

Beyond High-Entropy Exploration: Correctness-Aware Low-Entropy Segment-Based Advantage Shaping for Reasoning LLMs

Anonymous ACL submission

Abstract

Reinforcement Learning with Verifiable Rewards (RLVR) has become a central approach for improving the reasoning ability of large language models. Recent work studies RLVR through token entropy, arguing that high-entropy tokens drive exploration and should receive stronger updates. However, they overlook the fact that most of a reasoning trajectory consists of low-entropy segments that encode stable and reusable structural patterns. Through qualitative and quantitative analyses, we find that the overlap of low-entropy segments across correct responses strongly correlates with model accuracy, while overlaps involving incorrect responses exhibit stable but unproductive patterns. Motivated by these findings, we propose LESS, a correctness-aware reinforcement framework that performs fine-grained advantage modulation over low-entropy segments. LESS amplifies segments unique to correct responses, suppresses those unique to incorrect ones, and neutralizes segments shared by both, while preserving high-entropy exploration in the underlying RL algorithm. Instantiated on top of GRPO and GSPO, LESS not only improves accuracy over strong RL baselines across three backbones and six math benchmarks, but also achieves stronger robustness of the performance floor.

1 Introduction

The reasoning capability of Large Language Models (LLMs) plays a central role in tasks such as mathematics (DeepSeek-AI et al., 2025; Chen et al., 2024), programming (5 Team et al., 2025; Wei et al., 2025b; Da et al., 2025), science problem-solving (M2 Team et al., 2025; Sellergren et al., 2025; Jing et al., 2025), and legal analysis (Zhang et al., 2025a,b). Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as an effective approach for improving reasoning reliability, where the correctness of the final answer is used as a reward signal to update the model. Representative

RLVR methods (Shao et al., 2024; Yu et al., 2025) typically apply policy updates uniformly across all tokens in a generated sequence.

Recent studies have argued that different parts of a reasoning sequence contribute differently to the final outcome through the lens of token entropy, and that RLVR training should take this into account. They observe that high-entropy tokens often correspond to exploratory reasoning steps, where the model tests alternative solution paths. For example, Cui et al. (2025) show that training encourages the model to explore uncertain reasoning branches; Zhang et al. (2025) encourage diversity in correct attempts by adjusting update strength in high-entropy regions; and Cheng et al. (2025) show that increasing entropy can improve the ability to search for solutions. Most notably, Wang et al. (2025) study demonstrates that only a small subset of tokens with high entropy disproportionately influence reasoning outcomes, suggesting that RL training should focus attention on these regions.

While existing entropy-based approaches have shown promising results, they focus almost entirely on high-entropy tokens, treating these points as the main drivers of exploration in reasoning. This overlooks that most of a reasoning sequence is composed of low-entropy segments (e.g., Wang et al. (2025) treat about 80% of tokens in a response as low-entropy), which form the stable structural scaffold that shapes how the solution is carried out. To examine the role of these low-entropy segments, we conduct both qualitative and quantitative analyses on RLVR-trained models, which is elaborated in §2. From the preliminary studies, we have the following observations. **First**, as shown in Fig 1(a), correct responses share consistent low-entropy segments that reflect coherent and productive reasoning steps, while incorrect responses also display their own repeated low-entropy patterns that represent stable but unproductive reasoning habits. **Second**, cross-dataset evaluation

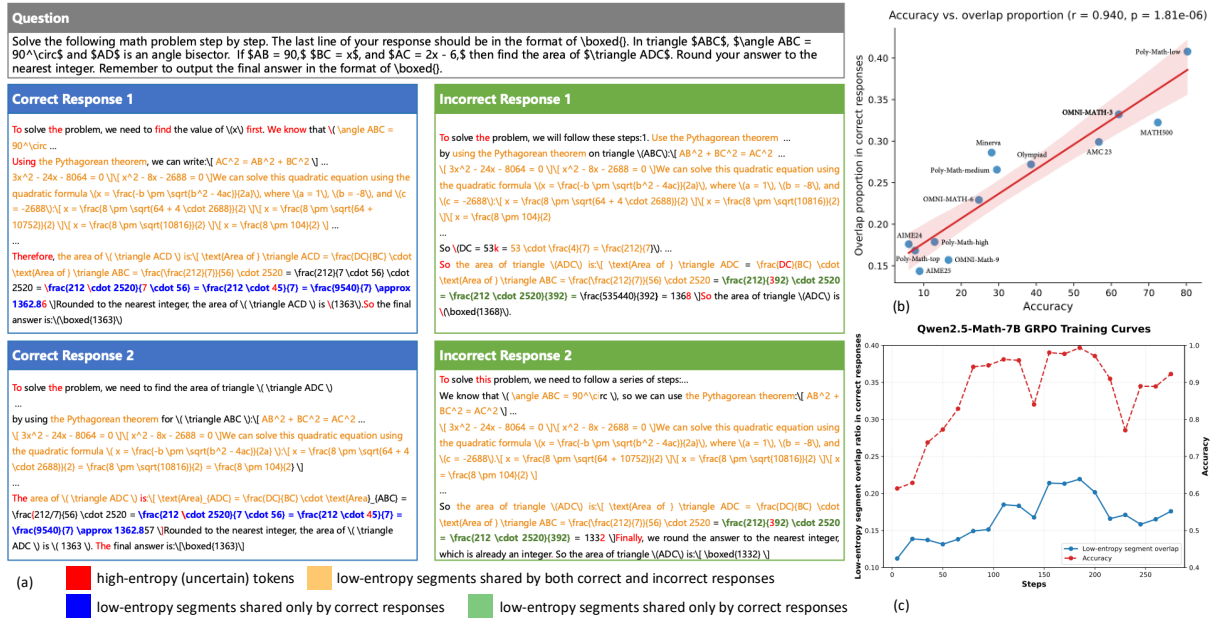


Figure 1: Low-entropy analysis reveals stable reasoning behaviors. **Left:** a case study where correct and incorrect responses exhibit shared and distinct low-entropy segments. **Right-top:** Across math benchmarks, accuracy strongly correlates with low-entropy segment overlap in correct responses. **Right-bottom:** During GRPO training of Qwen2.5-Math-7B, both accuracy and low-entropy overlap rise together, showing that performance gains emerge alongside the stabilization of reasoning patterns.

in Fig 1(b) shows a strong positive correlation between Qwen2.5-Math-7B accuracy and the overlap of low-entropy segments in correct responses, indicating that this relationship holds across tasks and model setups. **Third**, from the training dynamics in Fig 1(c), we observe that the model accuracy and the overlap of low-entropy segments across correct responses increase together, showing that improvements in reasoning ability are accompanied by the consolidation of shared structural patterns. These observations suggest that simply emphasizing high-entropy regions is insufficient, and the treatment of low-entropy segments is directly related to whether useful or harmful reasoning routines are reinforced.

Building on these observations, we introduce **Low-Entropy Segment Shaping (LESS)**, a reinforcement learning with verifiable rewards framework that treats low-entropy structure as an explicit training signal. LESS inserts an entropy-aware segmentation step into the policy update. For each generated trajectory, it splits the trajectory into high-entropy exploration tokens and contiguous low-entropy segments, and then aggregates how often each segment appears in correct versus incorrect responses within a rollout group. These statistics are used to rescale token-level advantages in a structured way: segments that occur only in correct trajectories receive amplified positive ad-

vantages, segments that occur only in incorrect trajectories receive amplified negative advantages, segments that co-occur in both are neutralized, while high-entropy tokens keep their original RL updates. In this way, LESS strengthens reusable reasoning structure and suppresses repeated failure patterns without harming exploration. The framework is agnostic to the underlying RLVR objective. In this work, we instantiate LESS on top of two widely used multi-sample training methods—Group Relative Policy Optimization (GRPO) (Shao et al., 2024) and Group Sequence Policy Optimization (GSPO) (Zheng et al., 2025). Extensive experiments conducted on six reasoning benchmark demonstrate LESS outperforms popular baseline across almost all tasks and model scales (1.5B, 7B math-tuned, and 7B base). In particular, it yields notable improvements on AIME24/25 and AMC23, where stable multi-step reasoning is essential. Moreover, compared to vanilla GRPO, LESS markedly reduces the worst-case dispersion among sampled responses. This aligns with the core goal of LESS—to reinforce beneficial structural segments and suppress misleading ones—ultimately producing more stable and reliable policy updates.

