with higher entropy, (b) LOW-PPL prompts yield higher rewards and validation accuracy early on, and (c) HIGH-PPL prompts maintain exploration and ultimately close the performance gap. For Llama3.1-8B (Figure 10), a similar pattern appears, though overall performance remains limited as the validation dataset is too challenging for Llama models.

3.3 Comparison with Random Sampling

To test whether PPL-based grouping offers value beyond chance, we also compare with a random 20% subset. As shown in Figure 4, the random group consistently falls between the LOW-PPL and HIGH-PPL groups. This indicates that PPL is a non-trivial, policy-intrinsic signal that distinguishes easy prompts with high-pass rates from more challenging ones, providing a valuable basis for the design of online batch selection.

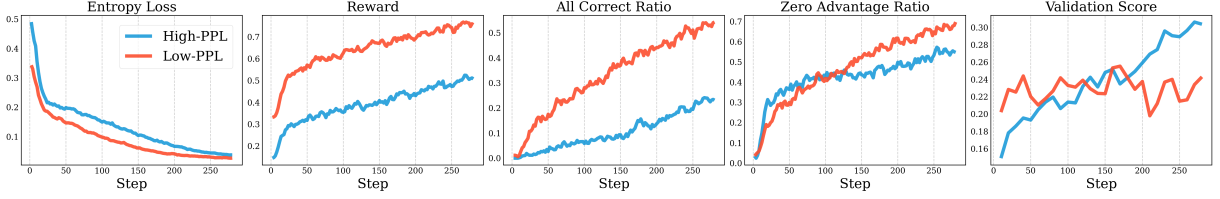


Figure 3: Training dynamics of LOW-PPL vs. HIGH-PPL prompts on Qwen2.5-Math-7B. (a) HIGH-PPL prompts have higher entropy. (b) LOW-PPL prompts have more reward gains. (c) LOW-PPL prompts reach higher all-correct ratios faster. (d) LOW-PPL prompts show higher zero-advantage ratios in the later stage. (e) HIGH-PPL prompts eventually outperform LOW-PPL prompts.

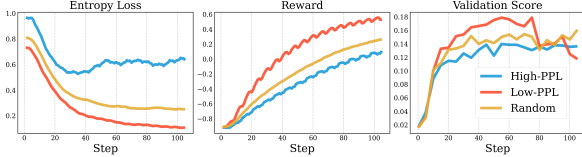


Figure 4: Comparison among LOW-PPL, HIGH-PPL, and Random Subsets. *Random* lies between the two, showing that PPL-based grouping provides a meaningful pruning signal.

4 PREPO: PPL-Schedule Relative-Entropy Policy Optimization

Building on the preliminary results, PREPO integrates a perplexity-driven schedule with entropy-based rollout weighting, forming a unified framework to improve efficiency in RLVR.

4.1 General Online Batch Selection

Let $\mathcal{B} = \{x_i\}_{i=1}^N$ denote the candidate batch at a training step. The goal of online batch selection is to design a mapping

$$\Phi : [0, 1] \rightarrow 2^{\mathcal{B}}, \quad \rho \mapsto \mathcal{I}_\rho, \quad (1)$$

where $\rho \in [0, 1]$ denotes the normalized training progress, and $\mathcal{I}_\rho \subseteq \mathcal{B}$. The mapping Φ is required to (i) explicitly depend on ρ , so that the distribution of selected samples evolves with training; (ii) the sub-batch size is fixed during training, i.e., $\forall \rho, |\mathcal{I}_\rho| = K$.

4.2 PPL-Schedule Online Batch Selection

For a prompt $x_i = (x_{i,1}, \dots)$, we measure its perplexity under the policy π_ρ at training progress ρ as

$$P_i(\rho) = \exp \left(-\frac{1}{|x_i|} \sum_{t=1}^{|x_i|} \log \pi_\rho(x_{i,t} | x_{i,<t}) \right), \quad (2)$$

where π_ρ is the model distribution at progress ρ . As π_ρ is parameterized by θ and evolves throughout

training, $P_i(\rho)$ provides a model-based measure for active data selection. We then define the PPL-schedule sub-batch as

$$\mathcal{I}_\rho = \{ \sigma(j) : l(\rho) \leq j \leq l(\rho) + K - 1 \}, \quad (3)$$

where σ is the permutation that sorts \mathcal{B} by ascending $P_i(\rho)$. The starting index $l(\rho)$ is given by a linear schedule

$$l(\rho) = \lfloor \rho \cdot (N - K) \rfloor, \quad (4)$$

so that \mathcal{I}_ρ shifts smoothly from LOW-PPL to HIGH-PPL prompts. While linear scheduling is the simplest case, a nonlinear² can also be used. In general, the PPL-schedule serves as an online data-selection procedure that shifts from more to less in-domain prompts as training progresses.

4.3 Relative Entropy Weighting

As shown in Section 3.2, we empirically find that training on LOW-PPL prompts accelerates reward improvement but also leads to a rapid collapse of entropy, thereby reducing exploration. To mitigate this effect during the PPL schedule, we introduce a sequence-level relative-entropy weighting scheme that adaptively emphasizes uncertain rollouts.

The token-level entropy of a rollout is defined as $H_t = -\sum_{v \in \mathcal{V}} \pi_\theta(v | o_{<t}, x) \log \pi_\theta(v | o_{<t}, x)$, where \mathcal{V} is the vocabulary. For rollout i , the sequence-level entropy is the average across its tokens

$$\bar{H}_i = \bar{H}(o_i | x) = \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} H_t. \quad (5)$$

The batch-average entropy over B rollouts is

$$\bar{H} = \frac{1}{B} \sum_{k=1}^B \bar{H}_k. \quad (6)$$

²We define a general family of nonlinear schedules as $l(\rho) = \lfloor \rho^\alpha \cdot (N - K) \rfloor$, where $\alpha \in \mathbb{R}^+$. For example, choosing $\alpha = \frac{1}{2}$ yields a square-root schedule that transitions more quickly toward higher-PPL prompts. One (e.g., quadratic or exponential)

The relative weight assigned to rollout i is then given by

$$w_i = \frac{\bar{H}_i}{\bar{H}}. \quad (7)$$

This formulation is scale-invariant, as a rollout’s contribution depends only on its entropy relative to the batch mean. Intuitively, this design enables the model to *seek uncertainty within certainty* during the PPL schedule. While LOW-PPL prompts early in training often yield confident (low-entropy) responses, relative weighting amplifies the impact of less confident (higher-entropy) rollouts, thereby preserving exploration throughout training.

4.4 Objective Function

The PREPO objective integrates PPL-schedule filtering with relative-entropy weighting as below.

$$\mathcal{J}_{\text{PREPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{I}_\rho, \{o_i\}_{i=1}^G \sim \pi_{\text{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G w_i \cdot \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \min \left(s_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(s_{i,t}(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}) \hat{A}_{i,t} \right) \right] \quad (8)$$

where \mathcal{I}_ρ is the PPL-schedule-filtered batch of prompts at training progress ρ , w_i encodes the relative entropy of rollout i at the current micro-batch, $s_{i,t}(\theta)$ is the token-level importance ratio, $s_{i,t}(\theta) = \frac{\pi_\theta(o_{i,t}|x, o_{i,<t})}{\pi_{\text{old}}(o_{i,t}|x, o_{i,<t})}$, and $\hat{A}_{i,t}$ is the group-based advantage estimate, $\hat{A}_{i,t} = \frac{r_i - \text{mean}(\{r_j\}_{j=1}^G)}{\text{std}(\{r_j\}_{j=1}^G)}$.

5 Experiments

5.1 Setups

We benchmark PREPO against three baselines: Random selection, Dynamic Sampling (DS) [Yu et al., 2025], and GRESO [Zheng et al., 2025] (detailed in Appendix D). Our evaluation covers multiple models: Qwen2.5-7B [Team, 2024], Qwen2.5-Math-1.5B and Qwen2.5-Math-7B [Yang et al., 2024], Qwen3-4B (non-thinking) [Yang et al., 2025], and Llama3.1-8B [Dubey et al., 2024]. For training data, we use DAPO-Math-17K [Yu et al., 2025] and MATH500 [Lightman et al., 2023a] with a total of 18,417 training samples.

Training and Evaluation. All models are trained using the ver1 [Sheng et al., 2025], with vLLM [Kwon et al., 2023] employed for rollout generation to ensure efficient inference. For the Qwen

models, we evaluate them on four benchmarks, including AIME25 [Art of Problem Solving, 2025] (30 samples), AIME24 [Art of Problem Solving, 2024] (30 samples), MATH500 [Lightman et al., 2023b] (500 samples), and OlympiadBench [He et al., 2024] (8,476 samples), which cover a diverse range of mathematical reasoning challenges. The Llama model is evaluated on MATH500 [Lightman et al., 2023b] and GSM8K [Cobbe et al., 2021] (1,319 samples). We evaluate all models using *pass@1 (avg16)*, i.e., the accuracy of the top-1 response averaged over 16 generations, with temperature 1. We evaluate models every 50 training steps and report the best average performance on all benchmarks.

Experiment Configuration. For the Qwen2.5-Math models, we use a maximum context length of 4096 tokens, matching their supported limit. For Qwen3-4B and Llama3.1-8B, we set the context length to 32,768 tokens. Rollouts are generated with temperature as 1 using vLLM, producing 8 responses per prompt. For PREPO, Random, and GRESO, we adopt an online selection ratio of $K/N = 20\%$ at each training step, where the candidate batch size (N) is 1280 and the actual batch size (K) is fixed at 256. For GRESO, we set the targeted zero-variance percentage as 50%. For DS, the candidate batch size is 384. The mini-batch size (B) is 64 for all experiments. For all experiments, we set the clipping thresholds to $\epsilon_{\text{low}} = 0.2$ by default, with a larger $\epsilon_{\text{high}} = 0.28$ as the upper bound. The actor model is optimized with AdamW using a constant learning rate of 1×10^{-6} , momentum parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$, and a weight decay of 0.01. Following Yu et al. [2025], we omit the KL-divergence regularization term. Training is applied only to the actor parameters and parallelized with Fully Sharded Data Parallel. All experiments are conducted on 32 GPUs.

