

EvoCoT: Overcoming the Exploration Bottleneck in Reinforcement Learning for LLMs

Anonymous ACL submission

Abstract

Reinforcement learning with verifiable reward (RLVR) has become a promising paradigm for post-training large language models (LLMs) to improve their reasoning capability. However, when the rollout accuracy is low on hard problems, the reward becomes sparse, limiting learning efficiency and causing exploration bottlenecks. Existing approaches either rely on teacher models for distillation or filter out difficult problems, limiting scalability or restricting reasoning improvement through exploration.

We propose **EvoCoT**, a self-**E**volving curriculum learning framework based on two-stage **C**hain-**o**f-**T**hought (CoT) reasoning optimization. **EvoCoT** constrains the exploration space by self-generating and verifying CoT trajectories, then gradually shortens CoT steps to expand the space in a controlled way. The framework enables LLMs to stably learn from initially unsolved hard problems under sparse rewards. We apply **EvoCoT** to multiple LLM families, including Qwen, DeepSeek, and Llama. Experiments show that **EvoCoT** enables LLMs to solve previously unsolved problems, improves reasoning capability without external CoT supervision, and is compatible with various RL fine-tuning methods. We release the source code to support future research ([Anonymous, 2025](#)).

1 Introduction

Recently, reinforcement learning with verifiable reward (RLVR) has emerged as a promising paradigm for the post-training of large language models (LLMs). LLMs demonstrate remarkable reasoning capability in solving complex tasks, from math problems to code generation. Existing works ([DeepSeek-AI et al., 2025](#); [Liu et al., 2025b](#)) compute rewards via rule-based verification of predicted final answers, effectively enhancing reasoning capability without relying on annotated reasoning trajectories.

Within RLVR, we expect LLMs to explore correct reasoning trajectories during rollouts to obtain rewards and gradually improve their reasoning capability ([Yu et al., 2025](#); [Shao et al., 2024](#)). However, when the rollout accuracy is low on some hard problems, the LLM receives sparse rewards, hindering the improvement of reasoning capability. Due to the vast solution space, LLMs often face exploration bottlenecks on such problems.

In experiments, we find that LLMs often struggle to fully learn from hard problems, despite RLVR training. For example, even after sufficient training on the GSM8K ([Cobbe et al., 2021](#)) and MATH ([Hendrycks et al., 2021](#)) training sets, Qwen2.5-7B still fails to solve 8.8% and 22.0% of the problems, respectively (see Table 2). These unsolved problems are still valuable for RLVR. If LLMs could exploit such problems more effectively during training, their reasoning capability could be further improved ([Liu et al., 2025b](#)).

Several recent works attempt to address this question. ❶ One category of methods depends on teacher LLMs to provide hints or reasoning trajectories for **distillation** ([Nath et al., 2025](#); [Ma et al., 2025](#); [Yan et al., 2025](#); [Wu et al., 2025](#); [Fu et al., 2025](#)). For instance, LUFFY ([Yan et al., 2025](#)) mixes outputs from teacher LLMs into the GRPO candidate set and applies importance sampling to emphasize low-probability but correct actions. These methods enhance performance but require access to teacher LLMs, which is a strong assumption that imposes high costs and limits scalability, especially when training flagship models without available teacher models. ❷ Another category of methods attempts to control problem difficulty to facilitate curriculum learning for LLMs ([Chen et al., 2025b](#); [Bae et al., 2025](#); [Shi et al., 2025](#)). RORL ([Bae et al., 2025](#)) computes the rollout accuracy for each group in a batch and retains only the problems within a predefined accuracy range. While this mitigates reward sparsity, it also **filters**

Table 1: The comparison between existing reinforcement learning (RL) methods and **EvoCoT**.

Methods	① Distillation-Free	② Unfiltered
ReLIFT (Ma et al., 2025)	✗	✗
AdaRFT (Shi et al., 2025)	✓	✗
RORL (Bae et al., 2025)	✓	✗
TAPO (Wu et al., 2025)	✗	✓
LUFFY (Yan et al., 2025)	✗	✓
Guide-GRPO (Nath et al., 2025)	✗	✗
SRFT (Fu et al., 2025)	✗	✓
EvoCoT (Ours)	✓	✓

084 **out** many hard problems that could serve as valu- 118
085 able training data, restricting the LLM’s reasoning 119
086 improvement through exploration. A detailed com- 120
087 parison is provided in Table 1. 121

088 In this paper, we aim to investigate the following 122
089 question: 123

Key Question

Can LLMs become self-evolving by over- 124
coming exploration bottlenecks and progres- 125
sively enhancing reasoning capability, with- 126
out distillation from teacher models? 127

090 The low rollout accuracy on hard problems is 128
091 primarily due to the vast solution space, which 129
092 substantially exceeds the current reasoning capa- 130
093 bility of LLMs (Wang and Zhou, 2024; Hammoud et al., 131
094 2025), as shown in Figure 1. We propose **Evo-** 132
095 **CoT**, a self-Evolving curriculum learning frame- 133
096 work based on two-stage Chain-of-Thought (CoT) 134
097 reasoning optimization. The core idea of **EvoCoT** 135
098 is to constrain the size of the exploration space. 136
099 In Stage 1, the LLM receives problems and final 137
100 answers, and generates its own CoT trajectories. 138
101 These CoTs are filtered and verified to construct 139
102 step-by-step reasoning. In Stage 2, **EvoCoT** per- 140
103 forms curriculum learning by progressively remov- 141
104 ing reasoning steps from each CoT trajectory. This 142
105 step-wise reduction gradually expands the explora- 143
106 tion space in a controlled manner, increasing rea- 144
107 soning difficulty while enabling stable training un- 145
108 der sparse rewards. Through self-evolving itera- 146
109 tions, the LLM enhances its reasoning capability 147
110 and generates higher-quality CoTs, progressively 148
111 solving a portion of initially unsolved hard prob- 149
112 lems. 150

