# ENHANCING REASONING IN LARGE LANGUAGE MODELS VIA ENTROPY-AWARE SELF-EVOLUTION

**Anonymous authors**Paper under double-blind review

### **ABSTRACT**

Large language models (LLMs) have exhibited remarkable reasoning capabilities. However, when self-evolution frameworks are employed to further enhance these models, a key challenge lies in balancing correctness, which ensures reliable supervision, and exploration, which promotes diverse reasoning trajectories. To address this dilemma, we propose an **entropy-aware self-evolution framework** that integrates verifier feedback with both sequence-level and token-level entropy. Our approach incorporates two key strategies: (i) high-entropy selection of verified trajectories to provide informative yet reliable signals; and (ii) entropy-aware rethinking, which revisits uncertain reasoning steps to uncover alternative solutions. Theoretically, we establish the connection between entropy and the expected supervised fine-tuning loss, showing that high-entropy trajectories yield stronger learning signals. Empirically, experiments across multiple reasoning benchmarks demonstrate that our framework consistently improves both reliability and exploratory capacity over strong baselines. With the assistance of the proposed framework, InternLM2.5-1.8B achieves an improvement of **8.27**% and surpasses the strong baseline by 1.82% on the GSM8K task, as measured by Pass@16. Our results highlight entropy as a principled driver of self-improvement, enabling LLMs to evolve toward models that are not only more accurate but also more exploratory.

### 1 Introduction

Large language models (LLMs) have shown impressive reasoning capabilities across tasks such as mathematical problem solving, code generation, and scientific discovery (OpenAI, 2024; DeepSeek-AI, 2025; Zhu et al., 2025). Despite these successes, traditional training methods often rely on static datasets and may not fully exploit the models' potential for iterative improvement. A growing trend, known as self-evolution, addresses this by generating new training trajectories and fine-tuning models iteratively on them (Wang et al., 2022; Xu et al., 2025; Zhou et al., 2025). While this approach supports scalable iterative self-improvement, it faces a fundamental dilemma: models must balance **correctness** (ensuring generated trajectories are valid and high-quality) with **exploration** (encouraging diverse and novel reasoning paths that might reveal new insights).

Existing approaches to self evolution typically lean towards one side of this trade-off. Verifier-based or reinforcement learning with verifiable rewards (RLVR) methods (Lambert et al., 2025; Shao et al., 2024) prioritize correctness by filtering out invalid trajectories and aligning models with reliable supervision. However, these methods often bias learning toward low-perplexity, deterministic reasoning paths, thereby diminishing exploration and leading to convergent behaviors (Yue et al., 2025). Conversely, exploration-driven strategies based on entropy, perplexity, or trial-and-error sampling (Wang et al., 2025b; Li et al., 2025; Deng et al., 2025) encourage diversity, but correctness is not guaranteed, producing noisy or misleading training signals. Consequently, despite significant progress, current self-evolution frameworks struggle to balance correctness and exploration effectively.

To address the correctness—exploration trade-off, we present an entropy-aware self-evolution framework. Our key insight is that verified high-entropy trajectories not only furnish reliable supervision but also, by leveraging their intrinsic uncertainty, illuminate alternative reasoning paths that warrant exploration. By exploiting entropy at both the sequence and token level, and integrating verifier

feedback, our framework achieves a principled balance between correctness—providing dependable learning signals—and exploration—enabling diverse and informative data generation. Specifically, the framework employs two complementary strategies: (i) **High-Entropy Selection**, which prioritizes trajectories with high uncertainty yet verified correctness to supply both informative and reliable training signals; and (ii) **Entropy-Aware Revisiting of Reasoning Steps**, which identifies high-uncertainty reasoning positions for truncation and regeneration, uncovering alternative solutions and promoting exploratory reasoning. Experiments across different models and tasks demonstrate the superiority of our proposed method, surpassing the strong baseline by **1.44%-5.52%** at average performance on four math reasoning tasks. Our contributions are as follows:

- We propose a novel high-entropy trajectory selection strategy that balances correctness and exploration, addressing a key limitation of prior low-perplexity-biased frameworks.
- We introduce an entropy-aware rethinking mechanism that revisits uncertain reasoning steps, systematically enriching solution diversity while preserving reliability.
- We provide both theoretical analysis, establishing the link between sequence-level entropy
  and expected supervised fine-tuning loss, and extensive empirical validation on reasoning
  benchmarks, demonstrating that our framework consistently improves both reliability and
  exploratory capacity compared to strong baselines.

### 2 RELATED WORK

Self-Evolution with Data Synthesis and Selection. Existing self-evolution approaches for LLMs have explored a variety of strategies for data synthesis and selection. Prior work on data synthesis for self-evolution has relied on heuristic filtering (Wang et al., 2022), confidence-based ranking (Huang et al., 2023), or similarity measures (Chen et al., 2024), while others incorporate external verifiers or interactive environments (Xu et al., 2025; Zhou et al., 2025). Although these strategies improve correctness, they often sacrifice data diversity, leading to convergent trajectories in later training stages. Recent uncertainty-aware approaches leverage entropy (Wang et al., 2025b), perplexity (Li et al., 2025), or exploration-driven sampling (Deng et al., 2025) to encourage diversity, but lack fine-grained utilization of trajectory entropy dynamics. In contrast, our method combines an external verifier with both trajectory-level and token-level entropy guidance, ensuring correctness while systematically enriching diversity and exploration, thus achieving a balanced and robust self-evolution process.