In summary, our contributions are three-fold: (1) We introduce a segment-level perspective on RLVR that distinguishes low-entropy segments

141 by correctness, revealing stable structural patterns
142 in LLM reasoning. (2) We propose LESS, a plug-
143 in algorithm that reweights token-level advantages
144 by segment statistics, amplifying low-entropy seg-
145 ments unique to correct trajectories, suppressing
146 those unique to incorrect ones, and neutralizing
147 shared segments. (3) We show that LESS consis-
148 tently improves accuracy across six mathematical
149 reasoning benchmarks and three backbones, while
150 also improving robustness under worst@K and re-
151 ducing variance across sampled rollouts.

152 2 Preliminary Analysis

153 We examine low-entropy segments as an indicator
154 of stable reasoning behavior in LLMs and study
155 how these signals relate to model correctness and
156 performance.

157 We begin by conducting a qualitative experiment
158 to visualize the entropy structure of multiple re-
159 sponses produced for the same question. Specifi-
160 cally, we analyze the token-level entropy patterns
161 of Qwen2.5-Math-7B on the mathematical reason-
162 ing dataset (Hendrycks et al., 2021), which enables
163 us to separate responses into stable and unstable re-
164 gions and to identify the parts of the model’s reason-
165 ing that remain consistently preserved across differ-
166 ent outputs. As shown in Fig. 1 (a), high-entropy
167 tokens (in red) mark unstable regions where the
168 model varies its reasoning, while low-entropy seg-
169 ments reveal stable structures that the model consis-
170 tently reuses. Within these low-entropy segments,
171 we observe three distinct patterns. (1) Segments
172 shared only by correct responses (in blue) corre-
173 spond to productive reasoning steps that reliably
174 support the correct solution. (2) Segments shared
175 only by incorrect responses (in green) reflect stable
176 but unproductive reasoning habits, for example, the
177 repeated computation “ $\frac{212}{392} \cdot 2520 = \frac{212 \cdot 2520}{392} = ?$ ”.
178 (3) Segments shared by both correct and incorrect
179 responses (in orange) capture general reasoning
180 components that are stable but not predictive of cor-
181 rectness—for instance, invoking “the Pythagorean
182 theorem,” which provides a common derivation
183 framework but is not the source of the subsequent
184 correct or incorrect calculations. This evidence
185 shows low-entropy segments encode structured rea-
186 soning behaviors that differentiate effective and
187 ineffective model responses.

188 To test whether the qualitative patterns extend
189 beyond a single example, we measure the over-
190 lap of low-entropy segments across correct re-

191 sponses for a range of math benchmarks, includ-
192 ing AIME24, AIME25, AMC23, MATH500, Min-
193 erva, and the Omni-MATH series. As shown in
194 Fig.1 (b), benchmark accuracy is strongly cor-
195 related with the degree of low-entropy overlap
196 across correct responses (Pearson $r = 0.94$ and
197 p -value = $1.81e^{-6}$). Benchmarks with higher ac-
198 curacy, such as MATH500 and Omni-MATH-3, ex-
199 hibit clear clustering toward higher overlap ratios,
200 while lower-accuracy benchmarks show weaker
201 consistency in their stable segments. The fitted
202 regression line further highlights this trend, indi-
203 cating that stronger task performance is associated
204 with more consolidated reasoning structure and
205 greater reuse of stable low-entropy patterns. Simi-
206 lar positive correlations (Appendix A.6) are ob-
207 served across several other backbones.

208 We further examine how these patterns evolve
209 during learning. Using GRPO training of Qwen2.5-
210 Math-7B, we track both accuracy and low-entropy
211 segment overlap over training steps. As shown
212 in Fig. 1 (c), the two trajectories rise together
213 throughout training: early stages display low accu-
214 racy and fragmented low-entropy structure, while
215 later stages show increasing stability in low-entropy
216 segments alongside improved accuracy. This syn-
217 chronous growth suggests that the model’s reason-
218 ing becomes more consistent as training progresses
219 and that stable low-entropy segments emerge as the
220 model acquires more reliable reasoning routines.
221 These results confirm that low-entropy overlap re-
222 flects not only final performance but also the de-
223 velopmental trajectory of the model’s reasoning
224 behavior. We observe the same co-evolution pat-
225 tern (Appendix A.5) on Qwen2.5-Math-1.5B and
226 Qwen2.5-7B.

227 These results show that low-entropy segments
228 provide a reliable signal for understanding and
229 guiding model reasoning. They capture stable com-
230 putational routines that distinguish correct from
231 incorrect behavior, reflect the degree of structural
232 consistency across benchmarks, and track the de-
233 velopment of reasoning stability during learning.
234 These observations suggest low-entropy segments
235 can serve as informative targets for optimization,
236 enabling the model to strengthen productive reason-
237 ing routines while suppressing unproductive ones.

238 3 Methodology

239 Motivated by these findings, we propose Low-
240 Entropy Segment Shaping (LESS), an RLVR frame-

work that improves reasoning stability by reshaping token-level advantages using statistics of low-entropy segments across sampled responses. LESS is compatible with standard RLVR algorithms and can be used as a plug-in module, and we instantiate it on top of GRPO and GSPO in this work.

3.1 LESS: Low-Entropy Segment Shaping

Given an input question q , the policy generates a group of responses $\mathcal{G} = \{O_1, \dots, O_G\}$, LESS detects low-entropy segments and shapes corresponding advantages as follows:

Entropy-based segment extraction. For a $O_i = [t_1, \dots, t_{|O_i|}]$, the entropy of token t_j is

$$\mathcal{H}_{t_j} = - \sum_{v \in V} \pi_{\theta_{\text{old}}}(v | x, O_{i < j}) \log \pi_{\theta_{\text{old}}}(v | x, O_{i < j}). \quad (1)$$

Following Wang et al. (2025), we compute an entropy threshold τ_i for each response O_i as the h -quantile of its token entropies \mathcal{H}_{t_j} . We then treat high-entropy tokens as isolated positions and group consecutive low-entropy tokens into contiguous spans. A minimum length μ is used to filter out trivial low-entropy spans (such as punctuation or very short frequent phrases). This gives three types of entropy-based structures:

$$\begin{aligned} \mathcal{S}_i^{\text{high}} &= \{t_j \in O_i \mid \mathcal{H}_{t_j} \geq \tau_i\}, \\ \mathcal{S}_i^{\text{frag}} &= \{O_i[a:b] \mid b - a + 1 < \mu, \forall t_j \in [a, b] : \mathcal{H}_{t_j} < \tau_i\}, \\ \mathcal{S}_i^{\text{seg}} &= \{O_i[a:b] \mid b - a + 1 \geq \mu, \forall t_j \in [a, b] : \mathcal{H}_{t_j} < \tau_i\}, \end{aligned} \quad (2)$$

where $\mathcal{S}_i^{\text{high}}$ collects individual high-entropy tokens, $\mathcal{S}_i^{\text{frag}}$ contains short low-entropy fragments that are likely uninformative, and $\mathcal{S}_i^{\text{seg}}$ contains longer low-entropy segments that we regard as structured reasoning candidates. Then, we aggregate how often each $\mathcal{S}_i^{\text{seg}}$ appears in correct versus incorrect responses within a rollout group. Let N_r and N_w denote the number of correct and incorrect responses in \mathcal{G} . For a low-entropy segment σ , we count its frequency over the group:

$$\begin{aligned} n_r(\sigma) &= |\{i \mid \text{correct}_i = 1 \wedge \sigma \in \mathcal{S}_i^{\text{seg}}\}|, \\ n_w(\sigma) &= |\{i \mid \text{correct}_i = 0 \wedge \sigma \in \mathcal{S}_i^{\text{seg}}\}|. \end{aligned} \quad (3)$$

Advantage shaping. LESS modifies the advantage assigned to each token t_j in O_i as:

$$\hat{A}_j^i = \begin{cases} A_i, & t_j \in \mathcal{S}_i^{\text{high}}, \\ A_i/N_r, & t_j \in \mathcal{S}_i^{\text{frag}}, \text{correct}_i = 1, \\ A_i/N_w, & t_j \in \mathcal{S}_i^{\text{frag}}, \text{correct}_i = 0, \\ 0, & \sigma_j^i \in \mathcal{S}_i^{\text{seg}}, n_r > 0, n_w > 0, \\ (n_r/N_r)A_i, & \sigma_j^i \in \mathcal{S}_i^{\text{seg}}, n_r > 0, n_w = 0, \\ (n_w/N_w)A_i, & \sigma_j^i \in \mathcal{S}_i^{\text{seg}}, n_r = 0, n_w > 0. \end{cases} \quad (4)$$

This rule: (i) preserves exploratory high-entropy behavior, (ii) reinforces stable segments unique to correct responses, (iii) penalizes those unique to incorrect responses, (iv) ignores ambiguous segments shared by both groups. The full pseudocode of LESS and its time complexity analysis are given in Appendix A.1.