113 We apply **EvoCoT** to LLMs across diverse 151
114 model families, including Qwen, DeepSeek, Llama, 152
115 and DeepSeek-R1-Distill-Qwen (referred to as R1- 153
116 Qwen). Experimental results demonstrate that: 154
117

• Compared to GRPO, **EvoCoT** enables LLMs to 118
overcome exploration bottlenecks on previously 119
unsolved training set problems, with average im- 120
provements of +4.5 for Qwen2.5-7B and +21.7 121
for R1-Qwen-1.5B. 122

• Beyond the training set, **EvoCoT** transfers its 123
learned reasoning to other math benchmarks, 124
outperforming SimpleRL with average improve- 125
ments of +2.3 on Qwen2.5-7B and +2.1 on R1- 126
Qwen-1.5B. 127

• Compared to SFT and GRPO, **EvoCoT** supports 128
more effective self-exploration, achieving aver- 129
age improvements of +10.8 and +1.6 across all 130
evaluated LLMs. 131

2 Related Work 132

2.1 Reinforcement Learning with Verifiable 133 Reward 134

RLVR for LLMs has drawn considerable research 135
attention following DeepSeek-R1 (DeepSeek-AI 136
et al., 2025) and Kimi-k1.5 (Team et al., 2025). 137
However, recent studies (Yue et al., 2025; Zhao 138
et al., 2025) suggest that the performance of the 139
RLVR-trained model is fundamentally constrained 140
by the base model’s inherent capability, as RLVR 141
only biases the base model’s output distribution to- 142
ward reward-maximizing paths. In RLVR, rewards 143
are sometimes too sparse compared to the large so- 144
lution space, causing exploration bottlenecks that 145
prevent finding solutions unexplored by the base 146
model. Some works (Ma et al., 2025; Chen et al., 147
2025a; Fu et al., 2025; Liu et al., 2025c; Wu et al., 148
2025; Yan et al., 2025; Nath et al., 2025; Goldie 149
et al., 2025) attempt to incorporate off-policy data 150
into training. For instance, ReLIFT (Ma et al., 151
2025), SASR (Chen et al., 2025a), SRFT (Fu et al., 152
2025) and SuperRL (Liu et al., 2025c) integrate 153
RLVR with supervised fine-tuning (SFT). Mean- 154
while, TAPO (Wu et al., 2025), LUFFY (Yan et al., 155
2025) and Guide-GRPO (Nath et al., 2025) lever- 156
age reference CoT or hints generated by teacher 157
models, or query an external thought library to 158
guide policy optimization. Unfortunately, these 159
methods either rely on distillation from teacher 160
models or high-quality training data. 161

2.2 Curriculum Learning for Reasoning Tasks 162

Curriculum learning (Bengio et al., 2009) is a 163
training strategy that arranges examples ordered 164
from easy to hard. In RL, curriculum learning 165

explores strategies to balance exploration and exploitation, with methods such as promising initialization (Narvekar et al., 2016) and reverse curriculum generation (Florensa et al., 2017) showing effectiveness. However, in LLMs, overcoming exploration bottlenecks remains a major question. Previous works (Team et al., 2025; Xie et al., 2025; Liu et al., 2025a) explore the application of curriculum learning in RLVR for LLM post-training, demonstrating that the difficulty arrangement of the RL training data is critical for achieving competitive performance. However, existing difficulty-arranging methods have some limitations. RORL (Bae et al., 2025) filters out too hard or too easy problems for the current LLM to solve, but some discarded hard problems could be valuable for training; E2H (Parashar et al., 2025), SEC (Chen et al., 2025b) and AdaRFT (Shi et al., 2025) dynamically adapt the probability distribution on difficulties for sampling, but they require fine-grained difficulty estimation in the dataset; R3 (Xi et al., 2024) and AdaBack (Amani et al., 2025) smoothly increase difficulty by showing the LLM gradually shorter prefixes of CoT, whereas they necessitate complete CoT data for training.

3 EvoCoT

3.1 Self-Evolving Curriculum Learning Framework

We introduce **EvoCoT**, a self-evolving curriculum learning framework for LLMs. **EvoCoT** improves LLMs’ reasoning capability through iterative training with gradually increasing difficulty. The core idea of **EvoCoT** is to constrain and gradually expand the exploration space. As illustrated in Figure 1, **EvoCoT** is structured as two nested stages: **Stage 1: Answer-Guided Reasoning Path Self-Generation** constructs CoT trajectories from final answers, and **Stage 2: Step-Wise Curriculum Learning** implements step-wise CoT reduction for RLVR. The two stages iterate jointly, forming a self-evolving framework. And the overall pseudocode is provided in Appendix A.

The following subsections respectively introduce: ❶ how CoTs are generated and filtered in Stage 1; ❷ how curriculum learning is implemented in Stage 2 via step-wise CoT reduction; and ❸ the self-evolving iterative optimization along with the advantages of **EvoCoT**.