Reinforcement Learning using Verifiable Rewards. With the increasing adoption of reinforcement learning in LLM training, Reinforcement Learning with Verifiable Rewards (RLVR) (Lambert et al., 2025) has emerged as a promising paradigm for enhancing reasoning in LLMs. Similar to our study, RLVR can be viewed as a self-evolution framework that integrates external verifiers. Notably, models such as OpenAI o1(OpenAI, 2024) and DeepSeek-R1(DeepSeek-AI, 2025) exemplify the effectiveness of this approach. In particular, DeepSeek-R1 employs the GRPO (Shao et al., 2024), which eliminates reliance on a reward model and has inspired a range of extensions such as DAPO(Yu et al., 2025) and VAPO(Yue et al., 2025). However, recent analyses indicate several limitations: post-RL models often exhibit reduced exploration compared to their base counterparts(Yue et al., 2025); and correct rewards may still be entangled with erroneous reasoning steps, leading to noisy training signals(Yee et al., 2024; Wan et al., 2025; Wen et al., 2025). Similar to some works on RL with an entropy perspective(Wang et al., 2025a; Cheng et al., 2025), our method leverages entropy-driven self-evolution to preserve exploration ability, operates effectively in domain-specific tasks without requiring long nature language CoTs, and employs a robust external verifier to ensure correctness, thereby avoiding reinforcement of spurious reasoning.

# 3 Метнор

As shown in Figure 1, We propose an entropy-aware self-evolution framework for LLMs, composed of three stages: (1) **Trajectory Exploration** — generating candidate reasoning trajectories to probe the task space, (2) **Trajectory Rethinking** — revisiting uncertain reasoning steps to diversify problem-solving paths, and (3) **Trajectory Selection** — curating informative trajectories to enhance both training signal and model exploration ability.

The central advantage of this design lies in its explicit focus on high-entropy samples, which are indicative of epistemic uncertainty and exploratory potential. By prioritizing such samples and leveraging verifier feedback, our framework not only improves data quality but also systematically encourages the model to explore alternative reasoning paths. The pipeline is iterated for I steps, starting with a base model  $\pi_0$  at iteration i=0.

### 3.1 Entropy Measures for Model Trajectories.

We quantify uncertainty in model-generated trajectories using token-level and sequence-level entropy.

**Local uncertainty**: We utilize the token-level entropy to capture local uncertainty and inform high-entropy truncation and revisiting during trajectory refinement. Formally, the token-level entropy at position t is defined as

$$H_t = -\sum_{i=1}^{V} p_{\theta}(v_i|\boldsymbol{y}_{< t}, \boldsymbol{x}) \log p_{\theta}(v_i|\boldsymbol{y}_{< t}, \boldsymbol{x}),$$
(1)

where  $p_{\theta}(v_i|\mathbf{y}_{< t}, \mathbf{x})$  is the model's predictive probability for token  $v_i$  given prefix  $\mathbf{y}_{< t}$  and input  $\mathbf{x}$ . A low  $H_t$  indicates that the model's predictions are concentrated on a small set of tokens, reflecting high confidence, while high  $H_t$  reflects multiple plausible alternatives, creating branching points that can decisively influence the trajectory.

**Global uncertainty**: We utilize the sequence-level entropy that aggregates token-level uncertainties to measure global unpredictability of a trajectory  $y = (y_1, \dots, y_T)$ :

$$H_{\text{seq}}(\boldsymbol{y} \mid \boldsymbol{x}) = \frac{1}{T} \sum_{t=1}^{T} H_t.$$
 (2)

Trajectories with high  $H_{\rm seq}$  contain multiple positions with substantial uncertainty, indicating both higher exploratory potential and richer information content. Conversely, low  $H_{\rm seq}$  trajectories correspond to more deterministic generations. Sequence-level entropy thus provides an effective criterion for selecting uncertainty and exploratory trajectories in supervised fine-tuning (SFT).

In out framework, token-level entropy identifies critical positions for trajectory refinement, while sequence-level entropy selects high-information trajectories for SFT. By leveraging both, the model benefits from trajectories that are both exploratory and informative, thereby enhancing the task-specific performance of LLMs.

### 3.2 Trajectory Exploration

We start by broadly exploring the solution space, allowing the model to generate candidate trajectories while quantifying their uncertainty. Let  $\mathcal{D}$  denote a task-specific dataset comprising instruction-answer pairs  $(\boldsymbol{x},a)$ . At iteration i, the current model  $\pi_i$  generates K trajectories for each input  $\boldsymbol{x}$ :  $\{\boldsymbol{y}_k\}_{k=1}^K \sim \pi_i(\cdot \mid \boldsymbol{x})$ . For each trajectory  $\boldsymbol{y}_k$ , we compute its sequence-level entropy:  $h_k = H_{\text{seq}}(\boldsymbol{y}_k \mid \boldsymbol{x})$ . Each trajectory is then verified by an external checker (Xu et al., 2025), yielding a correctness label:  $r_k = \text{validator}(\boldsymbol{y}_k, a), r_k \in \{0, 1\}$ . The final quadruple is stored as  $T_k = (\boldsymbol{x}, \boldsymbol{y}_k, h_k, r_k)$ . All positively verified trajectories are aggregated into the *exploration pool*:

$$\mathcal{P}_{i}^{+} = \{ T_{k} \mid r_{k} = 1 \}_{k=1}^{K} \cup \mathcal{P}_{i-1}^{+}, \quad \mathcal{P}_{-1}^{+} = \varnothing.$$
 (3)

This pool serves as the foundation for subsequent trajectory selection.