3.2 Instantiations

LESS is designed as a generic advantage-shaping framework and can be applied to current RLVR methods that computes token- or sequence-level advantages. In this work, we instantiate LESS using the GRPO (Shao et al., 2024) and GSPO (Zheng et al., 2025). These methods are particularly suitable for our setting because it (i) generates a group of responses for each query, allowing entropy-based statistics to be computed across samples, and (ii) performs stable clipped-ratio updates that interact well with our advantage shaping.

LESS with GRPO. For an input query q , the policy produces G responses with rewards $\{r_1, \dots, r_G\}$. GRPO standardizes these rewards to obtain group-relative advantages:

$$A_i = \frac{r_i - \text{mean}(r_{1:G})}{\text{std}(r_{1:G})}. \quad (5)$$

GRPO then updates the policy by maximizing a clipped likelihood-ratio objective regularized by a KL constraint toward a reference policy:

$$J_{\text{GRPO}}(\theta) = \mathbb{E} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} (\min(\alpha_i A_i, \tilde{\alpha}_i A_i) - \kappa_i) \right]. \quad (6)$$

$$\text{where } \alpha_i = \frac{\pi_{\theta}(o_i | x)}{\pi_{\theta_{\text{old}}}(o_i | x)}, \quad \tilde{\alpha}_i = \text{clip}(\alpha_i, 1 - \epsilon, 1 + \epsilon),$$

$$\kappa_i = \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}), \quad (6)$$

GRPO’s group-wise credit assignment is well aligned with LESS, since the same group of responses used for reward normalization is also used by LESS to compute low-entropy statistics. Replacing A_i in Eq. 6 with the shaped advantage \hat{A}_j^i in Eq. 4 yields our LESS-GRPO training objective.

LESS with GSPO. GSPO differs from GRPO by performing importance weighting and clipping at the sequence level. Given the same group of responses $\{o_i\}_{i=1}^G$, GSPO defines the sequence-level importance ratio and optimizes

$$J_{\text{GSPO}}(\theta) = \mathbb{E} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} (\min(\beta_i A_i, \tilde{\beta}_i A_i) - \kappa_i) \right]. \quad (7)$$

$$\text{where } \beta_i = \left(\frac{\pi_{\theta}(o_i | x)}{\pi_{\theta_{\text{old}}}(o_i | x)} \right)^{\frac{1}{|o_i|}}, \quad \kappa_i = \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}),$$

$$\tilde{\beta}_i = \text{clip}(\beta_i, 1 - \epsilon, 1 + \epsilon), \quad (7)$$

Similarly, to integrate LESS with GSPO, we retain GSPO’s sequence-level importance weighting but replace the scalar group advantage A_i Eq.7 with token-level shaped advantages in Eq.4.

4 Experiments

We answer the following research questions with experiments: **RQ1**: How does LESS affect performance across diverse benchmarks when applied on top of standard RLVR algorithms (GRPO and GSPO), compared with strong baselines? **RQ2**: How does LESS influence the training dynamics of LLM reasoning, compared with GRPO, in terms of accuracy growth, stability. **RQ3**: How do LESS and GRPO differ in the evolution of entropy-based reasoning structures during training? **RQ4**: Does LESS improve worst-case reasoning robustness compared with GRPO across different model sizes? **RQ5**: How sensitive is LESS to the minimum segment-length μ , and how does varying μ affect the stability and final accuracy of reinforcement-learning-based reasoning?

4.1 Experimental Setup

Datasets and evaluation metrics. Following (Shen, 2025), we train the models on the MATH dataset (Hendrycks et al., 2021), which contains 7,500 problems spanning algebra, geometry, counting, probability, number theory, and other areas. The dataset is widely adopted in LLM reasoning research due to its breadth and the step-wise reasoning it elicits, making it particularly suitable for entropy-based structure analysis.

In terms of the evaluation, we assess the trained models on a suite of standard mathematical reasoning benchmarks: MATH500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), OlympiadBench (He et al., 2024), AIME’24, and AIME’25 (LI et al., 2024). These datasets collectively cover varying difficulty levels and reasoning types, allowing us to examine whether LESS consistently improves reasoning stability. For all benchmarks except AIME, we report accuracy under greedy decoding, which is commonly used in math reasoning evaluation. For AIME’24/25, we follow prior works (Yu et al., 2025; Zheng et al., 2025; Yue et al., 2025) and compute the avg@32 accuracy by averaging predictions over 32 sampled rollouts. This protocol reduces the variance introduced by integer-answer formats and ensures fair comparison across RL-trained models.

Backbone LLM and baselines. We evaluate LESS on three Qwen2.5 family: Qwen2.5-Math-1.5B, Qwen2.5-Math-7B, and the general-purpose Qwen2.5-7B model (Qwen et al., 2024). These models allow us to test LESS across (i) different parameter scales and (ii) models with and without domain-specific pretraining.

We compare LESS against strong RLVR systems including: **GRPO** (Shao et al., 2024), the canonical multi-sample policy-gradient method and the underlying backbone of many reasoning RL pipelines; **Forking Tokens** (Wang et al., 2025), an approach that identifies repeated reasoning fragments to adjust token-level credit assignment. **KL-Cov** (Cui et al., 2025), an entropy-based mechanism that modulates KL penalties using covariance between reward and token log-probs. These baselines represent the closest lines of work involving token-level structure, multi-sample variance reduction, and entropy-informed regularization, making them the most relevant comparisons for evaluating LESS.

Implementation details. We provide comprehensive implementation details in Appendix A.3. Our code is available at <https://anonymous.4open.science/r/LESS-2543/>.

4.2 Overall Performance

To answer **RQ1**, we conduct experiments on seven reasoning benchmarks across three Qwen2.5 backbones. Table 1 shows that LESS consistently outperforms all RL baselines on every backbone. On the math-specialized 1.5B model, LESS (GRPO) attains the highest average score, with clear gains on challenging tasks such as AIME24, AIME25, and AMC23, indicating that entropy-aware advantage shaping improves reliability over token-level clipping and forking-based updates. For the stronger 7B math model, LESS (GRPO) further raises the average to 47.7, with the largest margins on MATH500 and OlympiadBench, suggesting that the method strengthens multi-step reasoning structure rather than only local symbolic steps. On the 7B base model without math specialization, LESS (GRPO) still brings consistent gains, showing that the approach generalizes beyond math-aligned checkpoints. When instantiated with GSPO, LESS also yields improvements over vanilla GSPO across two 7B backbones, supporting its role as a generic credit-shaping module. Overall, LESS delivers the best average performance in every setting, improving both easy and hard rea-

Table 1: Overall performance on mathematical reasoning benchmarks. AIME24/25 are evaluated with avg@32; other benchmarks use @avg1. Bold numbers are the best in each column. ‘†’ indicates the model significantly outperforms all baseline models with paired t-tests at $p < 0.05$ level.