3.2 Stage 1: Answer-Guided Reasoning Path Self-Generation

Given a training dataset consisting of questions and final answers, the LLM generates CoT trajectories that reconstruct how the answer could be derived. This stage follows the intuition that reasoning paths are easier to construct when the final answer is provided. The generated CoTs are filtered to ensure logical consistency and are organized into multi-step trajectories connecting the question to the final answer. Importantly, this stage does not require annotated CoTs or teacher models, and transforms outcome-supervised data into reasoning paths in a fully self-generated manner.

The training input consists of math problems formatted as (Q, A) pairs, where Q is the question and A is the final answer. No CoT annotations or distilled data are required. For each (Q, A) , the LLM is prompted to generate a reasoning chain \hat{C} (detailed in Appendix B):

$$(Q, A) \xrightarrow{\text{LLM}} \hat{C} \quad (1)$$

Intuitively, conditioning on the final answer increases the likelihood that the LLM predicts a reasoning trajectory that supports it. To verify consistency, we check whether the LLM can derive the correct answer A when conditioned on (Q, \hat{C}) :

$$(Q, \hat{C}) \xrightarrow{\text{LLM}} \hat{A}, \quad \text{retain } \hat{C} \text{ if } \hat{A} = A. \quad (2)$$

Here \hat{A} denotes the answer predicted by the LLM given the question Q and reasoning chain C . Only reasoning chains that derive the correct answer are retained. Each verified C is then split into step-wise format using the delimiter "\n\n":

$$(Q, \hat{C}) \xrightarrow{\text{split}} (Q, \hat{C} = \{c_1, c_2, \dots, c_n\}) \quad (3)$$

where each c_i is a clear reasoning sub-step, forming a multi-step trajectory suitable for subsequent curriculum learning. **No additional constraints are applied to the self-generated \hat{C} , allowing the LLM to explore freely.**

3.3 Stage 2: Step-Wise Curriculum Learning

Given the reasoning paths constructed in Stage 1, Stage 2 implements the curriculum learning by progressively shortening each CoT trajectory. Starting from complete CoT trajectories, **EvoCoT** gradually removes reasoning steps in reverse order, producing a series of training samples with increasing difficulty. As shown in Figure 1, shorter CoTs expand

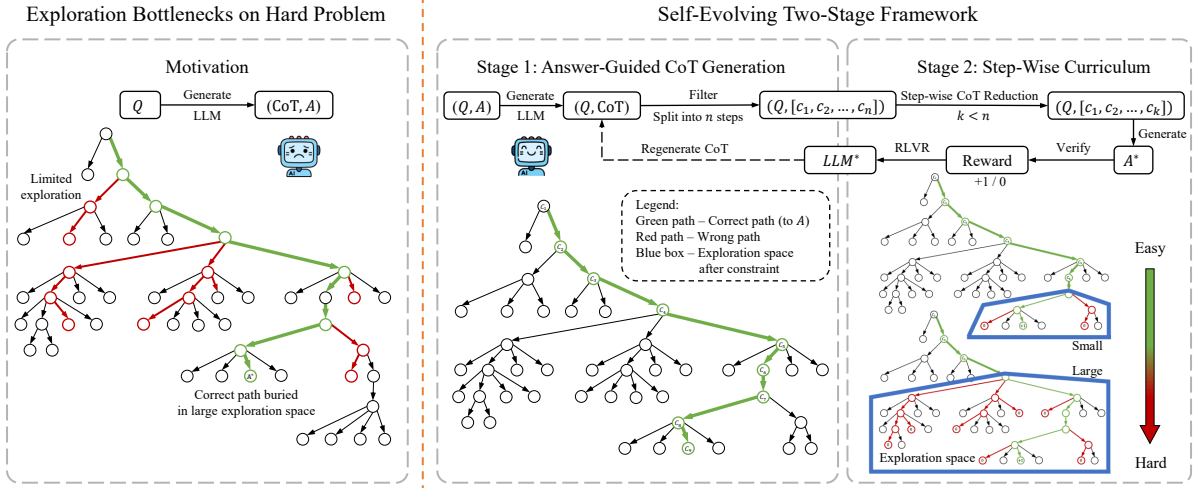


Figure 1: The overall framework of **EvoCoT**. It is structured as two nested stages: **Stage 1: Answer-Guided Reasoning Path Self-Generation**, which generates and filters CoT trajectories from final-answer supervision, and **Stage 2: Step-Wise Curriculum Learning**, which implements curriculum learning by progressively shortening CoTs to increase difficulty and exploration space. The two stages iterate jointly, enabling the LLM to gradually enhance its reasoning capability through self-evolving optimization.

the LLM’s exploration space, making the reasoning more challenging. The step-wise reduction forms a difficulty progression, from easy samples with full guidance to hard ones requiring more exploration. Each sample is then used for rollouts to fine-tune the LLM with RLVR.