5.2 Results

PREPO achieves consistent two- to three-fold rollout reduction. As shown in Table 2 below, PREPO substantially lowers rollout usage relative to all baselines while maintaining or surpassing accuracy. On Qwen2.5-7B, PREPO reduces rollouts from over 1M (DS: 1,040K; Random: 716K; GRESO: 680K) to 304K. On Qwen2.5-Math-1.5B, it cuts the budget from several million (DS: 3.6M; Random: 3.0M; GRESO: 2.5M) down to 1.1M. Similarly, PREPO reduces rollouts to 540K on

Table 2: Performance comparison (%) on Qwen models. For each model, the top row reports the base model’s performance. Best results are highlighted in bold or underlined. K = thousand, M = million.

Method	AIME25	AIME24	MATH	Olympiad	Avg ↑	# Rollouts ↓
<i>Qwen2.5-7B</i>	1.25	4.17	72.26	33.09	27.69	–
+ DS	7.92	<u>17.55</u>	75.50	38.91	34.97	1040K
+ Random	6.98	16.41	75.70	38.47	34.39	716K
+ GRESO	9.22	10.83	<u>76.65</u>	<u>42.07</u>	34.59	680K
+ PREPO (Ours)	<u>10.21</u>	16.09	76.30	39.85	35.61	304K
<i>Qwen2.5-Math-1.5B</i>	3.54	10.21	55.76	27.41	24.23	–
+ DS	10.83	<u>25.83</u>	76.40	23.33	34.10	3.6M
+ Random	<u>20.00</u>	16.67	76.25	30.50	35.86	3.0M
+ GRESO	15.38	20.00	<u>76.65</u>	24.17	34.16	2.5M
+ PREPO (Ours)	<u>20.00</u>	16.67	76.25	<u>32.00</u>	36.23	1.1M
<i>Qwen2.5-Math-7B</i>	9.17	20.80	72.26	39.56	35.45	–
+ DS	13.33	<u>33.33</u>	<u>81.35</u>	30.17	39.55	1664K
+ Random	10.00	26.67	77.80	<u>43.26</u>	39.45	905K
+ GRESO	<u>18.33</u>	25.83	77.80	26.83	37.46	654K
+ PREPO (Ours)	12.81	26.15	77.85	41.58	39.59	540K
<i>Qwen3-4B</i>	30.00	53.33	94.10	52.67	57.53	–
+ DS	63.33	66.67	95.10	58.00	70.78	688K
+ Random	60.00	70.00	96.00	59.33	71.33	553K
+ GRESO	56.67	69.17	96.40	57.33	69.89	472K
+ PREPO (Ours)	<u>66.67</u>	<u>80.00</u>	<u>96.60</u>	<u>60.67</u>	75.99	348K

Qwen2.5-Math-7B (vs. 1.66M for DS, 905K for Random, 654K for GRESO) and to 348K on Qwen3-4B (vs. 688K, 553K, and 472K, respectively).

Dynamic sampling and random selection are inefficient. Although DS sometimes yields accuracy gains, it consistently demands the largest rollout budget. For instance, on Qwen2.5-Math-1.5B, DS requires 3.6M rollouts to reach an average score of 34.10, whereas PREPO attains a higher 36.23 with only 1.1M rollouts. This confirms that DS wastes computation by discarding uninformative prompts only after rollouts are generated. Random selection occasionally performs comparably to or even better than DS, yet its rollout cost remains high. On Qwen2.5-Math-7B, for example, random selection consumes 905K rollouts to achieve 39.45, while PREPO surpasses it with 39.59 using just 540K rollouts. Overall, PREPO delivers both higher accuracy and lower cost, while also producing a more diverse set of problems than online random selection (see Appendix H.6).

GRESO improves efficiency but lags behind PREPO. GRESO reduces rollout demand by pre-filtering uninformative prompts, making it more

efficient than DS and Random. However, its accuracy often falls short of PREPO. For instance, on Qwen3-4B, GRESO achieves 69.89 accuracy with 472K rollouts, while PREPO reaches 75.99 with only 348K rollouts. This suggests that PREPO’s intrinsic exploration signals provide a more effective alternative to heuristic rollout filtering.

PREPO generalizes across model architectures. On Llama3.1-8B (see Table 3), PREPO once again achieves the strongest results, reaching an average of 36.55 with just 115K rollouts. In contrast, DS requires nearly five times as many rollouts, with substantially lower accuracy. This confirms that PREPO generalizes effectively across model families and scales, providing a consistent advantage in both performance and efficiency.

Table 3: Performance comparison (%) on Llama. For each model, the top row reports the base model’s performance. Best results are highlighted in bold or underlined. K = thousand.

Method	GSM8K	MATH	Avg ↑	# Rollouts ↓
<i>Llama3.1-8B</i>	9.53	6.05	7.79	–
+ DS	39.50	17.00	28.25	553K
+ Random	46.63	14.60	30.61	266K
+ GRESO	41.77	16.80	29.29	273K
+ PREPO (Ours)	<u>51.10</u>	<u>21.81</u>	36.55	115K

PREPO could be more effective when perplexity distributions are concentrated. We found that PREPO achieves large improvement when perplexity values are relatively compact across Qwen3-4B, and LLaMA3.1-8B. As shown in Table 4, the normalized standard deviation of those models remains below one, indicating less dispersed distributions that make perplexity-based filtering more reliable. In these cases, PREPO can better exploit the PPL-schedule to reduce rollouts.

Table 4: Normalized standard deviation of prompt perplexity (std/mean) across models. *Values above one indicate more dispersed distributions.*

Qwen2.5-7B	Qwen3-4B	LLaMA3.1-8B	Qwen2.5-Math-7B	Qwen2.5-Math-1.5B
0.73	0.65	0.75	1.23	1.02

Training Dynamics of PREPO versus the Random Baseline Figures 5 and 11 (see Appendix G) present a comparison of training dynamics between PREPO and random selection. PREPO shows a higher entropy loss, indicating stronger exploratory behavior throughout training. It also sustains a higher gradient norm while avoiding instability, suggesting more active yet controlled parameter updates. In terms of learning efficiency, PREPO reduces the proportion of rollouts with zero advantage, thereby providing more informative gradients for optimization. Furthermore, the average prompt length under PREPO decreases steadily, implying an adaptive shift from longer to shorter prompts over time. A similar trend is observed in the average response length, where PREPO generates longer outputs than the baseline across most steps on average, reflecting a longer thinking behavior.

5.3 Ablation Study

In this section, we conduct the ablation analysis to isolate the contribution of each component in PREPO. In addition to the linear PPL-schedule, we evaluate four variants: (1) linear PPL-schedule, (2) random filtering with relative-entropy weighting, (3) linear PPL-schedule combined with an entropy-regularization loss, and (4) a nonlinear PPL-schedule. The nonlinear variant follows the general form $l(\rho) = \lfloor \rho^\alpha (N - K) \rfloor$ with $\alpha = 0.5$, which produces a schedule that shifts more quickly toward higher-PPL prompts.

As shown in Tables 5 and 6, PREPO achieves the highest average performance across all model scales. Variants that remove or modify either component show consistent performance drops, indicat-

Table 5: Ablation study on Llama. Performance comparison (%) between PREPO and (1) linear PPL-schedule, (2) random filtering with relative-entropy, (3) linear PPL-schedule with entropy loss, (4) non-linear PPL-schedule with relative-entropy. *Best results are highlighted in bold or underlined.*

Model	Method	GSM8K	MATH	Avg \uparrow
Llama3.1-8B	PREPO	<u>51.10</u>	<u>21.81</u>	36.55
	w/o relative entropy	46.85	18.25	32.55
	w/o PPL-schedule	34.55	16.29	25.42
	w/ entropy loss	4.45	8.30	6.33
	w/ non-linear schedule	46.85	20.56	33.71

ing that both the perplexity-based schedule and the relative-entropy weighting contribute to PREPO’s effectiveness. The nonlinear schedule performs competitively but remains slightly below the linear schedule, suggesting that the linear form is already a stable and effective design choice.

5.4 Alternative Selection Ratio

Figure 6a and Table 8 (see Appendix G) show the effect of varying the online batch selection ratio (K/N) with N fixed for PREPO.

The 20% ratio achieves the best average performance while also requiring the fewest rollouts. A 30% ratio is competitive on some benchmarks but less efficient, while 25% does not improve overall accuracy. Smaller ratios (15%, 10%, 5%) degrade performance and consume more rollouts. Overall, 20% strikes the best balance between accuracy and efficiency, and is adopted as the default in our experiments.

5.5 Analysis of Relative-Entropy Weights

Effective Batch Sizes Relative-entropy weighting does not increase the effective batch size (i.e., $\frac{1}{B} \sum_i w_i$); it merely redistributes gradient contributions across sequences. Since the weights are normalized by \bar{H} , we obtain

$$\frac{1}{B} \sum_{i=1}^B w_i \cdot |o_i| = \frac{1}{B\bar{H}} \sum_{i=1}^B |o_i| \bar{H}_i = \frac{1}{B\bar{H}} \sum_{i=1}^B \sum_{t=1}^{|o_i|} H_{i,t} = \frac{\sum_{i=1}^B |o_i|}{B}. \quad (9)$$

The token-weighted average weight equals the average sequence length. If all sequences have equal length, then $\frac{1}{B} \sum_i w_i = 1$. As shown in Figure 6b, the effective batch size for PREPO on Qwen2.5-Math-1.5B stays close to B throughout training. At the early steps, LOW-PPL prompts yield low-entropy responses, so higher-entropy rollouts receive more weight, pushing the average above B . As training advances and higher-PPL

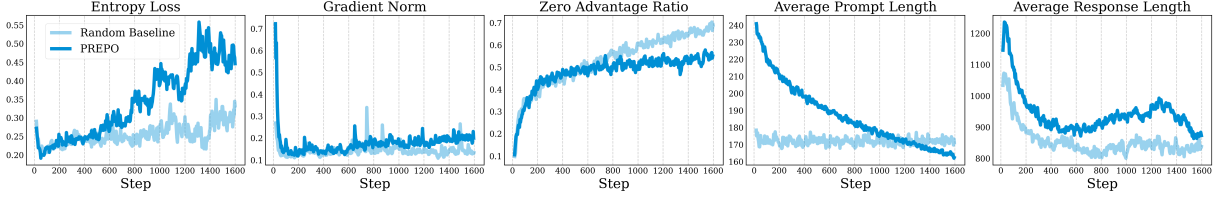


Figure 5: Full comparison between PREPO and random selection on Qwen2.5-Math-7B.