Method	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg
Qwen2.5-Math-1.5B							
Base LLM	6.2	3.8	37.5	58.6	15.8	26.5	24.7
Forking Tokens	21.6	7.0	60.0	75.6	29.0	38.5	38.6
KL-Cov	22.0	9.2	57.5	75.8	27.9	38.4	38.4
GRPO	21.6	6.6	57.5	74.8	27.2	39.6	37.8
LESS (GRPO)	26.2 †	12.4 †	62.5 †	75.2	30.8 †	39.6	41.1 †
Qwen2.5-Math-7B							
Base LLM	6.0	8.9	57.5	58.6	28.7	38.0	32.9
Forking Tokens	36.6	14.2	67.5	78.2	37.8	42.2	46.0
KL-Cov	35.2	14.1	70.0	78.6	38.2	43.7	46.6
GRPO	36.8	13.2	70.0	79.8	36.7	42.7	46.5
LESS (GRPO)	40.0 †	13.6	67.5	79.6	37.3	44.3	47.0
GRPO	33.3	13.8	65.0	79.8	38.6	44.9	45.9
LESS (GRPO)	36.0	15.6 †	70.0	81.6 †	37.8	45.7 †	47.7 †
Qwen2.5-7B							
Base LLM	3.2	5.2	37.5	53.4	18.0	23.9	23.5
Forking Tokens	20.2	12.3	62.5	77.8	37.5	40.6	41.8
KL-Cov	18.7	10.8	60.0	77.0	37.9	40.8	40.9
GRPO	18.3	11.2	62.5	78.4	36.3	40.2	41.1
LESS (GRPO)	19.2	11.6	62.5	77.8	37.8	41.8 †	41.8
GRPO	18.9	11.9	60.0	76.8	37.1	40.9	40.9
LESS (GRPO)	20.5 †	13.1 †	67.5 †	78.6 †	37.1	40.8	42.9 †

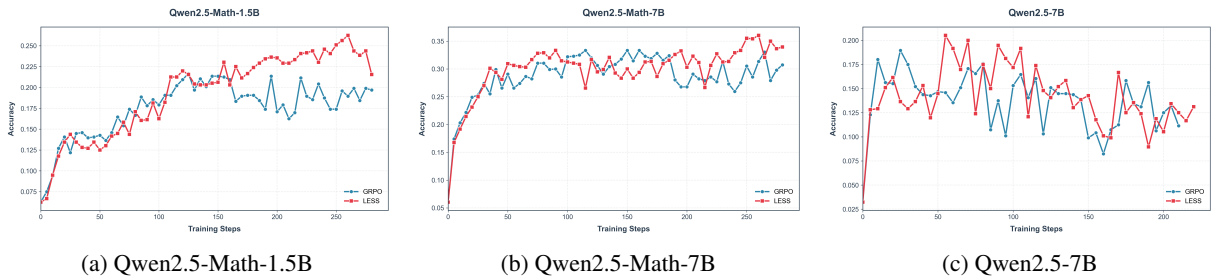


Figure 2: Training dynamics (accuracy over training) of GRPO and LESS across three backbones.

soning benchmarks. Unless otherwise stated, all detailed analyses in the following section use the GRPO-based instantiation of LESS.

4.3 Training Dynamics Analysis

To answer **RQ2**, we examine LESS and GRPO on how accuracy and entropy-based structure evolve throughout training. As shown in Fig. 2, across all three backbones, LESS exhibits a characteristic two-phase learning pattern. In the early stage, LESS improves slightly slower than GRPO because its advantage shaping reduces the update magnitude on low-entropy segments until the model accumulates enough evidence to distinguish productive from unproductive ones. However, as training progresses, LESS consistently surpasses GRPO

and maintains a higher accuracy plateau, indicating more stable policy improvement.

To answer **RQ3**, we analyze how the overlap of correct low-entropy segments and the entropy structure of incorrect responses evolve during reinforcement learning. Fig. 3 reports training dynamics on Qwen2.5-Math-7B; similar trends on Qwen2.5-Math-1.5B and Qwen2.5-7B are presented in Appendix A.4.

The top row of Fig. 3 shows that LESS consistently yields a higher overlap ratio of correct-only low-entropy segments as training progresses. This indicates that LESS explicitly amplifies structurally productive reasoning patterns that appear repeatedly in correct trajectories. In contrast, GRPO shows a flatter trend, suggesting that it does not

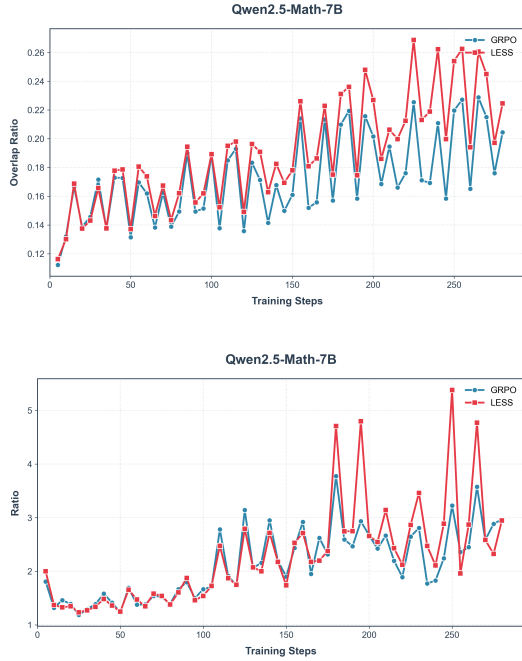


Figure 3: Training-dynamics comparison between LESS and GRPO on Qwen2.5-Math-7B. **Top:** Ratio of low-entropy segments that overlap exclusively among correct responses (higher is better). **Bottom:** Ratio between the entropy of incorrect responses and correct responses (higher indicates that incorrect answers remain exploratory). LESS consistently strengthens productive low-entropy structures while preventing premature entropy collapse in incorrect trajectories.

reliably consolidate these stable reasoning components. The clearer upward trajectory of LESS reveals that the model is progressively internalizing reusable, high-quality reasoning routines rather than relying on isolated or brittle solution paths.

The bottom row of Fig. 3 further shows that LESS maintains a higher entropy ratio between incorrect and correct responses, meaning that incorrect trajectories remain more uncertain. This separation is desirable, that is, LESS avoids prematurely stabilizing low-entropy segments that consistently lead to wrong answers, thereby reducing the risk of “locking in” systematic errors. GRPO, however, frequently collapses the entropy gap, causing incorrect responses to become low-entropy as well—an indication that harmful patterns are becoming entrenched in the policy.

These results show that LESS not only improves performance, but also progressively increases the overlap of correct low-entropy segments while keeping incorrect trajectories uncertain, creating a clear structural separation between productive and

Table 2: Worst-case reasoning performance (worst@ k) across three backbones. For each prompt, the worst-performing sample among k rollouts is selected and averaged over the dataset. LESS consistently improves worst-case accuracy across all settings.

Method	worst@32	worst@16	worst@8
<i>Qwen2.5-Math-1.5b</i>			
GRPO	6.8	8.0	10.8
LESS(GRPO)	12.9	15.6	18.2
<i>Qwen2.5-Math-7b</i>			
GRPO	13.4	17.4	20.6
LESS(GRPO)	21.2	22.2	24.3
<i>Qwen2.5-7B</i>			
GRPO	10.3	11.3	12.9
LESS(GRPO)	11.1	12.0	13.4

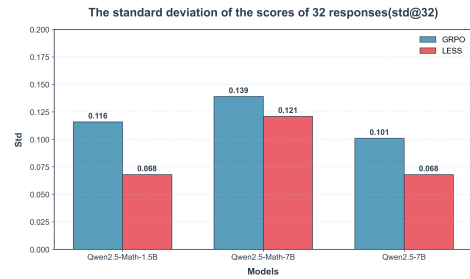


Figure 4: Standard deviation of 32 sampled responses (std@32). LESS reduces response-level variability across all backbones, indicating more stable and less volatile reasoning behavior compared with GRPO.

unproductive reasoning.

4.4 Robustness Under Worst-Case Sampling

To answer **RQ4** and examine the robustness of the learned policy, we employ the worst@ k metric, which selects the lowest-scoring output among k sampled responses and averages this worst-case score across prompts. This metric directly measures how the model behaves in its most vulnerable failure modes.