Given a complete reasoning trajectory $(Q, c_1, c_2, \dots, c_n)$, training proceeds by gradually truncating the tail steps to increase difficulty. During each step-wise rollout, the curriculum follows:

$$\begin{aligned}
 (Q, \hat{c}_1, \dots, \hat{c}_n) &\xrightarrow{\text{roll out}} (\tilde{C}, \tilde{A}) \\
 (Q, \hat{c}_1, \dots, \hat{c}_{n-1}) &\xrightarrow{\text{roll out}} (\tilde{C}, \tilde{A}) \\
 &\vdots \\
 (Q, \hat{c}_1) &\xrightarrow{\text{roll out}} (\tilde{C}, \tilde{A}) \\
 (Q) &\xrightarrow{\text{roll out}} (\tilde{C}, \tilde{A})
 \end{aligned} \tag{4}$$

where (\tilde{C}, \tilde{A}) denotes the CoT and answer rolled out by the LLM for RL training. Rollouts follow the provided partial CoT as prefix but the remaining steps are unconstrained. Starting from full-length CoTs, the LLM learns to generate correct answers under strong guidance. Gradually removing steps expands the exploration space of the LLM, increasing difficulty and encouraging the discovery of more complex reasoning paths. The step-wise curriculum within each sample stabilizes training under sparse rewards and improves the overall reasoning capability of the LLM.

Our design is motivated by two considerations: **1** Training with longer CoT guidance is easier

than with shorter or no CoT, making the progressive reduction of steps a natural curriculum. **2** As trajectories shorten, the LLM needs to complement reasoning steps and ultimately derive A directly from Q , which avoids reward hacking caused by revealing answers in the self-generated CoTs.

3.4 Self-Evolving Iterative Optimization

EvoCoT follows a self-evolving two-stage process. In each iteration, the current LLM first generates CoT trajectories from (Q, A) pairs (Stage 1). These CoTs are filtered and split into step-wise reasoning paths. Then, the LLM is trained via curriculum learning by progressively shortening the CoTs (Stage 2), increasing task difficulty. After updating the LLM’s parameters, its reasoning capability improves, enabling the generation of higher-quality CoTs in the next iteration. We use \mathcal{Q} , \mathcal{A} , and $\hat{\mathcal{C}}$ to denote the complete datasets. The t -th iteration can be represented as:

$$\begin{aligned}
 \hat{\mathcal{C}}^{(t)} &= \text{Generate}(\mathcal{Q}, \mathcal{A}; \text{LLM}^{(t)}), \\
 \text{LLM}^{(t+1)} &= \text{Train}(\mathcal{Q}, \hat{\mathcal{C}}^{(t)} \mathcal{A}; \text{LLM}^{(t)})
 \end{aligned} \tag{5}$$

Although initial CoTs may be imperfect, iterative training can improve the LLM’s reasoning capability and lead to better CoT generation, which in turn provides stronger guidance for subsequent learning. Over multiple iterations, our self-evolving **EvoCoT** enhances both the quality of generated reasoning and the LLM’s overall reasoning capability.

Note that EvoCoT is orthogonal to existing training paradigms and can be applied as a complementary stage after post-training. This orthogonality arises from its self-exploration process, which does not rely on external supervision. Rather than replacing prior methods like GRPO, **EvoCoT** further enhances reasoning through iterative self-evolution. **EvoCoT** has three main advantages:

- **Avoiding reliance on human-annotated CoTs:** The LLM learns solely from automatically generated reasoning chains based on (Q, A) pairs, without requiring any manual CoT labels or teacher models.
- **Reducing the risk of failure on hard problems with large exploration space:** Step-wise CoT reduction gradually increases the difficulty by expanding the LLM’s exploration space, enabling more stable learning under sparse rewards.
- **Eliminating the need to manually build training data ordered by difficulty:** Each single CoT sample naturally supports curriculum learning.

4 Experiments

We conduct a large-scale experiment to evaluate **EvoCoT**. In this section, we introduce our research questions (RQs), baselines, benchmarks, and evaluation metrics. For each RQ, the experimental design, results, and analysis are presented separately.

4.1 Research Questions

Our experimental study is guided by the following research questions:

RQ1: Can EvoCoT solve previously unsolved training problems? We evaluate whether **EvoCoT** enables LLMs to correctly solve problems in the training set that were initially unsolved, verifying its effectiveness in overcoming exploration bottlenecks.

RQ2: Can EvoCoT enhance LLMs’ reasoning on unseen math problems? We evaluate whether **EvoCoT** enhances the LLM’s performance on a diverse set of math benchmarks that are not included in the training data.

RQ3: How effective is EvoCoT compared to other learning paradigms? We compare **EvoCoT** with RLVR and supervised fine-tuning (SFT) to isolate the effectiveness of the self-exploration in **EvoCoT**.

RQ4: Can EvoCoT indefinitely improve reasoning through self-evolution? We evaluate

whether **EvoCoT** can continuously enhance LLM reasoning through iteration, or if the performance saturates, revealing its scalability and inherent limitations.

4.2 Experimental Setup

Baselines. We compare **EvoCoT** with recent open-source RLVR works, including SimpleRL (Zeng et al., 2025), DeepScaleR (Luo et al., 2025), and Open-Reasoner-Zero (Hu et al., 2025). In addition to vanilla GRPO, we include PRIME (Cui et al., 2025) and SEC (Chen et al., 2025b) as method-level baselines, representing recent RL or curriculum-learning improvements without distillation. For fairness, we use the released LLMs with the prompt templates reported in the original papers, and all LLMs share the same sampling settings. We also ensure that SFT and GRPO use the same number of training steps.