Table 6: Ablation study on Qwen. Performance comparison (%) between PREPO and (1) linear PPL-schedule, (2) random filtering with relative-entropy, (3) linear PPL-schedule with entropy loss, (4) non-linear PPL-schedule with relative-entropy. *Best results are highlighted in bold or underlined.*

Model	Method	AIME25	AIME24	MATH	Olympiad Bench	Avg \uparrow
Qwen2.5-Math-7B	PREPO	<u>12.81</u>	<u>26.15</u>	<u>77.80</u>	<u>41.58</u>	39.59
	w/o relative entropy	10.00	23.33	74.60	39.21	36.79
	w/o PPL-schedule	11.87	25.73	76.48	34.43	37.13
	w/ entropy loss	10.73	23.54	75.25	38.81	37.08
	w/ non-linear schedule	12.71	26.15	76.40	40.12	38.85
Qwen2.5-7B	PREPO	<u>10.20</u>	16.09	76.30	39.85	35.61
	w/o relative entropy	6.98	<u>16.41</u>	75.70	38.47	34.39
	w/o PPL-schedule	8.89	16.04	<u>79.41</u>	22.54	31.72
	w/ entropy loss	6.25	16.35	77.36	21.04	30.25
	w/ non-linear schedule	9.58	16.25	76.31	39.60	35.44
Qwen2.5-Math-1.5B	PREPO	<u>20.00</u>	<u>16.67</u>	<u>76.25</u>	<u>32.00</u>	36.23
	w/o relative entropy	10.21	15.68	72.10	30.50	32.12
	w/o PPL-schedule	11.81	13.12	70.21	30.02	31.29
	w/ entropy loss	6.46	16.56	73.77	30.99	31.95
	w/ non-linear schedule	9.15	15.83	75.85	30.26	32.77
Qwen3-4B	PREPO	<u>66.67</u>	<u>80.00</u>	<u>96.60</u>	<u>60.67</u>	75.99
	w/o relative entropy	64.77	72.70	90.54	59.06	71.77
	w/o PPL-schedule	64.99	77.73	95.03	57.86	73.90
	w/ entropy loss	61.10	75.17	88.39	60.54	71.28
	w/ non-linear schedule	61.45	74.47	92.37	<u>60.67</u>	72.24

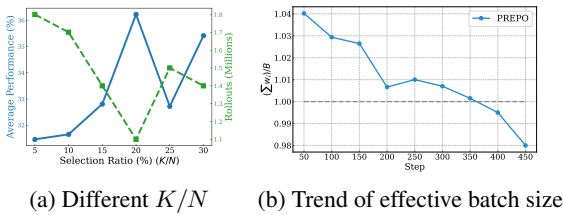


Figure 6: More ablation studies.

prompts are used, the overall entropy rises, and normalization shifts the average slightly below B .

6 Conclusion and Future Work

This study investigated how intrinsic data properties can improve the efficiency of RLVR training. Prompt perplexity enables a natural schedule from easier to harder prompts, while sequence-level relative entropy amplifies exploratory rollouts. Integrated in PREPO, these components reduce rollout

cost while maintaining or improving benchmark performance.

Beyond empirical gains, PREPO shows that RLVR can be guided by interpretable, policy-intrinsic cues rather than human heuristic or auxiliary models. This makes PREPO a practical and transparent recipe for compute-efficient RLVR, accessible even to researchers with limited resources. Future work may explore additional intrinsic signals (e.g., input token size) and combine data-driven exploration with system-level optimizations.

Appendix roadmap. Appendices B–C analyze prompt properties and training dynamics. Appendices D–F provide baseline definitions, implementation details, and evaluation procedures for reproducibility. Appendix G reports additional results, including ablations and sensitivity analyses. Appendix H discusses common concerns such as problem diversity and memorization risk.

7 Limitations

This study has several limitations that should be acknowledged. (1) response lengths were restricted to 32K tokens, leaving the applicability of PREPO to models generating substantially longer outputs an open question; and (2) the evaluation was limited to mathematical reasoning tasks, while its effectiveness in other domains remains to be explored.

References

Alon Albalak, Duy Phung, Nathan Lile, Rafael Rafailov, Kanishk Gandhi, Louis Castricato, Anikait Singh, Chase Blagden, Violet Xiang, Dakota Mahan, and Nick Haber. 2025. **Big-math: A large-scale, high-quality math dataset for reinforcement learning in language models.** *Preprint*, arXiv:2502.17387.

Art of Problem Solving. 2024. 2024 aime i. https://artofproblemsolving.com/wiki/index.php/2024_AIME_I. Accessed: 2025-06-09.

Art of Problem Solving. 2025. 2025 aime i. https://artofproblemsolving.com/wiki/index.php/2025_AIME_I. Accessed: 2025-06-09.

Sanghwan Bae, Jiwoo Hong, Min Young Lee, Hanbyul Kim, JeongYeon Nam, and Donghyun Kwak. 2025. Online difficulty filtering for reasoning oriented reinforcement learning. *arXiv preprint arXiv:2504.03380*.

Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu, Chao Wang, Cheng Zhu, Chengjun Xiao, Chengyu Du, Chi Zhang, Chu Qiao, Chunhao Zhang, Chunhui Du, Congchao Guo, Da Chen, Deming Ding, and 80 others. 2025a. **Minimax-m1: Scaling test-time compute efficiently with lightning attention.** *ArXiv*, abs/2506.13585.

Xiaoyin Chen, Jiarui Lu, Minsu Kim, Dinghuai Zhang, Jian Tang, Alexandre Piché, Nicolas Gontier, and Yoshua Bengio. 2025b. Self-evolving curriculum for llm reasoning. *arXiv preprint arXiv:2505.14970*.

Xinjie Chen, Minpeng Liao, Guoxin Chen, Chengxi Li, and Biao Fu. 2025c. From data-centric to sample-centric: Enhancing llm reasoning via progressive optimization. *arXiv preprint arXiv:2507.06573*.

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. 2025. Reasoning with exploration: An entropy perspective. *arXiv preprint arXiv:2506.14758*.

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. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, and 1 others. 2025. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407. 570

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*. 571

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, and 1 others. 2024. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*. 572

Ehsan Kamalloo, Shipeng Li, Shikun Li, Zhiqin Yang, Xinghua Zhang, Gaode Chen, Xiaobo Xia, Hengyu Liu, and Zhe Peng. 2025. Learnalign: Reasoning data selection for reinforcement learning in llms based on improved gradient alignment. *arXiv preprint arXiv:2506.11480*. 573

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Haoteng Zhang, and Ion Stoica. 2023. **Efficient memory management for large language model serving with pagedattention.** *Proceedings of the 29th Symposium on Operating Systems Principles*. 574

Kai Li, Xinggao Fan, and X. Liu. 2025. Not all rollouts are useful: Down-sampling rollouts in llm reinforcement learning. *arXiv preprint arXiv:2504.13818*. 575

Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023a. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*. 576

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023b. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*. 577

X Liu and 1 others. 2025. Improving data efficiency for llm reinforcement fine-tuning through difficulty-targeted online data selection and rollout replay. *arXiv preprint arXiv:2506.05316*. 578

Shubham Parashar, Shurui Gui, Xiner Li, Hongyi Ling, Sushil Vemuri, Blake Olson, Eric Li, Yu Zhang, 579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616	James Caverlee, Dileep Kalathil, and 1 others. 2025.	Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	670
617	Curriculum reinforcement learning from easy to	Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,	671
618	hard tasks improves llm reasoning. <i>arXiv preprint</i>	Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo:	672
619	<i>arXiv:2506.06632</i> .	An open-source llm reinforcement learning system	673
		at scale. <i>arXiv preprint arXiv:2503.14476</i> .	674
620	Yun Qu, Qi Wang, Yixiu Mao, Vincent Tao Hu, Björn	Enci Zhang, Xingang Yan, Wei Lin, Tianxiang Zhang,	675
621	Ommer, and Xiangyang Ji. 2025. Can prompt	and Qianchun Lu. 2025. Learning like humans: Ad-	676
622	difficulty be online predicted for accelerating rl	vancing llm reasoning capabilities via adaptive dif-	677
623	finetuning of reasoning models? <i>arXiv preprint</i>	ficulty curriculum learning and expert-guided self-	678
624	<i>arXiv:2507.04632</i> .	reformulation. <i>arXiv preprint arXiv:2505.08364</i> .	679
625	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Haizhong Zheng, Yang Zhou, Brian R Bartoldson,	680
626	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Bhavya Kailkhura, Fan Lai, Jiawei Zhao, and Beidi	681
627	Zhang, YK Li, Yang Wu, and 1 others. 2024.	Chen. 2025. Act only when it pays: Efficient rein-	682
628	Deepseekmath: Pushing the limits of mathematical	forcement learning for llm reasoning via selective	683
629	reasoning in open language models. <i>arXiv preprint</i>	rollouts. <i>arXiv preprint arXiv:2506.02177</i> .	684
630	<i>arXiv:2402.03300</i> .		
631	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	Yinmin Zhong, Zili Zhang, Bingyang Wu, Shengyu	685
632	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	Liu, Yukun Chen, Changyi Wan, Hanpeng Hu, Lei	686
633	Lin, and Chuan Wu. 2025. Hybridflow: A flexible	Xia, Ranchen Ming, Yibo Zhu, and 1 others. 2024.	687
634	and efficient rlhf framework. In <i>Proceedings of the</i>	Optimizing rlhf training for large language models	688
635	<i>Twentieth European Conference on Computer Sys-</i>	with stage fusion. <i>arXiv preprint arXiv:2409.13221</i> .	689
636	<i>tems</i> , pages 1279–1297.		
637	Kimi Team, Angang Du, Bofei Gao, Bawei Xing,		
638	Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun		
639	Xiao, Chenzhuang Du, Chonghua Liao, and 1 others.		
640	2025. Kimi k1. 5: Scaling reinforcement learning		
641	with llms. <i>arXiv preprint arXiv:2501.12599</i> .		
642	Qwen Team. 2024. Qwen2 technical report. <i>arXiv</i>		
643	<i>preprint arXiv:2407.10671</i> .		
644	Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shix-		
645	uan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin		
646	Yang, Zhenru Zhang, and 1 others. 2025. Beyond		
647	the 80/20 rule: High-entropy minority tokens drive		
648	effective reinforcement learning for llm reasoning.		
649	<i>arXiv preprint arXiv:2506.01939</i> .		
650	Zhenting Wang and Cui Guofeng. 2025. Dump: Auto-		
651	mated distribution-level data selection for rl. <i>arXiv</i>		
652	<i>preprint arXiv:2504.09710</i> .		
653	Mingqi Wu, Zhihao Zhang, Qiaole Dong, Zhiheng Xi,		
654	Jun Zhao, Senjie Jin, Xiaoran Fan, Yuhao Zhou, Yan-		
655	wei Fu, Qin Liu, and 1 others. 2025. Reasoning or		
656	memorization? unreliable results of reinforcement		
657	learning due to data contamination. <i>arXiv preprint</i>		
658	<i>arXiv:2507.10532</i> .		
659	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,		
660	Binyuan Hui, Bo Zheng, Bowen Yu, Chang		
661	Gao, Chengen Huang, Chenxu Lv, and 1 others.		
662	2025. Qwen3 technical report. <i>arXiv preprint</i>		
663	<i>arXiv:2505.09388</i> .		
664	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,		
665	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong		
666	Tu, Jingren Zhou, Junyang Lin, and 1 others. 2024.		
667	Qwen2. 5-math technical report: Toward mathe-		
668	matical expert model via self-improvement. <i>arXiv</i>		
669	<i>preprint arXiv:2409.12122</i> .		