### 3.3 TRAJECTORY RETHINKING

Prior work (Wang et al., 2025c; Gao et al., 2025) emphasizes that medium-difficulty and uncertain samples play a crucial role in self-training. To better exploit such informative cases, we introduce *trajectory rethinking*, which revisits high-entropy reasoning steps to encourage exploration of alternative solutions.

From the verified trajectories of this iteration  $\{T_k \mid r_k=1\}_{k=1}^K$ , we select the positive trajectory with the highest sequence-level entropy:  $\boldsymbol{y}^\star = \arg\max_{\boldsymbol{y}_k \in \mathcal{P}_i^+} H_{\text{seq}}(\boldsymbol{y}_k \mid \boldsymbol{x})$ . Let T be the length

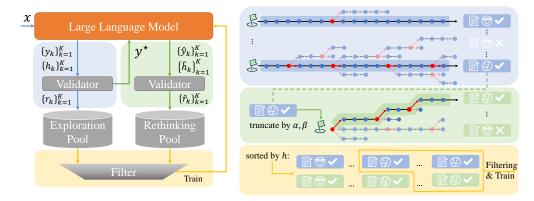


Figure 1: (**Left**) Pipeline shows our entropy-aware self-evolution framework. (**Right**) Three stages for the framework. Three background colors in the left—blue, green, and yellow—indicate the same stages as those in the right from top to bottom. The trajectory exploration stage, highlighted in blue, illustrates how the model explores and verifies candidate trajectories, as detailed in Section 3.2; The trajectory rethinking stage, highlighted in green, illustrates how we leverage the explored correct trajectories to truncate and regenerate, as detailed in Section 3.3. The trajectory selection stage, highlighted in yellow, selects highly exploratory and informative trajectories to enhance the model's capabilities, as detailed in Section 3.3. Through repeated iterations of this framework, we construct a set of trajectories that are both reliable and exploratory, which facilitates the enhancement of the model's task execution and exploratory capabilities. The three stages progressively transform raw trajectories into reliable yet diverse supervision signals.

of  $y^*$ . Token-level entropies  $H_t$  are used to identify uncertain positions. With hyperparameters  $\alpha \in (0,1)$  (fraction of top-entropy tokens) and  $\beta \in (0,1)$  (maximum truncation ratio), we define the candidate set:

$$\mathcal{I} = \{ t \mid t \le |\beta T|, \ y_t^* \in \text{Top}_{\alpha}(H_t) \}. \tag{4}$$

We then sample a truncation point:  $\tau \sim \operatorname{Uniform}(\mathcal{I})$ , and obtain the truncated prefix:  $\mathbf{y}_{\leq \tau}^{\star} = (y_1^{\star}, \dots, y_{\tau}^{\star})$ . Conditioned on  $(\mathbf{x}, \mathbf{y}_{\leq \tau}^{\star})$ , the model generates K continuations:  $\{\tilde{\mathbf{y}}_{k, > \tau}\}_{k=1}^{K} \sim \pi_i(\cdot \mid \mathbf{x}, \mathbf{y}_{\leq \tau}^{\star})$ , which are concatenated with the prefix to form *rethought trajectories*:  $\{\tilde{\mathbf{y}}_k\}_{k=1}^{K} = \{\mathbf{y}_{\leq \tau}^{\star} \oplus \tilde{\mathbf{y}}_{k, > \tau}\}_{k=1}^{K}$ . All rethought trajectories are verified, and positives are aggregated into the *rethinking pool*:

$$\tilde{\mathcal{P}}_{i}^{+} = \{ \tilde{T}_{k} = (\boldsymbol{x}, \tilde{\boldsymbol{y}}_{k}, \tilde{h}_{k}, \tilde{r}_{k}) \mid \tilde{r}_{k} = 1 \}_{k=1}^{K} \cup \tilde{\mathcal{P}}_{i-1}^{+}, \quad \tilde{\mathcal{P}}_{-1}^{+} = \varnothing.$$
 (5)

When no positively verified samples exist, we apply the procedure to the negative trajectory with the highest sequence-level entropy, so that high-entropy trajectories, regardless of their correctness, continue to drive exploration of alternative reasoning paths.

### 3.4 Trajectory Selection

During the self-evolution process, the contributions of different generated trajectories to model learning vary significantly. To maximize the utility of limited training resources, it is necessary to select trajectories that are both exploratory and information-rich from a large pool of candidates. The trajectory selection stage aims to aggregate and identify these critical trajectories to enhance the model's learning. By emphasizing high-entropy trajectories, this selection process encourages the model to explore uncertain regions of the solution space, thereby acquiring a more comprehensive reasoning experience.