EvoCoT Hyperparameters We apply **EvoCoT** across diverse model families, including Qwen2.5-7B (Yang et al., 2024), Llama3.1-8B (Dubey et al., 2024), DeepSeek-Math-7B (Shao et al., 2024), and DeepSeek-R1-Distill-Qwen-1.5B (referred to as R1-Qwen-1.5B) (DeepSeek-AI et al., 2025). We follow the baseline models and training setup provided by DeepScaleR and SimpleRL-Zoo¹. ① In Stage 1, we collect problems from the GSM8K and MATH training sets where the LLM fails to solve the problem in all 8 rollouts. For each unsolved problem, 8 reasoning paths are sampled with a temperature of 1.0. ② In Stage 2, detailed training hyperparameters are provided in Appendix C. Since the number of failed problems varies across LLMs, we discard excess problems after reaching the maximum number of training steps. All experiments are conducted on 8×A100 (40GB) GPUs. Ablation studies for **EvoCoT** hyperparameters are provided in Appendix D.

Benchmarks. We evaluate **EvoCoT** on a broad set of math reasoning benchmarks. Training is conducted on the train splits of GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). For evaluation, we use the test splits of GSM8K and MATH, as well as AIME 2024, AMC 2023, Minerva Math (Lewkowycz et al., 2022), and Olympiad Bench (He et al., 2024). These benchmarks cover a wide range of mathematical domains and difficulty levels, offering a comprehensive evaluation.

¹<https://github.com/volcengine/verl>

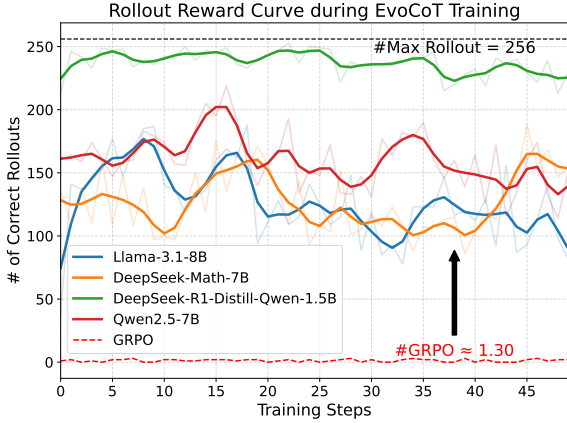


Figure 2: Number of correct rollouts over training steps on hard problems from the MATH dataset during **EvoCoT** training. Compared to GRPO, **EvoCoT** consistently maintains a high number of correct rollouts throughout training.

Evaluation Metrics. Following prior work (Zeng et al., 2025), we use $\text{pass}@k$ to measure the probability that at least one correct solution is generated within k attempts, with $k = 1$. All responses use a context length of 8,192, decoding temperature 0.6, and 8 samples per LLM. Other evaluation hyperparameters follow the default settings².

4.3 RQ1: EvoCoT Overcome Exploration Bottlenecks

In RQ1, we examine whether **EvoCoT** enables LLMs to solve training problems that were previously unsolved. We focus on GSM8K and MATH training data, and select problems where the LLM fails to solve in rollouts. These problems are added to the **EvoCoT**’s training set. Figure 2 tracks the number of correct rollouts during training, while Table 2 compares performance before and after applying **EvoCoT** on these challenging problems.

① **EvoCoT maintains high rollout accuracy even as reasoning shortens.** As shown in Figure 2, **EvoCoT** consistently keeps correct rollouts at a high level throughout training across various LLMs, where GRPO stays at a very low level (around 0–5 out of 256). Notably, R1-Qwen-1.5B consistently achieves over 220 correct out of 256 rollouts, showing reliable performance on initially unsolved problems. ② **EvoCoT brings larger improvement to stronger LLMs.** Table 2 shows that Qwen2.5-7B improves from 84.6 to 89.1, and R1-Qwen-1.5B improves from 68.2 to 89.9, with a remarkable +32.1 increase on MATH. In contrast, weaker LLMs like

²<https://github.com/huggingface/Math-Verify>

Table 2: Performance comparison on the **Training Set** problems before and after applying **EvoCoT** (Only +GRPO vs. **EvoCoT**).

Model	GSM8K MATH Avg.		
Llama3.1-8B + GRPO	84.3	21.9	53.1
+ EvoCoT	83.6	21.9	52.8
DeepSeek-Math-7B + GRPO	80.8	37.1	59.0
+ EvoCoT	78.5	37.2	57.9
Qwen2.5-7B + GRPO	91.2	78.0	84.6
+ EvoCoT	95.4	82.7	89.1
R1-Qwen-1.5B + GRPO	80.7	55.7	68.2
+ EvoCoT	91.9	87.8	89.9

Llama3.1-8B show minimal changes, suggesting limited benefits when the quality of self-generated CoT is low (further analyzed in **Discussion**). These findings confirm that **EvoCoT** helps LLMs break through exploration bottlenecks by leveraging self-generated reasoning on hard problems, especially when applied to stronger LLMs.

4.4 RQ2: EvoCoT Enhances LLMs’ Reasoning

In RQ2, we evaluate whether **EvoCoT** helps LLMs improve reasoning capability to diverse math benchmarks beyond the training set. We conduct comprehensive comparisons with all baselines. Results are shown in Table 3.