Appendices

A Disclaimer on LLM Usage	11
B HIGH-PPL and LOW-PPL Prompts	11
C More Training Dynamics of LOW-PPL and HIGH-PPL Prompts	13
D Description of Baseline Methods	13
E Code Implementation	13
F Details of Evaluation Framework	13
G Additional Experimental Results	14
H Discussion	14
H.1 Do the Results Generalize to Other Training Datasets?	14
H.2 What Does the PPL-schedule Contribute During Training?	14
H.3 What Does Relative Entropy Bring to Training?	16
H.4 Does Prompt Perplexity Change During Training?	16
H.5 What Is the Time Cost of PREPO?	16
H.6 How Does PREPO Affect Problem Diversity?	16
H.7 How Does PREPO Compare to Training Without Filtering?	16
H.8 Does PREPO Lead to Memorization of Training Data?	16
H.9 How Sensitive Is PREPO to Extreme Entropy Values?	16

A Disclaimer on LLM Usage 717

The use of LLMs is permitted as a general-purpose assistance tool. In this work, LLMs were employed solely for grammar correction and sentence rephrasing. Their role was restricted to improving clarity and style; therefore, they are not considered contributors to the research. 718
719
720
721
722
723

B HIGH-PPL and LOW-PPL Prompts 724

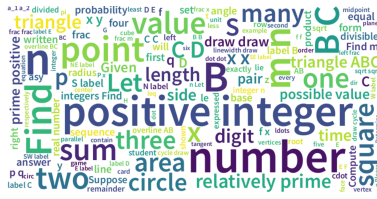
As illustrated in Figure 7, HIGH-PPL prompts exhibit a greater prevalence of non-English characters relative to LOW-PPL prompts. This pattern is consistently observed across both the Qwen2.5-series and Llama model families. 725
726
727
728
729

Example of LOW-PPL Problems

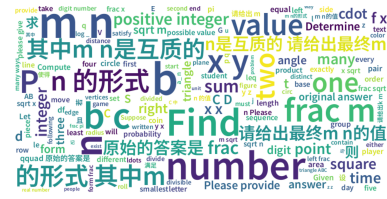
- Cube $ABCDEFGH$, labeled as shown below, has edge length 1 and is cut by a plane passing through vertex D and the midpoints M and N of \overline{AB} and \overline{CG} respectively. The plane divides the cube into two solids. Find the volume of the larger of the two solids. [asy] import cse5; unitsize(8mm); pathpen=black; pair A = (0,0), B = (3.8,0), C = (5.876,1.564), D = (2.076,1.564), E = (0,3.8), F = (3.8,3.8), G = (5.876,5.364), H = (2.076,5.364), M = (1.9,0), N = (5.876,3.465); pair[] dotted = A,B,C,D,E,F,G,H,M,N; D(A-B-C-G-H-E-A); D(E-F-B); D(F-G); pathpen=dashed; D(A-D-H); D(D-C); dot(dotted); label("A",A,SW); label("B",B,S); label("C",C,SE); label("D",D,NW); label("E",E,W); label("F",F,SE); label("G",G,NE); label("H",H,NW); label("M",M,S); label("N",N,NE); [/asy]The answer is in the form $x0cracmn$, where $\gcd(m, n) = 1$. Please provide the value of $m + n$.
- The number $a = \frac{p}{q}$, where p and q are relatively prime positive integers, has the property that the sum of all real numbers x satisfying

$$[x] \cdot \{x\} = a \cdot x^2$$

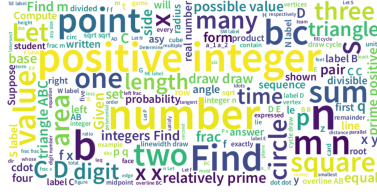
is 420, where $[x]$ denotes the greatest integer less than or equal to x and $\{x\} =$



(a) Low-PPL Prompts (Qwen2.5-7B)



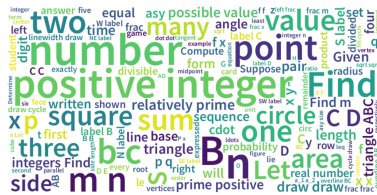
(b) HIGH-PPL Prompts (Qwen2.5-7B)



(c) Low-PPL Prompts (Qwen2.5-Math-1.5B)



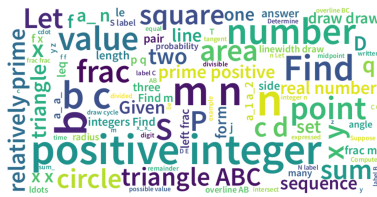
(d) HIGH-PPL Prompts (Qwen2.5-Math-1.5B)



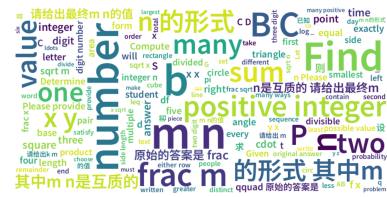
(e) Low-PPL Prompts (Qwen2.5-Math-7B)



(f) HIGH-PPL Prompts (Qwen2.5-Math-7B)



(g) Low-PPL Prompts (Llama3.1-8B)



(h) HIGH-PPL Prompts (Llama3.1-8B)

Figure 7: Wordcloud of the Most Frequent Words in Low-/HIGH-PPL Prompts

$x - \lfloor x \rfloor$ denotes the fractional part of x .
What is $p + q$?

- Let a, b be positive integers satisfying

$$\sqrt{\frac{ab}{2b^2 - a}} = \frac{a + 2b}{4b}$$

. Find $|10(a - 5)(b - 15)| + 8$.

- What is the greatest common divisor of $121^2 + 233^2 + 345^2$ and $120^2 + 232^2 + 346^2$?

Example of HIGH-PPL Problems

- At the Lexington High School, each student is given a unique five-character ID consisting of uppercase letters. Compute the number of possible IDs that contain the string "LMT".
- 设3阶实对称矩阵A的三个特征值分别为 $-1, -1, 2$, 且 $(1, 1, -1)^T$ 是特征值2所对应的特征向量. 记A中所有元素的平方和为 I , 则 $[I]=$ _____.
(English: Let the eigenvalues of A real symmetric matrix of order 3 be $-1, -1, 2$, respectively, and $(1, 1, -1)^T$ be the eigenvector corresponding to eigenvalue 2. For the sum of squares of all elements in the

A I , is $[I] = \underline{\hspace{2cm}}$.)

- 求六个元素的置换群 S_6 中6 阶元素的个数。(English: Find the number of elements of order 6 in the permutation group S_6 of six elements.)
- Three bells begin to ring simultaneously. The intervals between strikes for these bells are, respectively, $\frac{4}{3}$ seconds, $\frac{5}{3}$ seconds, and 2 seconds. Impacts that coincide in time are perceived as one. How many beats will be heard in 1 minute? (Include first and last.)

C More Training Dynamics of LOW-PPL and HIGH-PPL Prompts

For Qwen2.5-7B, Qwen2.5-Math-1.5B, and Llama3.1-8B, we found similar trends across training dynamics of LOW-PPL and HIGH-PPL groups, as shown in Figure 8, 9, and 10.

D Description of Baseline Methods

We compare PREPO against the following three baseline strategies.

- **Dynamic Sampling (DS).** Dynamic Sampling, as introduced in DAPO [Yu et al., 2025], dynamically filters out prompt groups whose generated responses all produce identical rewards (i.e., zero variance). In each training batch, DS resamples such uninformative prompt groups, thereby ensuring that the batch maintains a sufficient proportion of prompts that give meaningful gradient signals. However, it can still incur high rollout costs because many sampled prompts may remain uninformative until they are filtered.
- **Random.** The Random baseline uniformly selects prompts and associated rollouts without regard to historical feedback or variance in reward. All prompts are treated equally, so there is no mechanism to avoid rollouts on uninformative or zero-variance prompts. This method serves as a lower bound in terms of data selection sophistication.
- **GRESO.** GRESO (GRPO with Efficient Selective Rollout) [Zheng et al., 2025] is a lightweight, online, pre-rollout filtering approach. It uses statistics of reward dynamics over previous

epochs to predict which prompts are likely to be uninformative (e.g., zero variance among responses) and skips them before performing rollouts.

E Code Implementation

At each global step, PPL-schedule computes perplexity on the current batch of prompts and selects the subset that matches the current progress ratio. This removes the need for tuning and avoids extra computational cost. Here is the implementation:

```
def fit(self):
    .....
    # compute perplexity for each prompt
    batch = self.actor_rollout_wg.get_ppl_batch()
    h(batch)
    # compute training progress
    rho = self.global_steps /
    self.total_training_steps
    # determine the prompts to select
    start_idx = math.floor((N - K) * rho)
    idx = torch.argsort(batch['ppl'])[start_idx
    : start_idx + K]
    # select prompts based on the indices
    batch = batch.select_via_index(idx)
    .....
```

Similarly, the relative-entropy component uses only a few lines of code, requires no manual tuning, and adds no overhead beyond the entropy computation that is already part of the existing pipeline. Here is the implementation:

```
def update_policy(self, data: DataProto):
    .....
    # all return: (bsz, response_length)
    entropy, log_prob = self._forward_micro_batch(
    data, temperature=temperature)
    sample_entropy = verl_F.masked_mean(entropy,
    response_mask, axis=1).detach()
    batch_entropy = verl_F.masked_mean(entropy,
    response_mask).detach()
    weight = sample_entropy/batch_entropy
    advantages *= weight[:, None]
    .....
```

F Details of Evaluation Framework

We use the `_validate()`³ function to evaluate the model. In addition, the basic verifier is `math-verify`, and we add a LaTeX normalization package (`latex2sympy2-extended`) for the AIME and Olympiad bench.