As shown in Table. 2, across all three backbones, LESS consistently improves worst-case accuracy. For Qwen2.5-Math-1.5B and Qwen2.5-Math-7B, LESS achieves substantial gains, raising worst@32 by +6.1 and +7.8 points respectively, with positive margins maintained as k decreases. This pattern shows that LESS not only lifts average performance but also strengthens the weakest trajectories, suppressing brittle low-entropy patterns that GRPO tends to reinforce. Even on the non-math Qwen2.5-7B model, LESS produces steady improvements, indicating that its robustness effects generalize be-

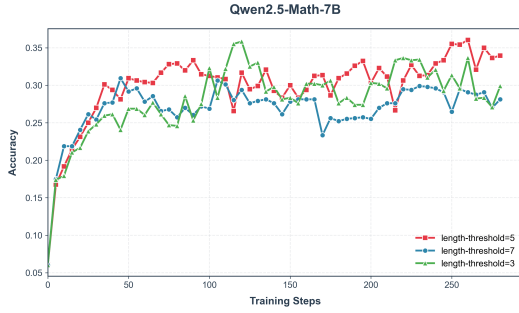


Figure 5: Effect of the low-entropy segment length threshold $\mu = \{3, 5, 7\}$ on training dynamics. We report accuracy over training steps.

yond specialized mathematical checkpoints.

The variance results shown in Fig. 4 further reinforce this finding: LESS consistently reduces the standard deviation of sampled rollouts ($\text{std}@32$), yielding more stable and predictable reasoning behavior. Together, the $\text{worst}@k$ and variance metrics demonstrate that LESS raises the floor of model performance while simultaneously mitigating response-level volatility.

These results show that LESS meaningfully improves robustness by raising the floor of model performance while simultaneously reducing instability across sampled rollouts.

4.5 Impact of the Segment Length

To answer **RQ5**, we study the sensitivity of LESS to the minimum segment-length threshold μ , we vary $\mu \in \{3, 5, 7\}$ and track training dynamics on Qwen2.5-Math-7B. As shown in Figure 5, $\mu = 5$ produces the most stable and highest final accuracy across the entire training trajectory. A very small threshold ($\mu = 3$) makes the model overly sensitive to short, noisy low-entropy fragments, causing the policy to reinforce many spurious local patterns and resulting in pronounced fluctuations. Conversely, a larger threshold ($\mu = 7$) filters out too many low-entropy segments, delaying the discovery of reliable reasoning motifs and slowing convergence.

The superior performance of $\mu = 5$ suggests that effective low-entropy guidance requires a balance: segments must be long enough to encode meaningful reasoning structure, yet short enough to capture fine-grained patterns that recur across correct trajectories. This indicates that LESS benefits from moderately sized structural units and is robust to reasonable choices of μ , but extremely small or large thresholds degrade the quality of structural signals made available to the policy.

5 Related Work

Token credit assignment. Vassoyan et al. (2025) identify critical tokens in chain-of-thought solutions— decision points where the model is likely to fail, and increase exploration around these tokens by adjusting the KL penalty. Lin et al. (2024) likewise locate tokens that strongly influence incorrect outcomes and show that editing or replacing these tokens can change the final decision. Other work (Chan et al., 2024; Xie et al., 2025; Guo et al., 2025; Wei et al., 2025a) addresses the coarse-grained nature of standard RL feedback by constructing dense, token-level rewards to resolve the credit assignment problem. These methods demonstrate that tokens within a trajectory should not be treated uniformly, but they still operate on local positions and do not capture how stable patterns repeat across multiple rollouts of the same question.

Entropy-based RL signals. Wang et al. (2025) split trajectories at high-entropy tokens and update only a subset of tokens, aiming to reduce over-optimization on already confident regions. Cui et al. (2025) further modulate the KL penalty based on token-level uncertainty, encouraging updates where the model is less certain and damping updates on low-entropy tokens. In these approaches, high-entropy tokens serve as a proxy for exploration, while low-entropy regions are treated as parts of the trajectory that should be protected from change. However, they do not distinguish low-entropy patterns that are consistently correct from those that encode repeated mistakes.

6 Conclusion

This paper presents a new perspective on training reasoning LLMs: reasoning should be guided at the level of low-entropy segments. Building on this insight, we propose LESS, a plug-and-play advantage-shaping framework that selectively amplifies reliable low-entropy reasoning segments and suppresses error-prone ones. Instantiated with GRPO and GSPO, LESS improves accuracy, stability, and robustness across multiple backbones and benchmarks. The framework consistently strengthens correct reasoning routines, preserves exploration on incorrect trajectories, and raises the worst-case performance of sampled rollouts. These results show that low-entropy structural signals offer a principled and effective handle for guiding RL training of reasoning models.

7 Limitations

We believe there is still room for improvement in our work. Our preliminary analysis and all main experiments are conducted on the Qwen2.5 family. We do not test LESS on other backbones such as Llama, so it is unclear whether the same entropy patterns and gains will hold more broadly. In addition, following Wang et al. (2025), we fix the entropy quantile and length threshold to extract low-entropy segments. We do not yet study how different quantiles, adaptive thresholds, or alternative segmentation rules would affect the learned segments and the final performance.

References

5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, and 152 others. 2025. *GLM-4.5: Agentic, Reasoning, and Coding (ARC) Foundation Models*. *arXiv e-prints*, arXiv:2508.06471.

Alex J. Chan, Hao Sun, Samuel Holt, and Mihaela van der Schaar. 2024. *Dense Reward for Free in Reinforcement Learning from Human Feedback*. *arXiv e-prints*, arXiv:2402.00782.

Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. 2024. *Not Everything is All You Need: Toward Low-Redundant Optimization for Large Language Model Alignment*. *arXiv e-prints*, arXiv:2406.12606.

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. 2025. *Reasoning with Exploration: An Entropy Perspective on Reinforcement Learning for LLMs*. *arXiv e-prints*, arXiv:2506.14758.

Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, Zhiyuan Liu, Hao Peng, Lei Bai, Wanli Ouyang, Yu Cheng, Bowen Zhou, and Ning Ding. 2025. *The Entropy Mechanism of Reinforcement Learning for Reasoning Language Models*. *arXiv e-prints*, arXiv:2505.22617.

Jeff Da, Clinton Wang, Xiang Deng, Yuntao Ma, Nikhil Barhate, and Sean Hendryx. 2025. *Agent-RLVR: Training Software Engineering Agents via Guidance and Environment Rewards*. *arXiv e-prints*, arXiv:2506.11425.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025.

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv e-prints*, arXiv:2501.12948.

Yiran Guo, Lijie Xu, Jie Liu, Dan Ye, and Shuang Qiu. 2025. *Segment Policy Optimization: Effective Segment-Level Credit Assignment in RL for Large Language Models*. *arXiv e-prints*, arXiv:2505.23564.

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. *OlympiadBench: A Challenging Benchmark for Promoting AGI with Olympiad-Level Bilingual Multimodal Scientific Problems*. *arXiv e-prints*, arXiv:2402.14008.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. *Measuring Mathematical Problem Solving With the MATH Dataset*. *arXiv e-prints*, arXiv:2103.03874.

Peiyuan Jing, Kinhei Lee, Zhenxuan Zhang, Huichi Zhou, Zhengqing Yuan, Zhifan Gao, Lei Zhu, Giorgos Papanastasiou, Yingying Fang, and Guang Yang. 2025. *Reason Like a Radiologist: Chain-of-Thought and Reinforcement Learning for Verifiable Report Generation*. *arXiv e-prints*, arXiv:2504.18453.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. *Solving Quantitative Reasoning Problems with Language Models*. *arXiv e-prints*, arXiv:2206.14858.

Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. 2024. *Numinamath*. <https://huggingface.co/AI-MO/NuminaMath-CoT>.

Zicheng Lin, Tian Liang, Jiahao Xu, Qiuzhi Lin, Xing Wang, Ruilin Luo, Chufan Shi, Siheng Li, Yujiu Yang, and Zhaopeng Tu. 2024. *Critical Tokens Matter: Token-Level Contrastive Estimation Enhances LLM’s Reasoning Capability*. *arXiv e-prints*, arXiv:2411.19943.

Ziru Liu, Cheng Gong, Xinyu Fu, Yaofang Liu, Ran Chen, Shoubo Hu, Suiyun Zhang, Rui Liu, Qingfu Zhang, and Dandan Tu. 2025. *GHPO: Adaptive Guidance for Stable and Efficient LLM Reinforcement Learning*. *arXiv e-prints*, arXiv:2507.10628.