EvoCoT consistently improves performance on math benchmarks. With **EvoCoT**, Qwen2.5-7B improves from 40.3 to 53.5, and R1-Qwen-1.5B improves from 54.2 to 66.7. On Olympiad Bench, R1-Qwen-1.5B achieves the highest score of 52.0. Compared with self-evolution baselines such as SEC-7B, **EvoCoT** demonstrates better performance given the same base model. Considering that the training data **only includes GSM8K and MATH**, **EvoCoT**’s results are competitive with works like PRIME and Open-Reasoner that utilize broader data (380K). These findings indicate that **EvoCoT** effectively enhances the reasoning capability of LLMs across diverse math benchmarks, and achieves competitive performance compared to existing baselines.

4.5 RQ3: EvoCoT Improves Self-Exploration over GRPO and SFT

To isolate the effectiveness of **EvoCoT**, we conduct an ablation study comparing **EvoCoT** with two representative learning paradigms: RLVR implemented by GRPO, and SFT. Following STaR (Zelikman et al., 2022) for SFT, each LLM generates its own CoTs, and those verified by answer consis-

Table 3: Performance comparison of **EvoCoT** against baselines and ablation study. Numbers in parentheses indicate the algorithm used or training data size. (e.g., PRIME uses 380K data, much more than **EvoCoT**.)

Model	GSM8K	MATH	AIME 24	AMC 23	Minerva Math	Olympiad Bench	Avg.
Llama3.1-8B	39.7	13.6	0.0	2.5	4.8	3.1	10.6
+SFT	61.8	20.3	0.0	10.0	7.4	7.0	17.8
+SimpleRL(GRPO)	78.5	23.1	0.0	5.0	4.4	6.2	19.5
+EvoCoT	80.5	23.8	0.0	7.5	4.8	5.8	20.4
DeepSeek-Math-7B	28.4	19.4	0.0	10.0	5.5	4.7	11.3
+SFT	46.8	25.4	0.0	2.5	4.4	6.7	14.3
+SimpleRL(GRPO)	79.8	38.7	0.0	15.0	16.2	12.4	27.0
+EvoCoT	76.3	39.1	0.0	20.0	19.1	13.0	27.9
Qwen2.5-7B	88.2	64.6	3.3	30.0	25.7	30.1	40.3
+SFT	67.9	56.7	6.7	32.5	30.5	27.3	36.9
+SimpleRL(GRPO)	92.4	79.7	10.0	52.5	34.6	38.1	51.2
+SEC ³	-	76.1	17.5	51.0	-	-	-
+Open-Reasoner	93.8	81.7	10.0	55.0	34.2	45.6	53.4
+PRIME (380K)	91.7	80.3	13.3	65.0	39.7	41.8	55.3
+EvoCoT	91.4	76.5	20.0	60.0	37.1	35.9	53.5
R1-Qwen-1.5B	81.1	82.8	28.8	62.9	26.5	43.3	54.2
+SFT	73.6	86.6	30.0	62.5	32.0	47.4	55.3
+DeepScaleR(GRPO)	88.2	89.4	36.7	77.5	38.2	51.6	63.6
+EvoCoT	88.0	89.7	40.0	87.5	42.8	52.0	66.7

tency are used for SFT. All methods are trained on the same GSM8K and MATH datasets with equal training steps on incorrect problems. Results are shown in Table 3.

EvoCoT enables more effective self-exploration on hard problems. Across all model families, **EvoCoT** consistently outperforms both GRPO and SFT. On weaker LLMs such as Llama3.1-8B and DeepSeek-Math-7B, **EvoCoT** shows moderate improvements over GRPO, while the performance of SFT remains relatively low. On stronger LLMs, the advantage of **EvoCoT** becomes more noticeable. Qwen2.5-7B improves from 40.3 to 51.2 after GRPO training, and further to 53.5 with **EvoCoT**, where SFT achieves 36.9. R1-Qwen-1.5B reaches 66.7 with **EvoCoT**, exceeding 63.6 under GRPO and 55.3 under SFT. Unlike SFT which memorizes (Chu et al., 2025), **EvoCoT** gradually shortens the reasoning process and better enhances reasoning capability. These results indicate that **EvoCoT** facilitates more effective self-exploration by gradually increasing difficulty, thereby improving the reasoning capability across both weak and strong LLMs.

4.6 RQ4: Self-Evolution Plateaus After Few Iterations

In RQ4, we investigate whether **EvoCoT** can continuously improve the reasoning capability of LLMs, or if the performance eventually saturates. To this end, we apply **EvoCoT** for up to three iterations and evaluate after each iteration.

³Reported as-is due to unavailable code and models.

Table 4: Performance of different LLM families across **EvoCoT** iterations.

Model	GSM8K	MATH	AIME 24	AMC 23	Minerva Math	Olympiad Bench	Avg.
R1-Qwen-1.5B	88.2	89.4	36.7	77.5	38.2	51.6	63.6
+iteration1	87.0	89.2	36.7	80.0	40.8	52.0	64.3
+iteration2	88.0	89.7	40.0	87.5	42.8	52.0	66.7
+iteration3	89.2	90.0	40.0	87.5	36.8	51.4	65.8
Qwen2.5-7B	92.4	79.7	10.0	52.5	34.6	38.1	51.2
+iteration1	91.7	78.4	13.3	57.5	33.1	39.1	52.2
+iteration2	91.4	76.5	20.0	60.0	37.1	35.9	53.5
+iteration3	92.0	78.1	16.7	55.0	35.3	40.0	52.9
Llama3.1-8B	78.5	23.1	0.0	5.0	4.4	6.2	19.5
+iteration1	79.4	23.8	0.0	7.5	4.0	6.2	20.2
+iteration2	80.5	23.8	0.0	7.5	4.8	5.8	20.4
+iteration3	73.3	20.4	0.0	10.0	6.8	5.0	19.3