³We adopt from the public Github repository: [line 569 in ray_trainer.py](#)

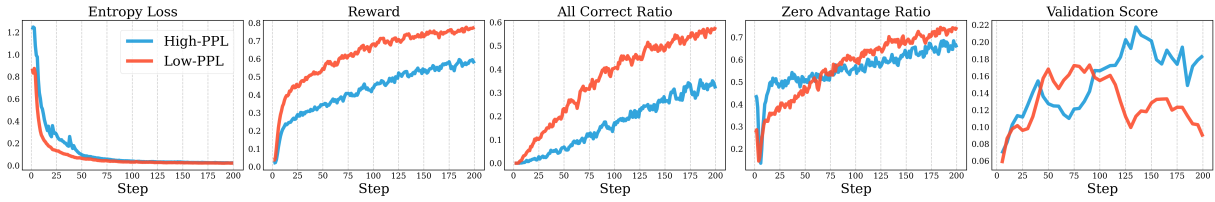


Figure 8: Training Dynamics of LOW-PPL vs. HIGH-PPL Prompts on Qwen2.5-7B.

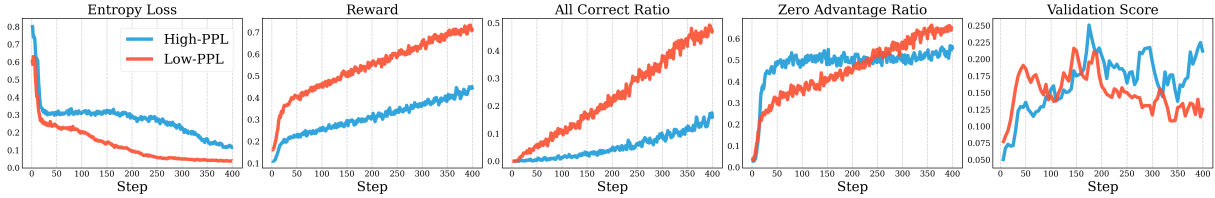


Figure 9: Training Dynamics of LOW-PPL vs. HIGH-PPL Prompts on Qwen2.5-Math-1.5B.

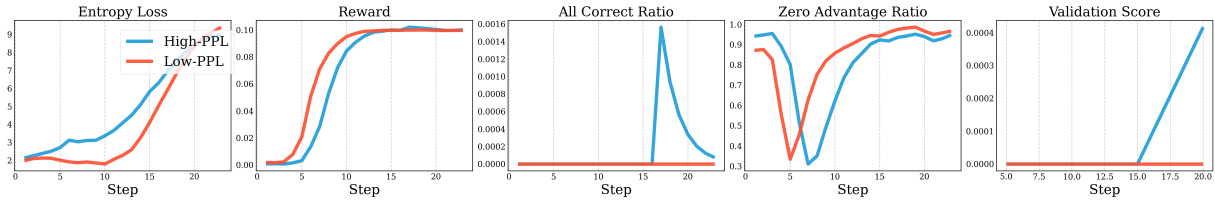


Figure 10: Training Dynamics of LOW-PPL vs. HIGH-PPL Prompts on Llama3.1-8B.

G Additional Experimental Results

In this section, we show the performance comparison between the random baseline and PREPO in Figure 11, the performance comparison across different evaluation metrics in Table 7 and the sensitivity analysis of the selection ratio in Table 8.

H Discussion

H.1 Do the Results Generalize to Other Training Datasets?

As DAPO-Math-17K contains both English and Chinese questions, prompt perplexity could correlate with language. To exclude the cofounder of language effects, we trained Qwen2.5-Math-7B on a purely English dataset selected from OpenR1-Math-220k⁴, and compared PREPO with several baselines under the same training budget. The results are shown in Table 9 and the Figure 12 below. PREPO again achieves lower rollout usage while achieving a higher average performance. This indicates that the gains of PREPO are not driven by language imbalance.

H.2 What Does the PPL-schedule Contribute During Training?

We compared three configurations on Qwen2.5-Math-7B training exclusively with HIGH-PPL prompts, exclusively with LOW-PPL prompts, and the PPL-schedule, which gradually transitions from Low- to HIGH-PPL prompts. The training dynamics are shown in Figure 13.

⁴<https://huggingface.co/datasets/open-r1/OpenR1-Math-220k>

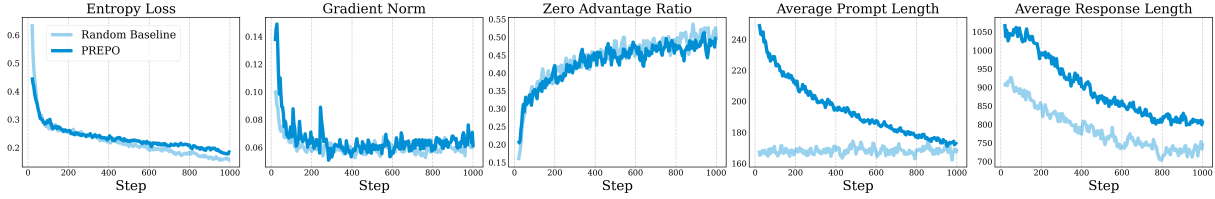


Figure 11: Full comparison between PREPO and random selection on Qwen2.5-Math-1.5B.

Table 7: Details of performance comparison (%) on Qwen models. Avg (avg@16/32/64) means the average of AIME25 (avg@16/32/64), AIME24 (avg@16/32/64), MATH, and Olympiad. Best results are highlighted in bold or underlined.

Method	AIME25 (avg@16)	AIME25 (avg@32)	AIME25 (avg@64)	AIME24 (avg@16)	AIME24 (avg@32)	AIME24 (avg@64)	Avg ↑ (avg@16)	Avg ↑ (avg@32)	Avg ↑ (avg@64)
<i>Qwen2.5-7B</i>	1.25	2.60	1.17	4.17	4.90	4.95	27.69	28.21	27.87
+ DS	7.92	7.08	7.19	<u>17.55</u>	15.00	14.90	34.97	34.76	34.13
+ Random	6.98	6.35	6.51	16.41	15.63	15.89	34.39	34.04	34.14
+ GRESO	9.22	7.29	7.50	10.83	13.65	13.39	34.39	34.92	34.90
+ PREPO (Ours)	<u>10.21</u>	<u>10.62</u>	<u>10.00</u>	16.09	<u>15.52</u>	<u>16.35</u>	35.61	35.72	35.63
<i>Qwen2.5-Math-1.5B</i>	3.54	3.75	4.48	10.21	8.33	8.96	24.23	23.81	24.15
+ DS	10.83	13.76	14.54	<u>25.83</u>	<u>21.72</u>	<u>23.04</u>	34.10	33.80	34.33
+ Random	<u>20.00</u>	17.45	18.81	16.67	15.42	16.50	35.86	34.91	35.52
+ GRESO	15.38	11.40	15.39	20.00	19.20	17.97	34.16	32.86	33.56
+ PREPO (Ours)	<u>20.00</u>	19.69	19.38	16.67	16.98	17.03	36.23	36.23	36.17
<i>Qwen2.5-Math-7B</i>	8.95	9.17	8.80	17.50	20.80	19.22	35.45	35.45	34.96
+ DS	13.33	14.27	14.84	<u>33.33</u>	<u>35.63</u>	<u>34.69</u>	34.10	40.36	40.26
+ Random	10.00	12.50	12.50	26.67	29.34	29.87	35.86	40.73	40.86
+ GRESO	<u>18.33</u>	<u>18.58</u>	<u>18.58</u>	25.83	24.38	25.00	37.46	36.90	37.05
+ PREPO (Ours)	12.81	14.41	15.93	26.15	29.17	29.25	39.59	40.75	41.15
<i>Qwen3-4B</i>	30.00	47.30	44.89	53.33	61.60	28.70	57.53	63.92	55.09
+ DS	63.33	53.90	54.80	66.67	63.50	65.30	70.78	67.63	68.30
+ Random	60.00	58.85	60.16	70.00	65.73	64.69	71.33	69.98	70.05
+ GRESO	56.67	58.33	57.24	69.17	71.77	69.79	69.89	70.96	70.19
+ PREPO (Ours)	<u>66.67</u>	<u>62.19</u>	<u>63.39</u>	<u>80.00</u>	<u>74.28</u>	<u>77.45</u>	75.99	73.44	74.53

Table 8: Performance Comparison (%) of PREPO with Different Selection Ratio (K/B , B fixed) on Qwen2.5-Math-1.5B. Best results are highlighted in bold or underlined.