Specifically, we rank both  $\mathcal{P}_i^+$  and  $\tilde{\mathcal{P}}_i^+$  in descending order of sequence-level entropy, obtaining  $\mathcal{R}_i^+$  and  $\tilde{\mathcal{R}}_i^+$ . From these, we select the top-N trajectories from the exploration pool:

$$\mathcal{T}_1 = \left\{ (\boldsymbol{x}, y_n) \mid n \le \min\left(N, |\mathcal{R}_i^+|\right), \ T_n \in \mathcal{R}_i^+ \right\}. \tag{6}$$

If  $|\mathcal{T}_1| < N$ , we fill the remainder from the rethinking pool:

$$\mathcal{T}_2 = \left\{ (\boldsymbol{x}, \tilde{\boldsymbol{y}}_n) \mid n \le \min\left(N - |\mathcal{T}_1|, |\tilde{\mathcal{R}}_i^+|\right), \ \tilde{T}_n \in \tilde{\mathcal{R}}_i^+ \right\}. \tag{7}$$

Supervised fine-tuning on the filtering trajectories. We fine-tune the model  $\pi_i$  on  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$  using maximum likelihood estimation (MLE) also known as the cross-entropy loss  $\mathcal{L}_{CE}$  to get nextiteration model  $\pi_{i+1}$ ,

$$\mathcal{L}_{CE} = -\sum_{(\boldsymbol{x}, \boldsymbol{y}) \sim \mathcal{T}_1 \cup \mathcal{T}_2} \log p_{\theta}(\boldsymbol{y} \mid \boldsymbol{x}). \tag{8}$$

# 3.5 Analysis of the Relationship Between Entropy and the Expected Supervised Loss

The defination of cross-entropy loss for SFT on one self-generated trajectory y is

$$\mathcal{L}_{CE}(\boldsymbol{y}|\boldsymbol{x}) = -\sum_{t=1}^{T} \log p_{\theta}(y_t \mid \boldsymbol{y}_{< t}, \boldsymbol{x}). \tag{9}$$

Its expectation over trajectories sampled from the model  $\pi_{\theta}(\cdot|\boldsymbol{x})$  can be expressed as

$$\mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\cdot | \boldsymbol{x})} [\mathcal{L}_{CE}(\boldsymbol{y} | \boldsymbol{x})] = -\sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\cdot | \boldsymbol{x})} [\log p_{\theta}(y_t \mid \boldsymbol{y}_{< t}, \boldsymbol{x})]$$
(10)

$$= \sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{y}_{< t} \sim \pi_{\theta}(\cdot | \boldsymbol{x})}[H_t]$$
 (11)

$$= T \cdot \mathbb{E}_{\boldsymbol{y} \sim \pi_{\theta}(\cdot \mid \boldsymbol{x})} [H_{\text{seq}}(\boldsymbol{y} \mid \boldsymbol{x})], \tag{12}$$

where the second equality follows from the definition of token-level entropy and the last equality from sequence-level entropy. This relationship shows that higher-entropy trajectories induce larger expected loss, producing stronger gradients and richer learning signals.

Overall, our method combines verifier guidance with entropy-aware trajectory selection. By explicitly exploiting high-entropy samples for both exploration and augmentation, the framework not only ensures training quality but also enhances the model's ability to explore and generalize across uncertain reasoning pathways. Through iterative self-evolution, the model progressively improves its task-specific reasoning performance.

# 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Datasets.** We evaluate the proposed framework on math reasoning tasks, using a Python executor as the validator. Reasoning tasks include: GSM8K(Cobbe et al., 2021), MATH(Hendrycks et al., 2021), GSM-Hard(Gao et al., 2023), SVAMP(Patel et al., 2021), and AsDiv(Miao et al., 2020). The training split of GSM8K, along with randomly selected samples from MATH, is used to construct the dataset with 13,492 samples for self-evolution. The test splits of GSM8K, GSM-Hard, SVAMP, and AsDiv are reserved for evaluation. In order to make use of the validator, we prompt the LLM to generate reasoning path with the format of executable python code.

**Training Details.** We use Qwen2.5-Instruct(Yang et al., 2024; Qwen, 2024), Llama3.2(Grattafiori et al., 2024; Meta, 2024) and InternLM-2.5(Cai et al., 2024) models for evaluation. At the first iteration, we utilize few-shot prompting to instruct the model to generate training samples as a cold start. The few-shot numbers for math reasoning tasks are set to 3. At each evolution iteration, the candidate trajectory size K is set to 5. The total iteration number I is set to 10 for InternLM2.5-1.8B, 7 for Llama3.2-1B and 7 for Qwen2.5-Instruct-1.5B. The top-N for trajectory augmentation is set to 10. Otherwise, we make use of the negative trajectories the same as the baseline (Xu et al., 2025). All the self-evolution experiments are implemented on  $4 \times RTX3090$  of 24GB VRAM.

### 4.2 MAIN RESULTS

Table 1 summarizes the evaluation results across four mathematical reasoning benchmarks. For reference, we include a few-shot baseline, while all other evaluations are conducted under the zero-

shot setting. To ensure fairness, all experiments adopt a consistent sampling strategy with top- p=0.95 and temperature =0.6. We further compare our approach with the Envisions framework (Xu et al., 2025) under identical conditions. To evaluate both accuracy and exploratory capacity, we use Pass@K as the primary metric, as it reflects the model's ability to produce correct solutions under multiple sampled attempts.

Overall Performance Improvements. Our method delivers substantial improvements over the base models and consistently outperforms Envisions, as shown in Tabel 1. On the held-in task GSM8K, InternLM2.5-1.8B achieves a remarkable 8.27% gain at Pass@16. Compared with Envisions, our method yields improvements of 1.82% and 4.39% at Pass@16 and Pass@128, respectively, along with an average performance gain of 2.57% when K ranges from 16 to 256. These results indicate that our approach not only strengthens task execution accuracy relative to the base models, but also enhances exploratory capacity when compared to existing frameworks.