M2 Team, Chengfeng Dou, Chong Liu, Fan Yang, Fei Li, Jiyuan Jia, Mingyang Chen, Qiang Ju, Shuai Wang, Shunya Dang, Tianpeng Li, Xiangrong Zeng, Yijie Zhou, Chenzheng Zhu, Da Pan, Fei Deng, Guangwei Ai, Guosheng Dong, Hongda Zhang, and 15 others. 2025. *Baichuan-M2: Scaling Medical Capability with Large Verifier System*. *arXiv e-prints*, arXiv:2509.02208.

693	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2024. Qwen2.5 Technical Report . <i>arXiv e-prints</i> , arXiv:2412.15115.	
694		
695		
696		
697		
698		
699		
700	Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau, Justin Chen, Fereshteh Mahvar, Liron Yatziv, Tiffany Chen, Bram Sterling, Stefanie Anna Baby, Susanna Maria Baby, Jeremy Lai, Samuel Schmidgall, and 62 others. 2025. MedGemma Technical Report . <i>arXiv e-prints</i> , arXiv:2507.05201.	
701		
702		
703		
704		
705		
706		
707		
708	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models . <i>arXiv e-prints</i> , arXiv:2402.03300.	
709		
710		
711		
712		
713		
714	Han Shen. 2025. On Entropy Control in LLM-RL Algorithms . <i>arXiv e-prints</i> , arXiv:2509.03493.	
715		
716		
717	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. HybridFlow: A Flexible and Efficient RLHF Framework . <i>arXiv e-prints</i> , arXiv:2409.19256.	
718		
719		
720		
721	Jean Vassoyan, Nathanaël Beau, and Roman Plaud. 2025. Ignore the KL penalty! boosting exploration on critical tokens to enhance RL fine-tuning . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 6108–6118, Albuquerque, New Mexico. Association for Computational Linguistics.	
722		
723		
724		
725		
726		
727	Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen Yu, Gao Huang, and Junyang Lin. 2025. Beyond the 80/20 Rule: High-Entropy Minority Tokens Drive Effective Reinforcement Learning for LLM Reasoning . <i>arXiv e-prints</i> , arXiv:2506.01939.	
728		
729		
730		
731		
732		
733		
734		
735	Quan Wei, Siliang Zeng, Chenliang Li, William Brown, Oana Frunza, Wei Deng, Anderson Schneider, Yuriy Nevmyvaka, Yang Katie Zhao, Alfredo Garcia, and Mingyi Hong. 2025a. Reinforcing Multi-Turn Reasoning in LLM Agents via Turn-Level Reward Design . <i>arXiv e-prints</i> , arXiv:2505.11821.	
736		
737		
738		
739		
740		
741	Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I. Wang. 2025b. SWE-RL: Advancing LLM Reasoning via Reinforcement Learning on Open Software Evolution . <i>arXiv e-prints</i> , arXiv:2502.18449.	
742		
743		
744		
745		
746		
747	Guofu Xie, Yunsheng Shi, Hongtao Tian, Ting Yao, and Xiao Zhang. 2025. CAPO: Towards Enhancing LLM Reasoning through Generative Credit Assignment . <i>arXiv e-prints</i> , arXiv:2508.02298.	
748		
749		
750		
	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. Qwen2.5-Math Technical Report: Toward Mathematical Expert Model via Self-Improvement . <i>arXiv e-prints</i> , arXiv:2409.12122.	751
		752
		753
		754
		755
		756
		757
		758
	Kun Yang, Zikang chen, Yanmeng Wang, and Zhigen Li. 2025. SSPO: Subsentence-level Policy Optimization . <i>arXiv e-prints</i> , arXiv:2511.04256.	759
		760
		761
	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. DAPO: An Open-Source LLM Reinforcement Learning System at Scale . <i>arXiv e-prints</i> , arXiv:2503.14476.	762
		763
		764
		765
		766
		767
		768
		769
	Yu Yue, Yufeng Yuan, Qiyang Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaye Chen, Chengyi Wang, Tiantian Fan, Zhengyin Du, Xiangpeng Wei, Xiangyu Yu, Gaohong Liu, Juncai Liu, Lingjun Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Chi Zhang, and 8 others. 2025. VAPO: Efficient and Reliable Reinforcement Learning for Advanced Reasoning Tasks . <i>arXiv e-prints</i> , arXiv:2504.05118.	770
		771
		772
		773
		774
		775
		776
		777
	Kepu Zhang, Guofu Xie, Weijie Yu, Mingyue Xu, Xu Tang, Yaxin Li, and Jun Xu. 2025a. Legal mathematical reasoning with LLMs: Procedural alignment through two-stage reinforcement learning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 1586–1598, Suzhou, China. Association for Computational Linguistics.	778
		779
		780
		781
		782
		783
		784
	Kepu Zhang, Weijie Yu, Zhongxiang Sun, and Jun Xu. 2025b. Syler: A framework for explicit syllogistic legal reasoning in large language models . In <i>Proceedings of the 34th ACM International Conference on Information and Knowledge Management, CIKM '25</i> , page 4117–4127, New York, NY, USA. Association for Computing Machinery.	785
		786
		787
		788
		789
		790
		791
	Xingjian Zhang, Siwei Wen, Wenjun Wu, and Lei Huang. 2025. EDGE-GRPO: Entropy-Driven GRPO with Guided Error Correction for Advantage Diversity . <i>arXiv e-prints</i> , arXiv:2507.21848.	792
		793
		794
		795
	Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization . <i>arXiv preprint arXiv:2507.18071</i> .	796
		797
		798
		799
		800
	Tianyu Zheng, Tianshun Xing, Qingshui Gu, Taoran Liang, Xingwei Qu, Xin Zhou, Yizhi Li, Zhoufutu Wen, Chenghua Lin, Wenhao Huang, Qian Liu, Ge Zhang, and Zejun Ma. 2025. First Return, Entropy-Eliciting Explore . <i>arXiv e-prints</i> , arXiv:2507.07017.	801
		802
		803
		804
		805
		806

A Appendix

A.1 Algorithm and Time Complexity Analysis

Algorithm 1 LESS: Low-Entropy Segment Shaping

```

1: Input:
2: Group of responses  $\mathcal{G} = \{O_1, \dots, O_G\}$ , token advantages  $\mathcal{A} = \{A_j^i \mid t_j \in O_i, i = 1, \dots, G\}$ , token entropies  $H = \{\mathcal{H}_j^i \mid t_j \in O_i, i = 1, \dots, G\}$ , correctness labels  $\{correct_1, \dots, correct_G\}$ , entropy quantile  $h$ , minimum segment length  $\mu$ .
3: Output: Shaped advantages  $\mathcal{A}' = \{\hat{A}_j^i \mid t_j \in O_i, i = 1, \dots, G\}$ 
4:  $N_r \leftarrow \sum_i \mathbb{I}[correct_i = 1]$ ;  $N_w \leftarrow \sum_i \mathbb{I}[correct_i = 0]$ 
5: for each response  $O_i$  do
6:   Compute entropy threshold  $\tau_i$  from  $\{\mathcal{H}_j^i\}_{t_j \in O_i}$  using quantile  $h$ 
7:   Segment  $O_i$  into  $\mathcal{S}_i^{\text{high}}$ ,  $\mathcal{S}_i^{\text{frag}}$ ,  $\mathcal{S}_i^{\text{seg}}$  using  $\tau_i$  and  $\mu$  (Eq. (2))
8: end for
9:  $\Sigma \leftarrow \emptyset$  {set of unique low-entropy segments}

10: for each response  $O_i$  do
11:   for all  $\sigma \in \mathcal{S}_i^{\text{seg}}$  do
12:     if no  $\sigma' \in \Sigma$  is a contiguous segments of  $\sigma$  then
13:        $\Sigma \leftarrow \Sigma \cup \{\sigma\}$ 
14:       Remove from  $\Sigma$  any  $\sigma'$  that is strictly contained in  $\sigma$ 
15:     end if
16:   end for
17: end for
18: for all  $\sigma \in \Sigma$  do
19:   Compute  $n_r(\sigma)$  and  $n_w(\sigma)$  according to Eq. (3)
20: end for
21: for each response  $O_i$  do
22:   for each token  $t_j \in O_i$  do
23:     Set  $\hat{A}_j^i$  according to Eq. (4)
24:   end for
25: end for
26: return  $\mathcal{A}' = \{\hat{A}_j^i\}$ 

```

The overall LESS procedure is summarized in Algorithm 1. Given a group of responses and their token-level entropies, we first compute an entropy threshold for each response and segment it into high-entropy tokens, short low-entropy fragments,

Table 3: Prompt template used for all experiments. {question} is replaced by the problem description.