Selection Ratio	AIME25	AIME24	MATH	Olympiad	Avg ↑	# Rollouts ↓
30%	<u>20.00</u>	<u>20.83</u>	75.85	25.00	35.42	1.4M
25%	13.33	20.00	73.55	24.00	32.72	1.5M
20%	<u>20.00</u>	16.67	<u>76.25</u>	<u>32.00</u>	36.23	1.1M
15%	13.33	20.00	<u>75.75</u>	22.17	32.81	1.4M
10%	15.83	19.17	70.25	21.33	31.65	1.7M
5%	19.17	16.67	69.85	20.17	31.46	1.8M

Table 9: Comparison of PREPO and baselines trained on all-English data (Model: Qwen2.5-Math-7B)

Method	AIME25	AIME24	MATH	Olympiad	Avg ↑	#Rollout ↓
DS	11.98	30.00	72.80	42.02	39.20	3256K
Random	10.79	28.25	<u>80.21</u>	36.72	38.99	2457K
GRESO	6.67	30.00	79.20	<u>45.63</u>	40.38	1331K
PREPO	<u>12.92</u>	<u>33.33</u>	78.05	44.10	42.10	860K

In terms of entropy loss, the PPL-schedule achieves a balance between the two extremes:

entropy decreases steadily but not excessively, thereby mitigating the risk of collapse. With re-

818
819

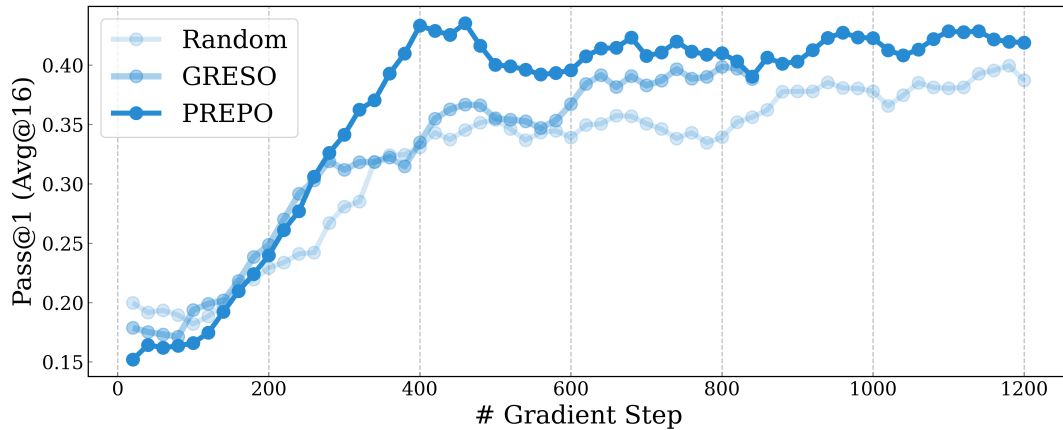


Figure 12: Comparison of performances across PREPO, GRESO, and Random (x-axis: training step, y-axis: validation score, model: Qwen2.5-Math-7B, training data: OpenR1-Math)

spect to the zero-advantage ratio, the PPL-schedule consistently maintains a lower value, ensuring that a greater proportion of rollouts remain informative throughout training.

H.3 What Does Relative Entropy Bring to Training?

As shown in Figure 14, PREPO with relative entropy further reduces the zero-advantage ratio, thereby improving sample efficiency.

Case Analysis. In Figure 15, we display examples of responses with their token-level entropy from the same mini-batch, where darker colors mean higher entropy, and each sequence is weighted by its relative entropy.

H.4 Does Prompt Perplexity Change During Training?

As shown in Figure 16, we computed the PPL range at each training epoch and observed that it remained relatively stable throughout training. In addition, the PPL of the prompts used for training exhibited minimal variation.

H.5 What Is the Time Cost of PREPO?

As shown in Figure 17, the time consumption of calculating prompt PPL is barely minimal compared to the rollout generation duration.

H.6 How Does PREPO Affect Problem Diversity?

The analysis indicates that PREPO selects a more diverse set of problems than random sampling. Specifically, PREPO achieves broader coverage

across Mathematics Subject Classification⁵ (MSC) categories during training, as illustrated in Figure 18.

H.7 How Does PREPO Compare to Training Without Filtering?

For Qwen2.5-Math-7B, we observe that PREPO attains performance comparable to training on the full dataset without any filtering, i.e., using 5 times rollouts per step, as shown in Figure 19. The “Less is More” pattern implies that efficient selection can reduce the amount of data needed for RLVR training.

H.8 Does PREPO Lead to Memorization of Training Data?

Following Wu et al. [2025], we trimmed 40% of each prompt to construct *partial problems*. We then evaluate models on these partial prompts and compute the average pass rate of 16 generations. As shown in Figure 20, the vast majority of partial problems have near-zero pass rate, with only a small fraction achieving non-trivial success. These distributions suggest that the model does not merely memorize the training data but instead requires the full problem context to solve tasks.

H.9 How Sensitive Is PREPO to Extreme Entropy Values?

Since $w_i = \bar{H}_i/\bar{H}$ is normalized by the batch mean, a rollout with $\bar{H}_j \gg \bar{H}$ has a large weight $w_j \gg 1$, while simultaneously reducing the weights of the others ($w_i \ll 1$ for $i \neq j$). The

⁵MSC : <https://zbmath.org/classification/>

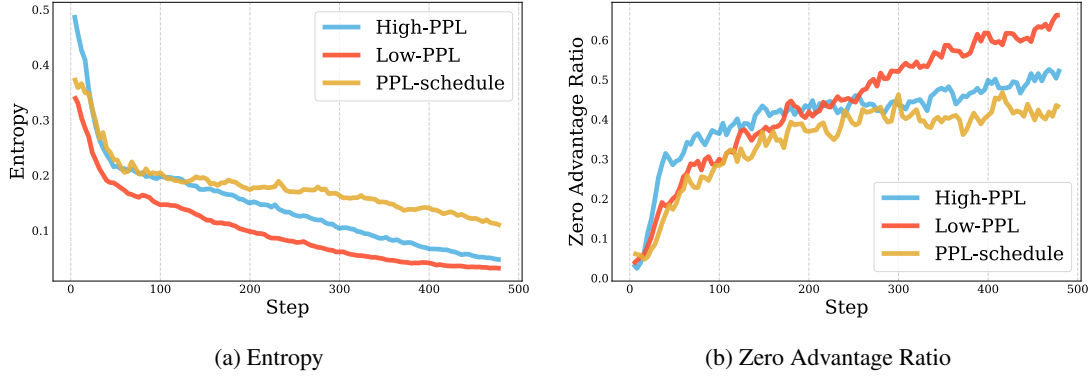


Figure 13: Comparison of Training Dynamics between PPL-schedule and Static PPL Selection (Low- and HIGH-PPL groups)

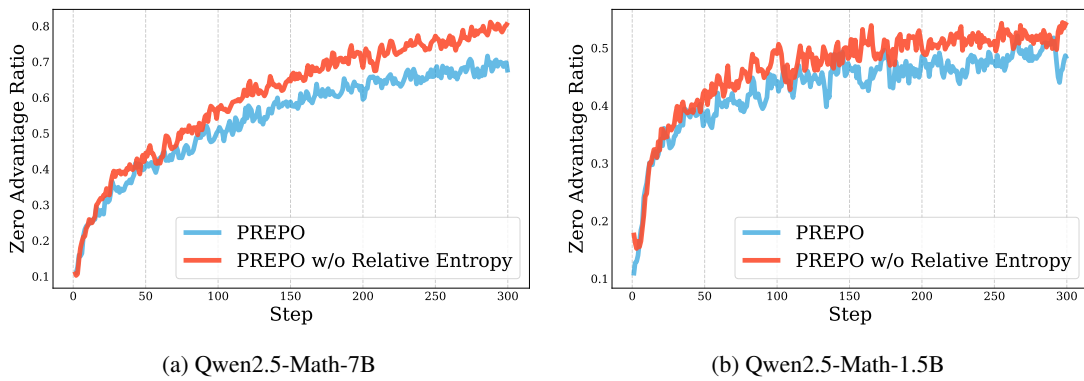


Figure 14: Comparison of zero advantage ratio between PPL-schedule and PREPO.

880 derivative

$$881 \quad \frac{\partial w_i}{\partial \bar{H}_j} = \begin{cases} \frac{1}{\bar{H}} - \frac{\bar{H}_j}{\bar{H}^2} \frac{|o_j|}{\sum_k |o_k|}, & i = j, \\ -\frac{\bar{H}_i}{\bar{H}^2} \frac{|o_j|}{\sum_k |o_k|}, & i \neq j, \end{cases}$$

882 clarifies this behavior: (1) for $i = j$, the two terms
883 partially cancel, so although w_j is already large, its
884 growth rate is moderated as \bar{H}_j increases further;
885 (2) for $i \neq j$, the derivative is negative, confirming
886 that a large \bar{H}_j suppresses the weights of all other
887 rollouts. Empirically, such outliers are rare. As
888 shown in Figure 21, the distribution of relative en-
889 tropy weights remains stable mainly across training
890 steps 100, 200, 300, and 400 for Qwen2.5-Math-
891 1.5B, where most rollouts clustered near one, with
892 negligible mass on extreme values.

To solve the problem, we need to understand the constraint and then systematically count the valid combination of 3 cards. We will break down the problem into manageable parts and use Python code to ensure the accuracy of our solution. **Understanding the problem** We have a set of 27 distinct cards, each with a shape, a letter, and a number. We need to choose 3 cards such that no two cards share two matching values. This means that for any two cards, they can have at most one matching value. **Counting the valid sets** We will iterate over all possible combinations of 3 cards and check if they satisfy the condition. If they do, we will count them as valid. We will use Python code to compute this. Here is the code to solve it: `python itertools.combinations(cards, 3) # Generate all combinations of 3 cards
for comb in combinations(cards, 3): # Iterate over all combinations
 valid = True # Assume the combination is valid
 for i in range(3): # Iterate over the 3 cards in the combination
 for j in range(i+1, 3): # Iterate over the remaining 2 cards
 # Check if the two cards share two matching values
 if comb[i][0] == comb[j][0] and comb[i][1] == comb[j][1]:
 valid = False # Invalid combination
 if valid: # Valid combination
 count += 1 # Increment the count of valid sets
print(count) # Output the number of valid sets` The number of valid sets is 1278. The number of ways to choose an unordered set of 3 cards from the pack such that no two of the chosen cards have two matching values is 1278.