Generalization to Held-out Benchmarks. To examine generalization, we conduct evaluations on GSM-Hard, AsDiv, and SVAMP (Table 1). Consistent with the observations on GSM8K, our method achieves clear gains over the base models and surpasses ENVISIONS on GSM-Hard and AsDiv. On GSM-Hard, InternLM2.5-1.8B improves by 7.21% and delivers an additional 1.44% average gain compared with ENVISIONS. On SVAMP and AsDiv, our method outperforms the baseline by 5.52% and 5.51% in average performance, respectively. These results demonstrate the strong generalization ability of our framework across diverse reasoning benchmarks. Moreover, on SVAMP, which is a relatively simple benchmark, InternLM2.5-1.8B already matches or exceeds the performance of self-evolution variants under few-shot settings. In contrast, our method better preserves the exploratory capacity of the base models, whereas ENVISIONS exhibits a noticeable decline.

**Generalization to Various Backbones.** We also compare our method with ENVISIONS on Llama 3.2-1B and Qwen 2.5-Instruct-1.5B. As shown in Figure 2, our method consistently outperforms ENVISIONS across tasks and backbones. Significantly, as illustrated in Figure 3, the performance improvements become more pronounced at larger K, highlighting that our evolutionary strategy effectively enhances the ability of models to explore diverse solution trajectories.

Table 1: Math Reasoning results of InternLM2.5-1.8B on four tasks.

	GSM8K			GSM-Hard			SVAMP			AsDiv		
	Pass@16	Pass@128	Avg	Pass@16	Pass@256	Avg	Pass@16	Pass@256	Avg	Pass@16	Pass@128	Avg
InternLM2.5-1.8	8 <i>B</i>											
Few-shot	63.53	84.00	73.73	52.84	74.68	60.93	84.30	95.70	89.52	76.01	84.68	80.00
Envisions	69.98	80.67	75.07	59.36	71.19	64.20	79.50	88.20	83.01	72.97	78.44	75.68
Ours	71.80	85.06	77.64	60.05	75.21	65.64	83.90	95.10	88.53	77.61	85.42	81.19
Δ	+1.82	+4.39	+2.57	+0.68	+4.02	+1.44	+4.40	+6.90	+5.52	+4.64	+6.98	+5.51

### 4.3 EVOLUTION PROGRESS FOR SELF-EVOLUTION FRAMEWORKS

As illustrated in Figure 4(**Left**), the iterative evolution curves of the self-training frameworks with InternLM2.5-1.8B as the LLM, demonstrate the progression of performance improvement. Compared with the ENVISIONS method, our framework exhibits a more pronounced performance improvement. Notably, while the performance of ENVISIONS tends to plateau after the fourth iteration, our method not only achieves superior results but also shows continued potential for further improvement. From Figure 4 (**Right**), it can be observed that under our framework, both the mean and variance of sequence-level entropy in the training dataset increase as the number of self-evolution iterations grows, exhibiting a trend in sharp contrast to that of the ENVISIONS method.

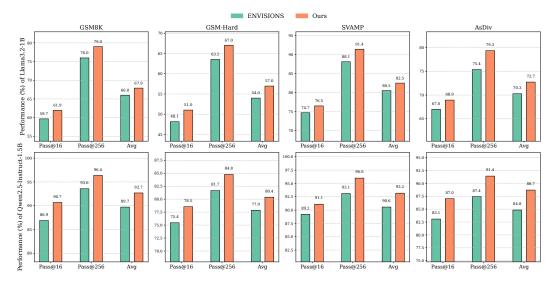


Figure 2: Math Reasoning evaluation of the Llama3.2-1B and Qwen2.5-Instruct-1.5B on the four tasks, compared with the existing method.

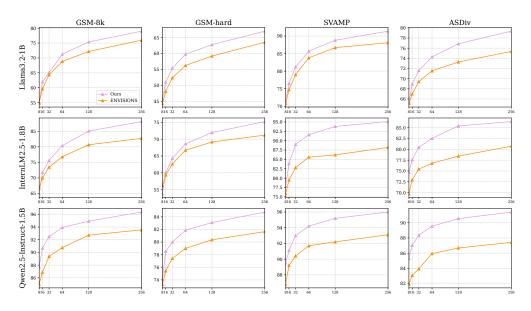


Figure 3: Pass@K performance of the LLMs with different self-evolution frameworks. The horizontal axis denotes K ranging from 8 to 256, and the vertical axis shows the corresponding Pass@K accuracy on the benchmarks.

### 5 ANALYSIS

# 5.1 HIGH-ENTROPY SELECTION ENHANCES TRAINING INFORMATION AND TRAJECTORY DIVERSITY

To further investigate the effect of out high-entropy selection strategy, we analyze the distribution of similarity scores and negative log probability of the selected trajectories for the last self-evolution iteration of three models.