Prompt Template

```

<|im start|>system
Please reason step by step, and put your final answer within \boxed{ }.
<|im end|>
<|im start|>user
{question}
<|im end|>
<|im start|>assistant

```

and longer low-entropy segments (Eq. (2)). We then build a set Σ of non-redundant low-entropy segments across the group by keeping only segments that are not strictly contained in longer ones. For each segment $\sigma \in \Sigma$, we count how many correct and incorrect responses it appears in (Eq. (3)), and finally assign a shaped advantage to every token based on its entropy category and the statistics of the segment it belongs to (Eq. (4)). The resulting token-level advantages \mathcal{A}' can be plugged into any group-based RLVR update.

In terms of complexity, when the batch size is B , the group size is G , and the maximum response length is L , the segmentation and shaping operations visit each token a constant number of times, giving a practical time complexity of $O(BGL)$. Under our main setting (batch size 512, group size 8, average response length about 800), LESS adds roughly 60 seconds of overhead in our implementation, which is small compared to the overall RL training time.

A.2 Prompt Template

We use a unified prompt template for all training and evaluation experiments, adapted from the official Qwen-Math template (Yang et al., 2024). The concrete format is shown in Table 3.

A.3 Implementation Details

We conduct experiments using the VeRL (Sheng et al., 2024) framework for reinforcement learning with LLMs on 8 NVIDIA A100-40G GPUs. The training setup includes a batch size of 512, a learning rate of 1×10^{-6} , and a clip range between 0.2 and 0.28. Each response sequence is up to 3k tokens in length. The mini-batch size is set to 32. The temperature is 1.0 for training and 0.1 for evaluation. Following prior RLVR work (Liu et al., 2025; Yang et al., 2025), we perform 8 rollouts per prompt and do not use entropy regularization or

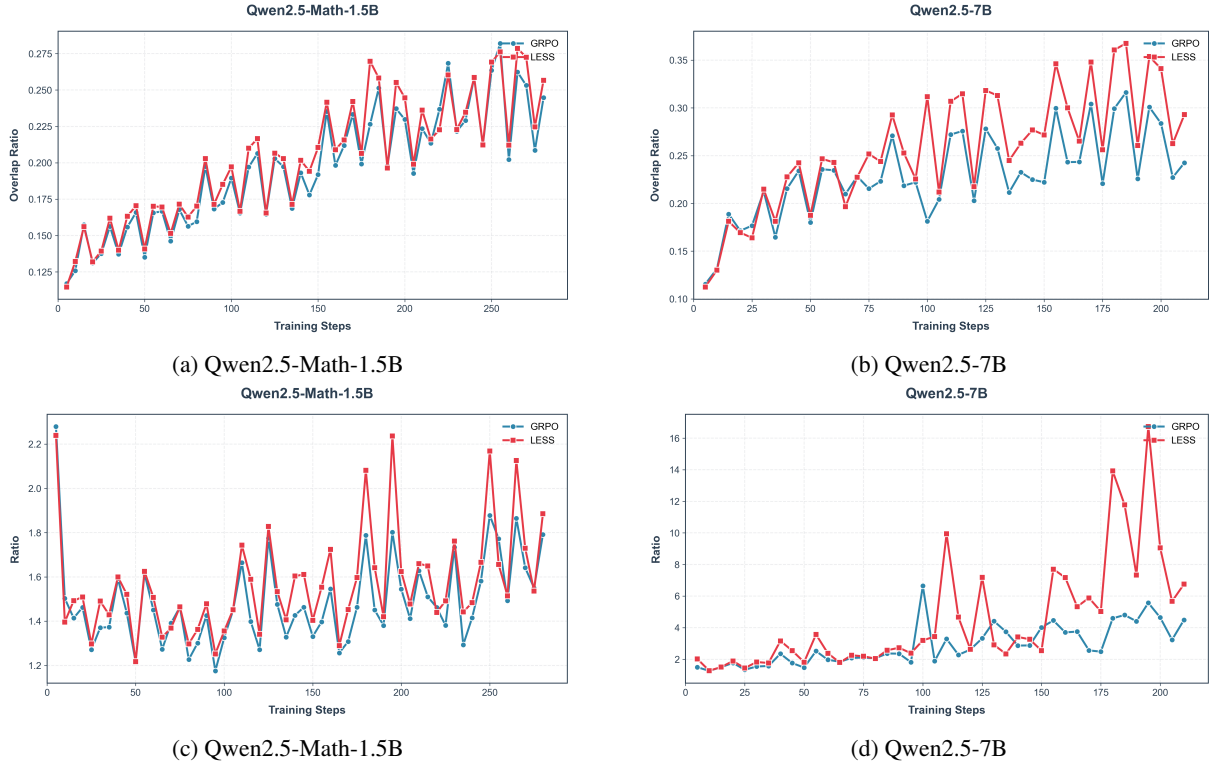


Figure 6: Training-dynamics comparison between LESS and GRPO across two model sizes. **Top:** Ratio of low-entropy segments that overlap exclusively among correct responses (higher is better). **Bottom:** Ratio between the entropy of incorrect responses and correct responses (higher indicates that incorrect answers remain exploratory).

852 KL penalties during training (KL coefficient = 0,
 853 entropy loss = 0), allowing us to isolate the effect
 854 of LESS from other types of regularization. This
 855 choice also reveals whether LESS alone can stabi-
 856 lize reasoning trajectories without relying on heavy
 857 KL anchoring. We use a unified prompt format for
 858 all experiments, and the exact template is provided
 859 in Appendix A.2.

860 A.4 Additional Training Dynamics Analysis

861 Fig. 6 presents the same training-dynamics analysis
 862 for Qwen2.5-Math-1.5B and Qwen2.5-7B. Across
 863 both backbones, LESS consistently increases the
 864 overlap of correct-only low-entropy segments dur-
 865 ing training and maintains a higher entropy ratio
 866 between incorrect and correct responses. These
 867 trends mirror the behavior observed on Qwen2.5-
 868 Math-7B, indicating that the structural effects of
 869 LESS are stable across model scales and are not
 870 specific to a single backbone.

871 A.5 Additional Overlap–Accuracy Curves

872 To check whether the correlation between low-
 873 entropy overlap and accuracy holds beyond
 874 Qwen2.5-Math-7B, we repeat the analysis in § 2
 875 on Qwen2.5-Math-1.5B and Qwen2.5-7B. Figure 7

876 shows that, under GRPO training, the overlap ra-
 877 tio of low-entropy segments in correct responses
 878 grows together with accuracy for both backbones.
 879 Early in training, both curves are low and noisy;
 880 as learning proceeds, the overlap becomes higher
 881 and smoother while accuracy also rises. These re-
 882 sults support our claim that low-entropy segment
 883 overlap tracks the formation of stable reasoning
 884 routines across different model sizes and pretraining setups.

885 A.6 Additional Accuracy–Overlap 886 Correlations

887 In § 2, we report that, for Qwen2.5-Math-7B,
 888 benchmark accuracy is strongly correlated with
 889 the overlap of low-entropy segments across cor-
 890 rect responses. To test the robustness of this phe-
 891 nomenon, we repeat the correlation analysis on four
 892 additional backbones: Qwen2.5-7B, DeepSeek-R1-
 893 Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-7B,
 894 and Qwen2.5-Math-1.5B-Oat-Zero.