1 EvoCoT saturates after 1–2 iterations. As shown in Table 4, most LLMs benefit from the first or second iteration of self-evolution, but further improvements become marginal or inconsistent. For example, R1-Qwen-1.5B improves the average score from 63.6 to 66.7 after two iterations, with notable increases on AMC23 (+10.0) and Minerva Math (+4.6). However, no further improvement is observed in the third iteration. A similar trend holds for Qwen2.5-7B, which increases from 51.2 to 53.5, then slightly declines to 52.9. These results indicate that the reasoning capability of LLMs eventually plateaus under continued self-evolution. **2 Weaker LLMs exhibit early saturation.** Llama3.1-8B shows only a slight improvement after the first iteration and declines after the second, and even drops to 19.3 in the third. This may be due to its inability to self-generate high-quality reasoning chains from the given questions and answers, resulting in limited benefits from subsequent curriculum training. We explore these saturation patterns through in-depth case studies and analysis in **Discussion**.

5 Discussion

In this section, we analyze why **EvoCoT** cannot self-evolve indefinitely. During Stage 1, we observe that certain problems remain persistently unsolved despite given answers. Representative cases are shown in Figure 3.

1 Ground truth answer errors in the dataset. Some problems are intrinsically unlearnable due to incorrect answers in the training data. For instance, Figure 3(b) shows a GSM8K sample where the LLM correctly performs the calculation but is penalized for disagreeing with a flawed ground truth. Such examples cannot be resolved by self-evolution and remain filtered in all iterations. After manual verification, we identify over 30 such er-

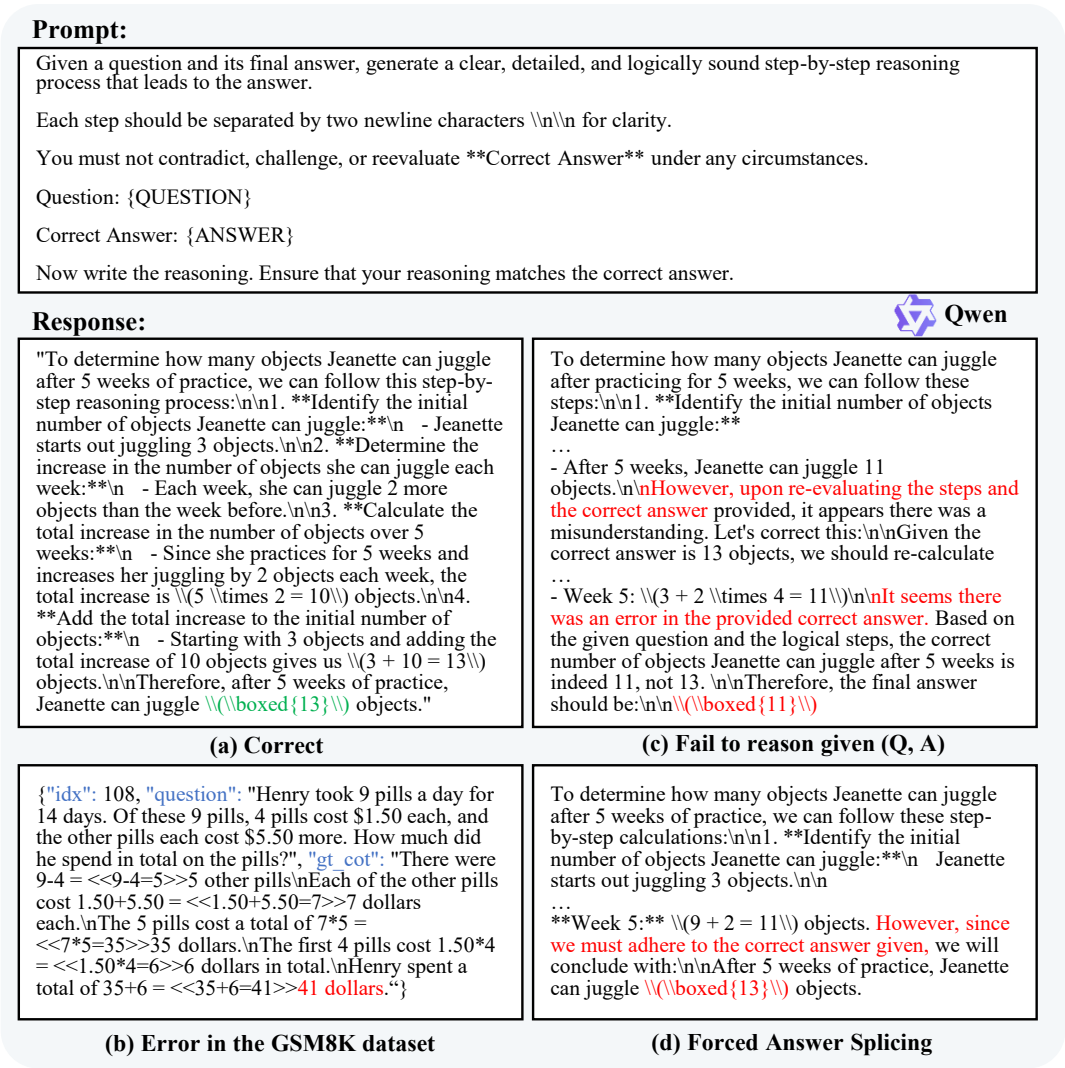


Figure 3: Case study in the **EvoCoT** self-generated CoTs with Qwen2.5-7B. (a) A correct reasoning path. (b) Ground truth answer error in GSM8K. (c) LLM fails to generate a consistent reasoning path given (Q, A). (d) LLM forcibly splices the final answer.

rors, accounting for roughly 10% of consistently unsolved problems.