The similarity score quantifies the alignment among generated trajectories, with higher values indicating greater overlap and lower values reflecting higher diversity. Formally, given a set of n trajectories  $(t_1, t_2, \ldots, t_n)$  corresponding to the same problem, we obtain their embeddings  $\{e_i\}_{i=1}^n$ 

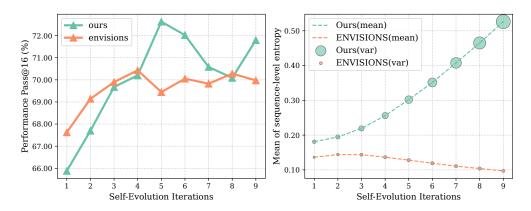


Figure 4: (**Left**) Performance evolution of two frameworks on InternLM-2.5-1.8B model. (**Right**) Mean and variance of sequence-level entropy of the SFT training datas for each evolution.

from a pretrained embedding model  $f(\cdot)$  (Zhang et al., 2025). The similarity score is computed as

$$\operatorname{Sim} = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{\substack{j=1 \\ j \neq i}}^{n} \langle f(t_i^q), f(t_j^d) \rangle$$

where  $f(t_i^q)$  and  $f(t_j^c)$  denote query-style and candidate-style embeddings of trajectory t, and  $\langle \cdot, \cdot \rangle$  denotes the inner product. See Appendix C for more details.

As shown in the top row of Figure 5, our method produces a wider distribution of similarity scores with a noticeable shift toward lower values compared to ENVISIONS, indicating that high-entropy selection promotes greater trajectory diversity. Meanwhile, the bottom row reveals that our approach selects trajectories with higher negative log probabilities, implying that the chosen samples carry more informative signals rather than being restricted to high-confidence outputs. Overall, these results demonstrate that high-entropy selection enhances both the information content and the diversity of the training data, which are crucial for improving the expertise and generalization capability of LLMs in self-evolution frameworks.

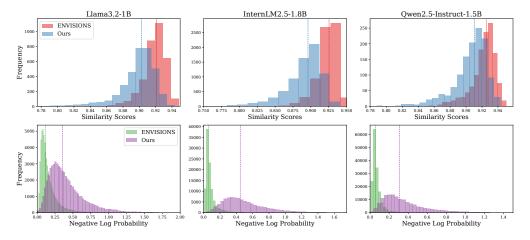


Figure 5: Histogram of Similarity Scores and Negative Log Probability of the trajectories selected for the last self-evolution iteration. The dashed lines in the figures denote the median.

#### 5.2 The role of Trajectory Rethink in Self-Evolution.

To analyze the role of the *Trajectory Rethink* stage within our framework, we conduct an in-depth investigation from three perspectives. First, we evaluate its impact on reasoning performance. Specifically, we compare the Qwen2.5-Instruct-1.5B model on the GSM-Hard, a relatively challenging

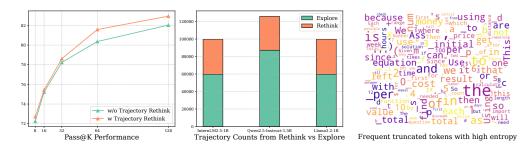


Figure 6: Analysis of Trajectory Rethink in self-evolution. (Left) Pass@K Performance: Incorporating rethink consistently improves performance across different values of K. (Middle) Trajectory Counts: Rethink and explore complement each other across different base models, leading to an increase in effective training samples. (Right) High-Entropy Tokens: The frequent occurrence of truncated tokens with high entropy indicates that rethink mitigates uncertainty and enhances trajectory diversity.

task, with and without the Trajectory Rethink strategy under small-batch training. As shown in Figure 6 (**Left**), incorporating Trajectory Rethink consistently improves pass@K accuracy, demonstrating its effectiveness. In contrast, the model without this stage relies solely on the self-refine strategy achieves inferior performance on exploration.

Moreover, we examine the contribution of Trajectory Rethink to trajectory diversity. Figure 6 (**Middle**) shows that this strategy accounts for more than one-third of the training trajectories generated during the evolution process, substantially enriching the diversity of the training data. This indicates that rethink contributes significantly to the breadth of explored reasoning paths.

Finally, we analyze the linguistic patterns associated with rethink. We visualize the most frequent truncated tokens with high entropy, as shown in Figure 6 (**Right**). Words such as "because", "since", and "then" often determine the direction of reasoning. Truncating trajectories at these critical tokens enables the model to rethink from pivotal decision forks, thereby facilitating more flexible and diverse reasoning. These analyses demonstrate that Trajectory Rethink is a crucial component of our self-evolution framework. It enhances the diversity of reasoning trajectories and encourages re-exploration from meaningful reasoning pivots, ultimately leading to richer and more informative training signals, particularly beneficial for challenging reasoning tasks.

### 6 CONCLUSION

We propose an entropy-aware self-evolution framework that enhances reasoning in large language models by strategically leveraging uncertainty to balance correctness and exploration. Integrating verifier feedback with sequence-level and token-level entropy, our method prioritizes high-entropy yet verified trajectories for training, ensuring reliable supervision while actively promoting diverse reasoning paths. Theoretical analysis shows that such trajectories yield stronger learning signals due to their higher expected loss, enabling more effective fine-tuning. Empirically, our approach achieves significant gains across multiple reasoning benchmarks. Notably, InternLM2.5-1.8B improves by **8.27%** on GSM8K at Pass@16 and surpasses the strong Envisions baseline by **4.39%** at Pass@128, with consistent gains on held-out tasks like GSM-Hard, SVAMP and AsDiv. Critically, performance improvements grow with larger sampling budgets, confirming enhanced exploration without sacrificing accuracy.

**Limitation** Our experiments are limited to models up to 1.8B parameters due to computational constraints; scaling to larger architectures (e.g., 7B+) remains untested. The framework's reliance on executable verifiers also restricts current applicability to math/code domains. Future work will address efficiency, entropy approximation, and extension to semantic reasoning tasks.