895 Figure 8 summarizes the results. For each back-
 896 bone, we compute four overlap ratios at the bench-
 897 mark level: (i) overlap among all responses, (ii)
 898 overlap among correct responses only, (iii) over-
 899 lap among segments shared by correct and incor-
 900 rect responses, and (iv) overlap among incorrect

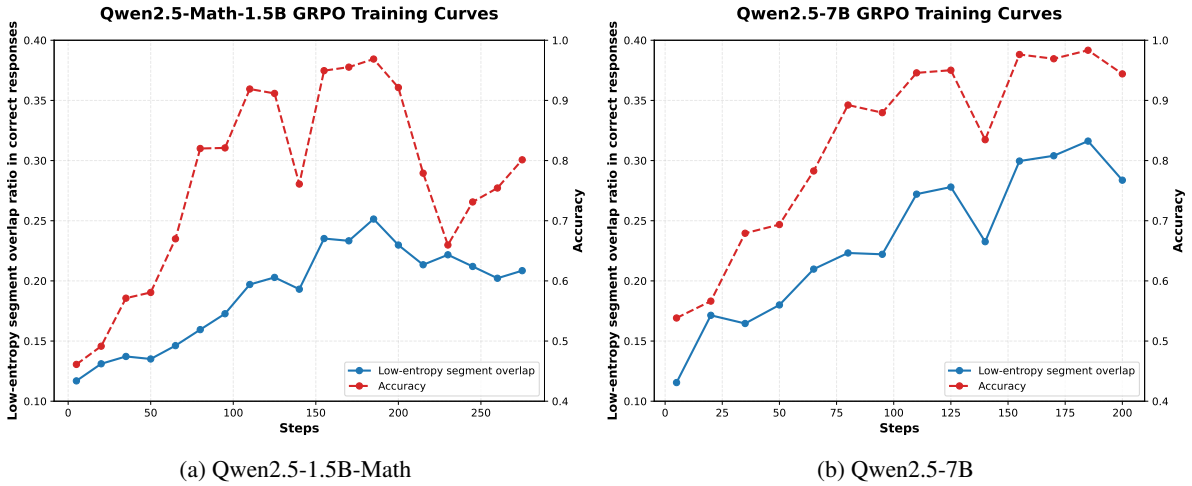
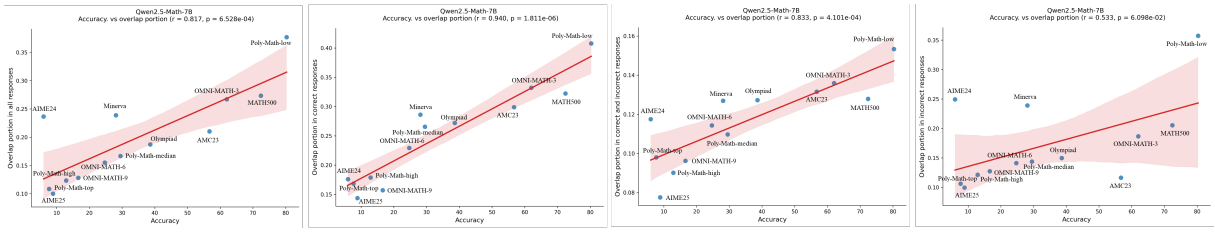


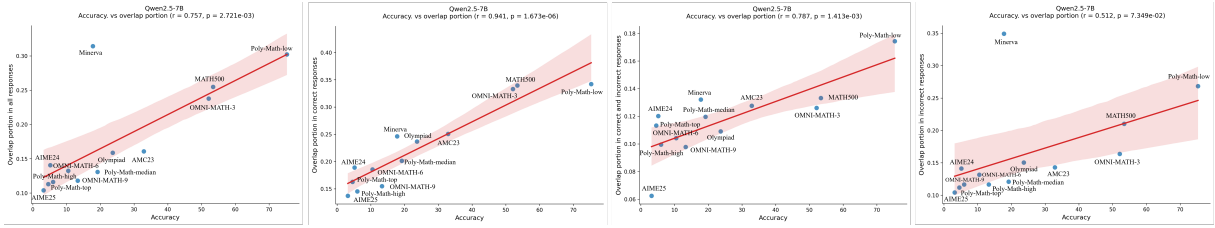
Figure 7: GRPO training curves on two additional backbones. We plot the overlap ratio of low-entropy segments in correct responses (left y-axis, blue) and accuracy (right y-axis, red) over training steps for (a) Qwen2.5-Math-1.5B and (b) Qwen2.5-7B. In both models, low-entropy overlap and accuracy increase in tandem, echoing the trend observed for Qwen2.5-Math-7B in the main text.

901 responses only. We then correlate each ratio with
 902 benchmark accuracy.

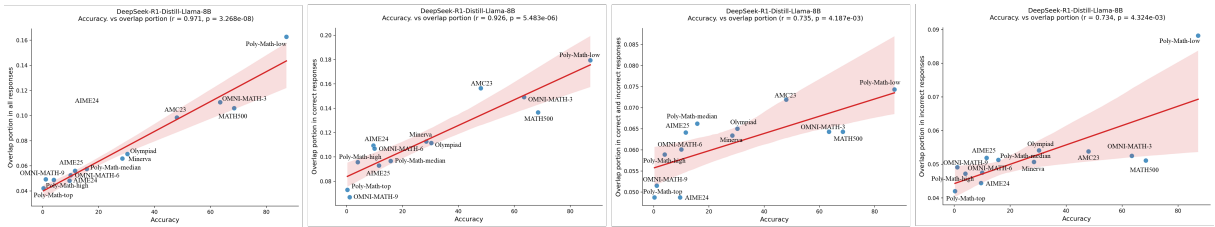
903 Across all five backbones, we observe a con-
 904 sistent pattern: overlap among correct responses
 905 shows the strongest positive correlation with ac-
 906 curacy, overlap among all responses and shared
 907 segments yields weaker but still positive correla-
 908 tions, and overlap among incorrect-only segments
 909 is weakly correlated or even negatively correlated.
 910 These additional results support our claim that sta-
 911 ble low-entropy structure in correct trajectories is a
 912 reliable indicator of reasoning quality, while over-
 913 lap driven by incorrect trajectories does not trans-
 914 late into better performance.



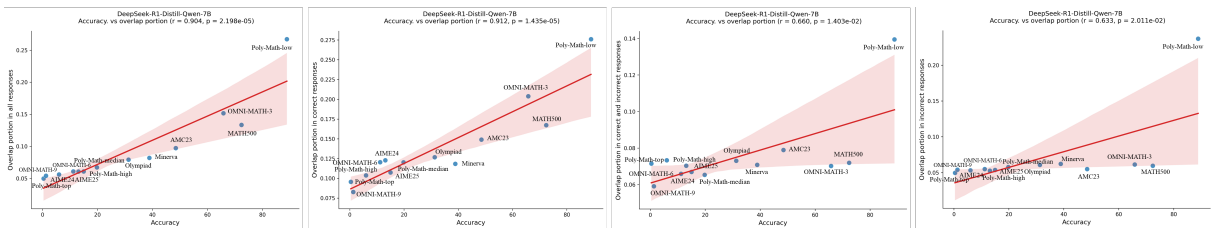
(a) Qwen2.5-Math-7B



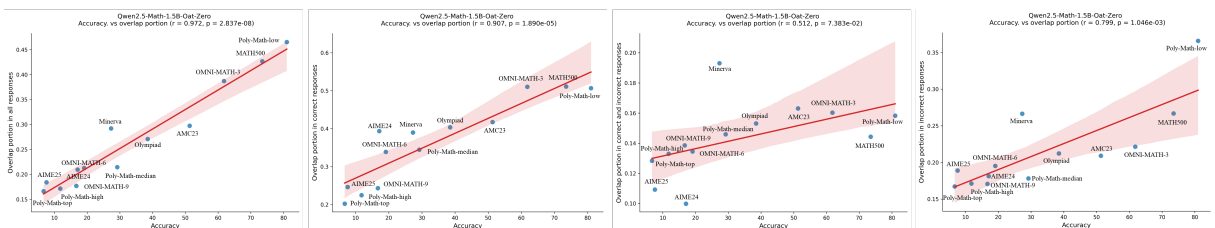
(b) Qwen2.5-7B



(c) DeepSeek-R1-Distill-Llama-8B



(d) DeepSeek-R1-Distill-Qwen-7B



(e) Qwen2.5-Math-1.5B-Oat-Zero

Figure 8: Additional correlations between accuracy and low-entropy segment overlap across backbones. Panels (a)–(e) report, for Qwen2.5-Math-7B, Qwen2.5-7B, DeepSeek-R1-Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-7B, and Qwen2.5-Math-1.5B-Oat-Zero, the Pearson correlations between benchmark accuracy and four overlap ratios: all responses, correct-only responses, segments shared by correct and incorrect responses, and incorrect-only responses. Each point is a benchmark; the red line and shaded area show the fitted regression and its confidence band. Across models, accuracy is most strongly aligned with overlap among correct responses.