(a) Correct answer
x1.65

Let's solve this problem step-by-step. We need to find the maximum value of $S_d \cdot n_S$, where $S_d \cdot n_S$ is the greatest common divisor (GCD) of $S \cdot n \cdot S$ and $S \cdot a \cdot (n+1)S$, and $S \cdot a \cdot n = 100 + n^2 \cdot 2$. First, let's express $S \cdot a \cdot n_S$ and $S \cdot a \cdot (n+1)S$: $\{a, n\} = 100 + n^2 \cdot 2$, $\{a, n+1\} = 100 + (n+1)^2 \cdot 2 = 100 + n^2 \cdot 2 + 2n + 1$. Now, we want to find the GCD of $S \cdot a \cdot n_S$ and $S \cdot a \cdot (n+1)S$: $\{d, n\} = \gcd(a \cdot n, a \cdot (n+1)) = \gcd(100 + n^2 \cdot 2, 100 + n^2 \cdot 2 + 2n + 1)$. Notice that $\{a, n+1\} - a \cdot n = (100 + n^2 \cdot 2 + 2n + 1) - (100 + n^2 \cdot 2) = 2n + 1$. Therefore, we can write $\{d, n\} = \gcd(100 + n^2 \cdot 2, 2n + 1)$. Since $S \cdot d \cdot n_S$ divides both $S \cdot 100 + n^2 \cdot 2S$ and $S \cdot 2n + 1S$, it must also divide any linear combination of them. In particular, it must divide $\{d, n\} = \gcd(100 + n^2 \cdot 2, 2n + 1) = \gcd(2n + 1, 100 - n)$. Now, let's consider $S \cdot 100 - n \pmod{d} = \{d, n\} = \gcd(100 - n, k(2n + 1))$ for some integer k . Rearranging gives $\{100 - n\} = k(2k + 1) + k$. To maximize $S \cdot d \cdot n_S$, we need to check the values of $S \cdot n_S \cdot 2n + 1S$ divides 101 . Since 101 is a prime number, the only divisors of 101 are 1 and 101 . Therefore, we need to check both cases. If $S \cdot 2n + 1 = 1S$, then $S \cdot n = 0S$, but $S \cdot n_S$ must be a positive integer. If $S \cdot 2n + 1 = 101S$, then $S \cdot n = 50S$. Let's verify that $S \cdot d \cdot \{50\} = 101S$: $\{a \cdot \{50\}\} = 100 + 50^2 \cdot 2 = 100 + 2500 = 2600$. $\{a \cdot \{51\}\} = 100 + 51^2 \cdot 2 = 100 + 2601 = 2701$. $\gcd(2600, 2701) = 101$. Therefore, the maximum value of $S \cdot d \cdot n_S$ is 101 . Let's confirm this using Python code: `python import math def gcd(a, b): return math.gcd(a, b) max_gcd = 0 for n in range(1, 1000): # We can use a large range to check for maximum gcd a_n = 100 + n**2 a_n1 = 100 + (n+1)**2 d_n = gcd(a_n, a_n1) if d_n > max_gcd: max_gcd = d_n print(max_gcd) # Output 101` The maximum value of $S \cdot d \cdot n_S$ as $S \cdot n_S$ ranges through the positive integers is indeed $401S$, occurring when $S \cdot n = 200S$. The value 101 was also considered, but it is not the maximum. Therefore, the correct answer is $\boxed{401}$.

(c) Correct answer
x0.73

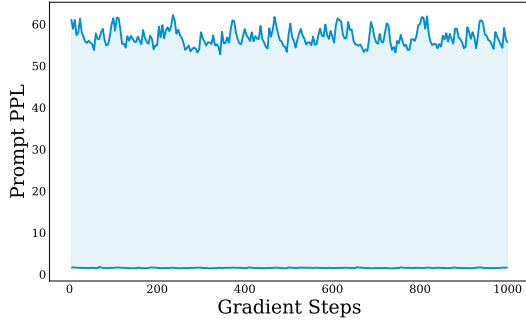
Let's solve this step-by-step. **Step 1:** Calculate the number of ways Winnie-the-Pooh can distribute the balloons. Winnie-the-Pooh wants to give a balloon of one of the seven colors of the rainbow to each of the 25 donkeys such that successive donkeys receive balloons of different colors and so that at least one balloon of each color is given to some donkey. This is a combinatorial problem where we need to ensure that at least one balloon of each color is given out. The first donkey can be given any of the 7 colors. The second donkey can be given any of the remaining 6 colors and so on. However, since we need to ensure that at least one balloon of each color is given out, we need to find a valid permutation of the 7 colors that can cover at least one full set of colors in the 25 positions. To calculate this, we can use the number of derangements (permutations where no element appears in its original position) and then make sure that every color is covered at least once. However, for large permutations where we have more positions than colors, we can use the fact that we can pick any permutation of the colors and then repeat them in a way that ensures all colors are covered at least once. The number of ways to do this is simply $7! \cdot \left(\frac{25}{7}\right)$ (the number of ways to arrange 7 colors) times the number of ways to fill the remaining positions (which is just placing the first 7 colors repeatedly in the new positions). Since we need to ensure that the first 7 donkeys have different colors, the number of ways to arrange the first 7 donkeys with different colors is $7!$. After the first 7 donkeys, the remaining 18 donkeys can have any of the 7 colors, but we need to ensure that the arrangement of the first 7 colors is maintained. This means that the remaining positions of the 7 donkeys can be filled in $(7 \cdot 18)!$ ways. So, the total number of ways for Winnie-the-Pooh is $7! \cdot \left(\frac{25}{7}\right) \cdot 7! \cdot 18!$. **Step 2:** Calculate the number of ways Eeyore can distribute the pots. Eeyore wants to give a pot of one of six colors (except red) to each of the 24 donkeys such that at least one pot of each color is given to some donkey. This is similar to the above problem but with fewer colors and fewer donkeys. The first donkey can be given any of the 6 colors. The second donkey can be given any of the remaining 5 colors and so on. Again, we need to ensure that at least one pot of each color is given out. The number of ways to arrange 6 pots in 24 positions (repetition allowed, and ensuring all 6 colors are used at least once) is similar to the above problem but with fewer colors. The number of ways to arrange 6 colors in 24 positions is $6! \cdot \left(\frac{24}{6}\right)$. The remaining 18 donkeys can be given any of the 6 colors, so the number of ways to fill these positions is $(6 \cdot 18)!$. So, the total number of ways for Eeyore is: $6! \cdot \left(\frac{24}{6}\right) \cdot 6! \cdot 18!$. **Step 3:** Compare the two results. We need to compare W and E to see who has more ways to implement their plan. We can calculate the ratio of W to E to find out how many times more Winnie-the-Pooh has compared to Eeyore. Let's calculate this using Python: `python import math W = math.factorial(7) * (7**18) # Calculate W (Winnie-the-Pooh's ways) E = math.factorial(6) * (6**18) # Calculate E (Eeyore's ways) ratio = W / E print(ratio) # Output 5.1000000000000000` The calculator shows that Winnie-the-Pooh has 5 times more ways to implement his plan than Eeyore. Therefore, the answer is $\boxed{5}$.

(b) Incorrect answer
x1.33

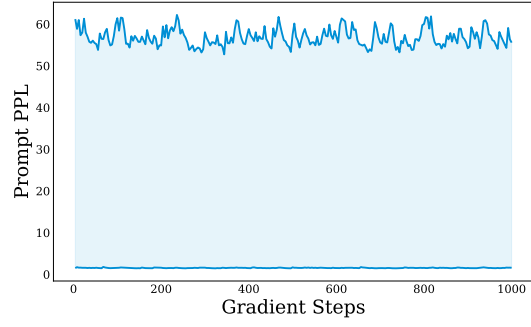
To solve this problem, we need to find the radius of the largest possible circle that fits inside or on the boundary of a quadrilateral $(ABCD)$ with given side lengths $(AB = 14)$, $(BC = 9)$, $(CD = 7)$, and $(DA = 12)$. This radius is known as the inradius of the quadrilateral. A quadrilateral that has an inscribed circle is called a tangential quadrilateral. For a tangential quadrilateral, the sum of the lengths of opposite sides must be equal. Let's check if this condition is satisfied for the given side lengths: $AB + CD = 14 + 7 = 21$, $BC + DA = 9 + 12 = 21$. Since $AB + CD = BC + DA$, the quadrilateral $(ABCD)$ is indeed a tangential quadrilateral. The area (K) of a tangential quadrilateral can be expressed in terms of its inradius (r) and its semiperimeter (s) : $K = r \cdot s$ where the semiperimeter (s) is given by: $s = \frac{AB + BC + CD + DA}{2} = \frac{14 + 9 + 7 + 12}{2} = 21$. To find the area (K) , we can use Brahmagupta's formula for the area of a cyclic quadrilateral. However, since we are unsure if the quadrilateral is cyclic, we will use another method. We can split the quadrilateral into two triangles by drawing a diagonal, and use the formula for the area of a triangle. Let's consider splitting the quadrilateral along diagonal (AC) . We can find the area of triangles (ABC) and (ACD) separately and sum them up. We will use Heron's formula for each triangle. Let's find the area of triangle (ABC) . The semiperimeter (s_1) of triangle (ABC) is: $s_1 = \frac{AB + BC + AC}{2}$. We do not know (AC) , but we can express the area in terms of (AC) and use Heron's formula: $K_{ABC} = \sqrt{s_1(s_1 - AB)(s_1 - BC)(s_1 - AC)}$. Similarly, for triangle (ACD) : $s_2 = \frac{AC + CD + DA}{2}$. $K_{ACD} = \sqrt{s_2(s_2 - AC)(s_2 - CD)(s_2 - DA)}$. The total area (K) of the quadrilateral is then: $K = K_{ABC} + K_{ACD} = \sqrt{s_1(s_1 - AB)(s_1 - BC)(s_1 - AC)} + \sqrt{s_2(s_2 - AC)(s_2 - CD)(s_2 - DA)}$. We can find the correct value of (AC) that maximizes the area. However, for simplicity, we can use Brahmagupta's formula for a tangential quadrilateral: $K = \sqrt{(s - AB)(s - BC)(s - CD)(s - DA)}$. Substituting the values: $K = \sqrt{(21 - 14)(21 - 9)(21 - 7)(21 - 12)}$. $K = \sqrt{7 \cdot 12 \cdot 12 \cdot 9} = \sqrt{10584} = 102.876$. Now, we can find the inradius (r) : $r = \frac{K}{s} = \frac{102.876}{21} = 4.9$. However, this does not match the form $(k \cdot \sqrt{m} + n)$. We need to check if the area calculator can be simplified to match the given form. Let's re-check the calculation: `python import sympy as sp # Given side lengths AB = 14 BC = 9 CD = 7 DA = 12 # Semiperimeter s = (AB + BC + CD + DA) / 2 # Area using Brahmagupta's formula for a tangential quadrilateral K = sp.sqrt((s - AB) * (s - BC) * (s - CD) * (s - DA)) # Inradius r = K / s print(K, r) # Output (84.00000000000000, 4.000000000000000)` The area (K) of the quadrilateral is (84) and the inradius (r) is (4) . The form $(k \cdot \sqrt{m} + n)$ seems to be a special case where $(k = 0)$ and $(n = 4)$, but the problem states that the answer should be in the form $(k \cdot \sqrt{m} + n)$. Given that (4) can be written as $(0 \cdot \sqrt{1} + 4)$, the values of (k) , (m) , and (n) are (0) , (1) , and (4) , respectively. Therefore, $(k + m + n = 0 + 1 + 4 = 5)$. The final answer is: $\boxed{5}$.