In summary, our entropy-aware self-evolution framework offers a principled, theoretically grounded, and empirically validated approach to enhancing both the reliability and exploratory capacity of LLMs. By treating uncertainty not as noise to be suppressed but as signal to be harnessed, we enable models to evolve into more capable, flexible, and robust reasoners.

### REFERENCES

Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. Internlm2 technical report, 2024.

Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), Proceedings of the 41st International Conference on Machine Learning, volume 235 of Proceedings of Machine Learning Research, pp. 6621–6642. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/chen24j.html.

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. Reasoning with exploration: An entropy perspective on reinforcement learning for llms, 2025. URL https://arxiv.org/abs/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. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.

DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

Jia Deng, Jie Chen, Zhipeng Chen, Daixuan Cheng, Fei Bai, Beichen Zhang, Yinqian Min, Yanzipeng Gao, Wayne Xin Zhao, and Ji-Rong Wen. From trial-and-error to improvement: A systematic analysis of llm exploration mechanisms in rlvr, 2025. URL https://arxiv.org/abs/2508.07534.

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. PAL: Program-aided language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 10764–10799. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/gao23f.html.

Zitian Gao, Lynx Chen, Haoming Luo, Joey Zhou, and Bryan Dai. One-shot entropy minimization, 2025. URL https://arxiv.org/abs/2505.20282.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco

541

542

543

544

546

547

548

549

550

551

552

553

554

558

559

561

562

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

582

583

584

585

586

588

592

Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L.

Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

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

Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1051–1068, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.67. URL https://aclanthology.org/2023.emnlp-main.67/.

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL https://arxiv.org/abs/2411.15124.

Haochen Li, Wanjin Feng, Xin Zhou, and Zhiqi Shen. GiFT: Gibbs fine-tuning for code generation. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pp. 12271–12284, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.599. URL https://aclanthology.org/2025.acl-long.599/.

Meta. LLaMA 3.2 model card. https://huggingface.co/meta-llama/Llama-3.2-1B, 2024. Accessed: 2025-09-10.

Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing english math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 975–984, 2020.

OpenAI. Learning to reason with llms. https://openai.com/index/learning-to-reason-with-llms/, 2024. [Accessed: 2025-05-01].

- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. naacl-main.168. URL https://aclanthology.org/2021.naacl-main.168.
  - Team Qwen. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
  - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.
  - Yue Wan, Xiaowei Jia, and Xiang Lorraine Li. Unveiling confirmation bias in chain-of-thought reasoning, 2025. URL https://arxiv.org/abs/2506.12301.
  - 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. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning, 2025a. URL https://arxiv.org/abs/2506.01939.
  - Xiaoxuan Wang, Yihe Deng, Mingyu Derek Ma, and Wei Wang. Entropy-based adaptive weighting for self-training, 2025b. URL https://arxiv.org/abs/2503.23913.
  - Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Lucas Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, Weizhu Chen, Shuohang Wang, Simon Shaolei Du, and Yelong Shen. Reinforcement learning for reasoning in large language models with one training example. *arXiv* preprint arXiv:2504.20571, 2025c.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 2022.
- Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Sam Bowman, He He, and Shi Feng. Language Models Learn to Mislead Humans via RLHF. In Y. Yue, A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *International Conference on Representation Learning*, volume 2025, pp. 74670–74692, 2025. URL https://proceedings.iclr.cc/paper\_files/paper/2025/file/b9a5a60573637f329b04dlbeda4cd404-Paper-Conference.pdf.
- Fangzhi Xu, Qiushi Sun, Kanzhi Cheng, Jun Liu, Yu Qiao, and Zhiyong Wu. Interactive evolution: A neural-symbolic self-training framework for large language models. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12975–12993, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.635. URL https://aclanthology.org/2025.acl-long.635/.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- Evelyn Yee, Alice Li, Chenyu Tang, Yeon Ho Jung, Ramamohan Paturi, and Leon Bergen. Dissociation of faithful and unfaithful reasoning in llms, 2024. URL https://arxiv.org/abs/2405.15092.

Qiying 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, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025. URL https://arxiv.org/abs/2503.14476.

Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?, 2025. URL https://arxiv.org/abs/2504.13837.

Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025.

Yifei Zhou, Sergey Levine, Jason Weston, Xian Li, and Sainbayar Sukhbaatar. Self-challenging language model agents, 2025. URL https://arxiv.org/abs/2506.01716.

Yao Zhu, Yunjian Zhang, Zizhe Wang, Xiu Yan, Peng Sun, and Xiangyang Ji. Patchwise cooperative game-based interpretability method for large vision-language models. *Transactions of the Association for Computational Linguistics*, 13:744–759, 2025.

### LLM USAGE

 We used large language models (LLMs) as auxiliary tools for writing assistance and language polishing. Specifically, LLMs were employed to improve readability, grammar, and presentation of the text. All research ideas, experimental designs, and scientific contributions are entirely the work of the authors. The authors take full responsibility for the content of this paper.

### A TRAINING DETAILS

The SFT training in our framework and baselines is conducted on  $4 \times RTX3090$  with a maximum length of 2,048. They are optimized and accelerated with Deepspeed Zero3 and FlashAttention2. We use the AdamW optimizer with a *Linear* learning rate of 2e-5. The training epoch is set to 1.