(d) Incorrect answer
x0.74

Figure 15: Token-Level Entropy of Sequences within a Mini-batch (with Relative-Entropy in Red).



(a) Qwen2.5-Math-7B



(b) Qwen2.5-Math-1.5B

Figure 16: Range of prompt PPL during training.

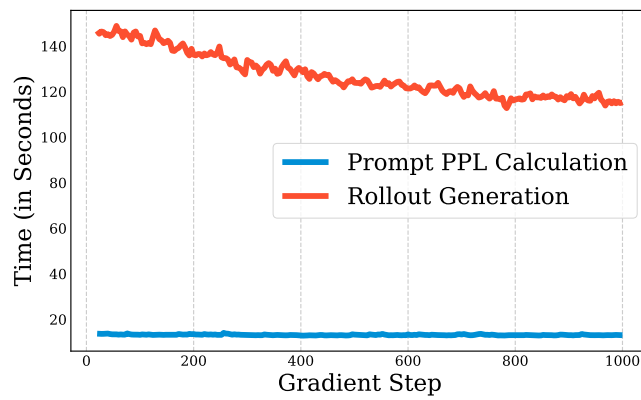
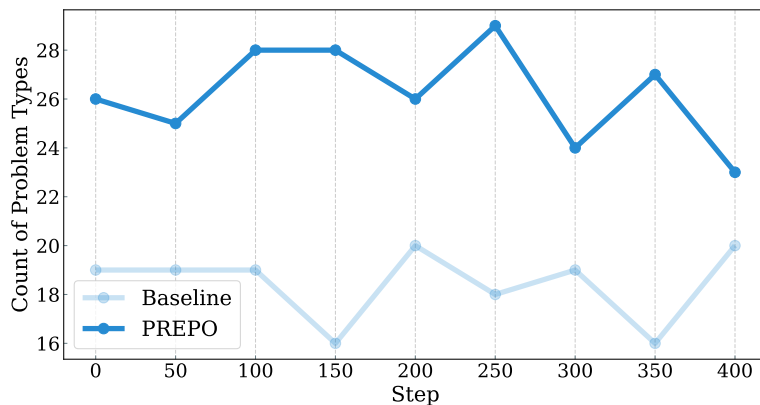
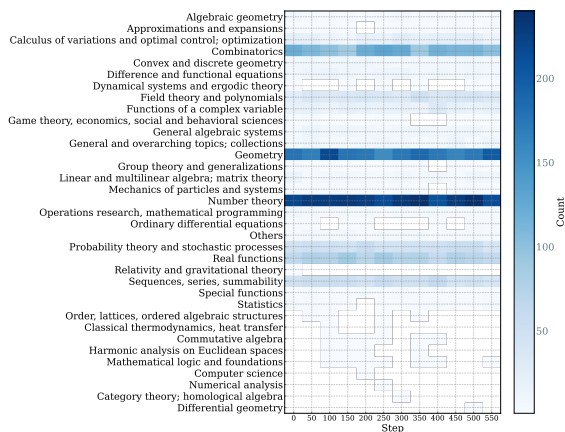


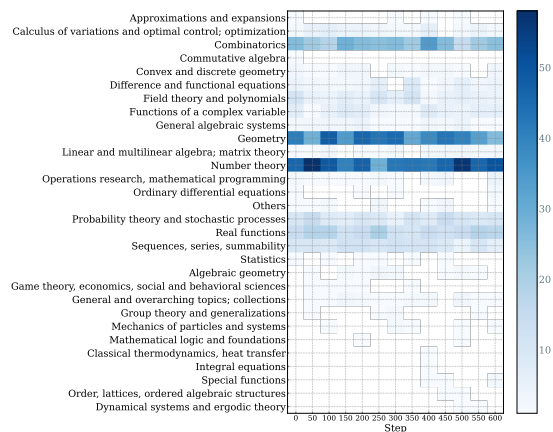
Figure 17: Comparison of calculating prompt PPL and rollout generation



(a) Unique number of MSC tags



(b) MSC frequency of PREPO



(c) MSC frequency of baseline (random)

Figure 18: Comparison of problem diversity between PREPO and baseline on Qwen2.5-Math-7B.

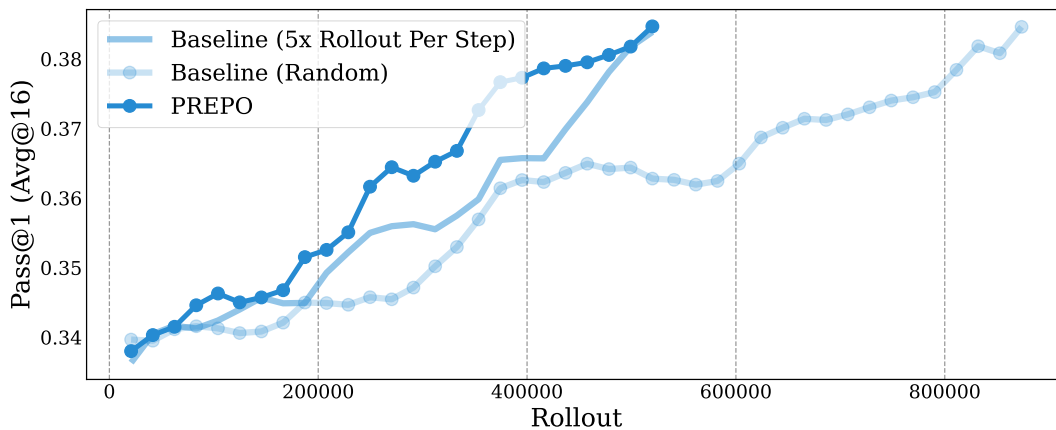


Figure 19: Comparison of PREPO and baselines (a) training w/o filtering (b) random selection.

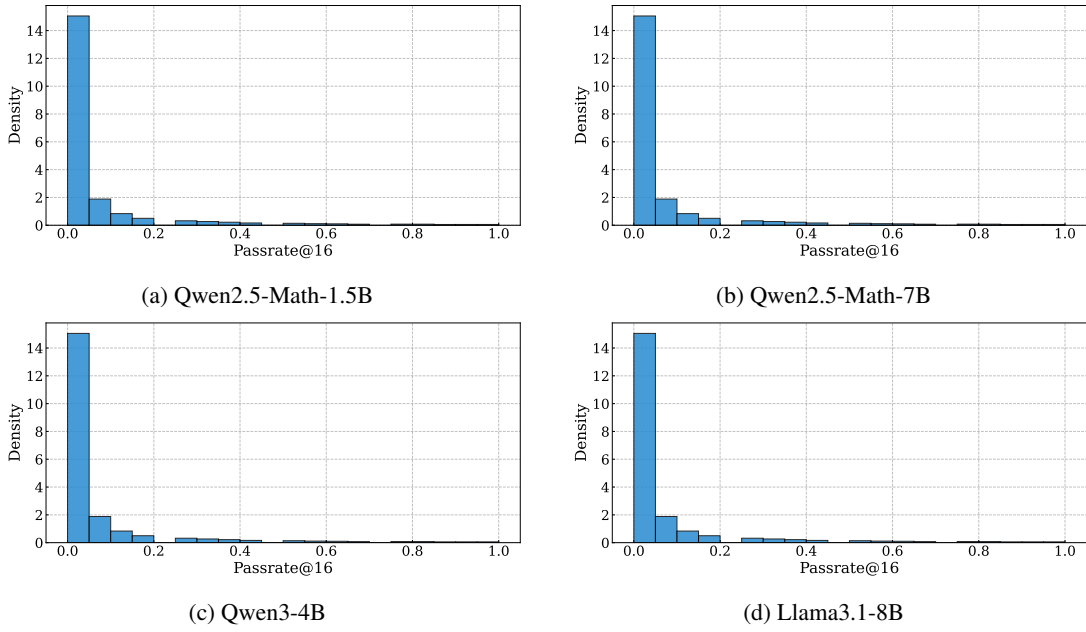


Figure 20: Distribution of passrate@16 over all partial prompts

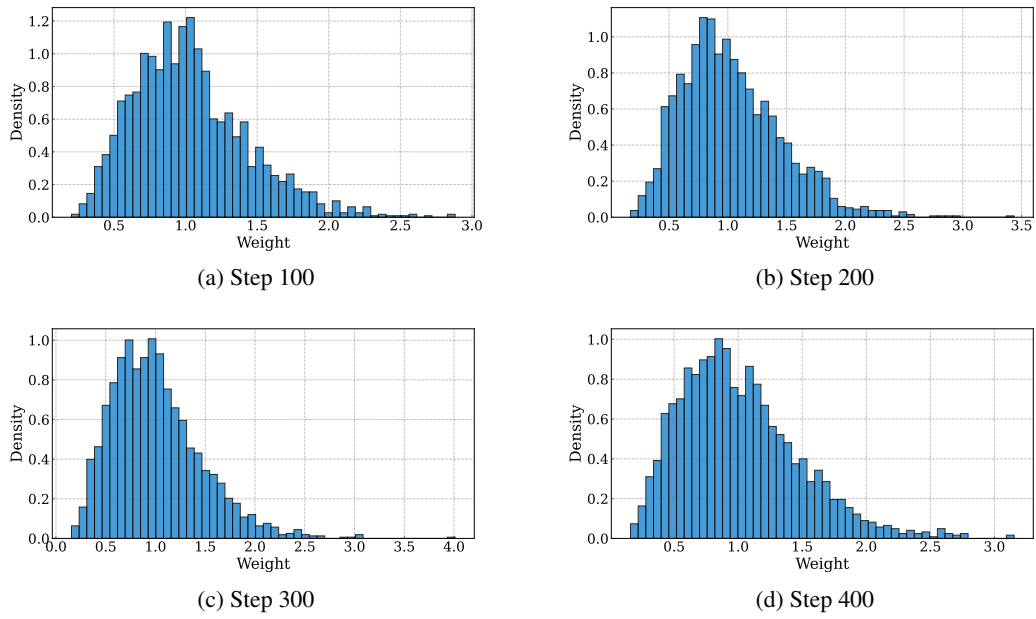


Figure 21: Weight Distribution During Training (Qwen2.5-Math-1.5B).