**Prompt Examples.** To guide the model towards generating executable Python code, we prepend the following prompt before each input:

```
Write Python code to solve the question.
```

We illustrate the few-shot prompts used in our experiments. The following shows the training-time few-shot prompt (MATH\_PROMPT\_FS) and the test-time prompt (MATH\_PROMPT\_FS\_TEST). The test-time prompt only contains the first example of training-time prompt.

Listing 1: Few-shot prompt for training (MATH\_PROMPT\_FS)

```
746
      The following are three examples for reference.
747
748
      Example 1:
749
      The question is: Olivia has $23. She bought five bagels for $3 each.
750
      How much money does she have left?
      The solution code is:
751
       '''python
752
      def solution():
753
           '''Olivia has $23. She bought five bagels for $3 each.
754
          How much money does she have left?'''
          money_initial = 23
          bagels = 5
```

```
756
    bagel_cost = 3
    money_spent = bagels * bagel_cost
758
    money_left = money_initial - money_spent
759
    result = money_left
760
    return result
761
    ... (Examples 2 and 3 omitted for brevity)
```

## B TEST TASKS AND BENCHMARK

Table 2 lists the benchmark tasks used in our experiments. Below we provide more detailed descriptions of each dataset: the types of math problems included, what makes them hard or easy, and an example from each.

### **B.1** Dataset Descriptions

- **GSM8K** (**Grade School Math 8K**) (Cobbe et al., 2021) This dataset contains approximately 8,500 linguistically diverse grade-school level word problems. Problems require between 2 to 8 reasoning steps and use basic arithmetic operations (addition, subtraction, multiplication, division). The problems are designed to be solvable without advanced mathematics, but test multi-step reasoning and managing intermediate fractional or decimal computations.
- **GSM-Hard** (Gao et al., 2023) A held-out or more challenging subset related to GSM8K, designed to test generalization under harder or out-of-distribution settings. It shares the same format but contains examples that are less similar to the training distribution.
- **SVAMP** (Patel et al., 2021) Consists of 1,000 math word problems constructed by applying perturbations to existing datasets (such as ASDiv), adding irrelevant information or changing problem structure to challenge robustness. Each problem typically has one unknown variable, with no more than two mathematical expressions.
- ASDiv (Miao et al., 2020) Contains 2,305 word problems spanning a variety of types, with greater lexical variety, more diverse wording, variable placements, and reasoning patterns.
   Problems vary from relatively simple to fairly complex, testing both arithmetic and reasoning about relationships.

### B.2 Example Instances

To illustrate the characteristics of different datasets, we present representative examples as follows:

### • GSM8K

Q: Janet's ducks lay 16 eggs per day. She eats 3 for breakfast and bakes with 4. She sells the remainder at the market for \$2 per egg.

A: 18

# • GSM-Hard

Q: A robe takes 2,287,720 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

*A*: 3,431,580

### SVAMP

Q: There are 87 oranges and 290 bananas. If the bananas are organized into 2 groups, how big is each group of bananas?

A: 145

#### ASDiv

Q: Seven red apples and two green apples are in the basket. How many apples are in the basket?

A: 9

Table 2: Benchmark tasks used in our experiments.

Domains	Task name	Is Held-out?	Test Samples	Max Length	Sources
	GSM8K		1,319	2,048	Cobbe et al. (2021)
Math Reasoning	GSM-Hard	$\checkmark$	1,319	2,048	Gao et al. (2023)
Main Reasoning	SVAMP	$\checkmark$	1,000	2,048	Patel et al. (2021)
	AsDiv	$\checkmark$	2,305	2,048	Miao et al. (2020)

### C COMPUTATION OF SIMILARITY SCORES

To evaluate the diversity of reasoning trajectories, we define a similarity score based on trajectory embeddings.

**Setup.** For each problem instance with at least 10 trajectories, we align datasets by intersecting their origin\_id sets. Each trajectory is embedded using Qwen/Qwen3-Embedding-0.6B, as  $f(\cdot)$ . Queries  $t^q$  are prefixed with a short instruction describing the task of retrieving logically equivalent trajectories, while candidate trajectories  $t^d$  are encoded directly. The instruction for retrieving query is:

task = 'Given a reasoning trajectory in code form, identify and retrieve
 those strictly similar in logic and structure'
return f'Instruction: {task}\nThe given trajectory: {query}'

This instruction guides the model to focus on logical and structural consistency rather than surface-level textual overlap

**Pairwise Similarity.** Let  $E \in \mathbb{R}^{n \times d}$  denote the embeddings of n trajectories. We compute the cosine similarity matrix

$$S = E \cdot E^{\top}$$
.

Self-similarities on the diagonal are masked out. The similarity score for an instance is then

$$\operatorname{Sim}_{\operatorname{instance}} = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{\substack{j=1\\ i \neq i}}^{n} \langle e_i, e_j \rangle,$$

where  $\langle e_i^q, e_j^d \rangle$  denotes cosine similarity between embeddings  $e_i^q = f(t_i^q)$  and  $e_j^d = f(t_j^d)$ .

**Dataset-Level Score.** The dataset-level similarity is the mean over all valid instances:

$$\operatorname{Sim}_{\operatorname{dataset}} = \frac{1}{|\mathcal{D}|} \sum_{k \in \mathcal{D}} \operatorname{Sim}_{\operatorname{instance}}^{(k)}.$$

**Visualization.** We plot histograms of similarity scores across datasets and mark the median with dashed lines, enabling analysis of both central tendency and diversity, as shown in Figure 5. Lower similarity reflects richer trajectory diversity, while higher similarity indicates redundancy.