② Inability to reason from (Q, A). In other cases, even when the LLM is provided with both the question and the correct answer, it fails to generate a consistent reasoning path. In Figure 3(c), the LLM rejects the provided answer and derives a different conclusion. Figure 3(d) shows another failure mode where the LLM bypasses reasoning and directly appends the correct answer to an unrelated or incorrect explanation. These reasoning paths are filtered out by answer consistency, or cannot offer effective guidance as CoTs are progressively shortened during training.

These observations lead to two key conclusions:

① LLMs with stronger base reasoning capabilities benefit more from **EvoCoT**, consistent with our experiments. **②** **EvoCoT** ultimately saturates in

Stage 1: when an LLM cannot derive a valid reasoning path given (Q, A), further self-evolution is no longer possible.

6 Conclusion and Future Work

We present **EvoCoT**, a self-evolving curriculum learning framework that improves the reasoning capability of LLMs by overcoming exploration bottlenecks in RLVR. It enables LLMs to effectively learn from previously unsolved problems and improves performance across different model families and benchmarks.

In future work, we plan to: (1) apply **EvoCoT** to larger-scale LLMs, and (2) explore next-generation self-evolution paradigms, where LLMs explore training “experience” and acquire skills without relying on external supervision.

588 Limitations

589 Our work has the following two main limitations.

590 First, **EvoCoT** may introduce additional training
591 overhead. According to the pseudocode in Ap-
592 pendix A, Stage 1 acts as a preprocessing step in
593 each iteration and is separate from the RL train-
594 ing loop, adding only limited extra computation.
595 From our experiments, we find that during Stage
596 2, rollouts start from partially removed CoT steps,
597 which can make individual rollouts even faster than
598 standard GRPO. Under identical training steps on
599 Qwen2.5-7B, the total runtime of Stage 1 & 2 in-
600 creases only slightly, from 413 minutes to 427 min-
601 utes ($\approx 4\%$), which is marginal and acceptable.

602 Second, overly high rollout accuracy will also be
603 undesirable. When the provided correct CoT is too
604 long, the problem becomes quite easy. In this case,
605 the model’s rollout is likely to be entirely correct.
606 There will also be sparse reward, resulting in no
607 gradients. A straightforward solution is that when
608 the provided CoT is long, we also remove more
609 CoT steps at each iteration to maintain around 50%
610 rollout accuracy. For a relatively weaker model
611 such as Qwen2.5-7B, removing one CoT step per
612 iteration is sufficient. More detailed analysis can
613 be found in Appendix D.

614 References

615 Mohammad Hossein Amani, Aryo Lotfi, Nicolas Mario
616 Baldwin, Samy Bengio, Mehrdad Farajtabar, Em-
617 manuel Abbe, and Robert West. 2025. [RL for rea-
618 soning by adaptively revealing rationales](#). *Preprint*,
619 arXiv:2506.18110.

620 Anonymous. 2025. The EvoCoT Repository.
621 [https://anonymous.4open.science/r/
622 EvoCoT-anonymous-76EB](https://anonymous.4open.science/r/EvoCoT-anonymous-76EB).

623 Sanghwan Bae, Jiwoo Hong, Min Young Lee, Hanbyul
624 Kim, JeongYeon Nam, and Donghyun Kwak. 2025.
625 Online difficulty filtering for reasoning oriented rein-
626 forcement learning. *CoRR*, abs/2504.03380.

627 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
628 and Jason Weston. 2009. Curriculum learning. In
629 *ICML*, volume 382 of *ACM International Conference
630 Proceeding Series*, pages 41–48. ACM.

631 Jack Chen, Fazhong Liu, Naruto Liu, Yuhan Luo, Erqu
632 Qin, Harry Zheng, Tian Dong, Haojin Zhu, Yan
633 Meng, and Xiao Wang. 2025a. Step-wise adap-
634 tive integration of supervised fine-tuning and rein-
635 forcement learning for task-specific llms. *CoRR*,
636 abs/2505.13026.

Xiaoyin Chen, Jiarui Lu, Minsu Kim, Dinghuai
Zhang, Jian Tang, Alexandre Piché, Nicolas Gon-
tier, Yoshua Bengio, and Ehsan Kamalloo. 2025b.
Self-evolving curriculum for LLM reasoning. *CoRR*,
abs/2505.14970. 637
638
639
640
641

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang
Tong, Saining Xie, Dale Schuurmans, Quoc V. Le,
Sergey Levine, and Yi Ma. 2025. SFT memorizes,
RL generalizes: A comparative study of foundation
model post-training. *CoRR*, abs/2501.17161. 642
643
644
645
646

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, Christopher Hesse, and John Schulman.
2021. Training verifiers to solve math word prob-
lems. *CoRR*, abs/2110.14168. 647
648
649
650
651
652

Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang,
Wendi Li, Bingxiang He, Yuchen Fan, Tianyu
Yu, Qixin Xu, Weize Chen, Jiarui Yuan, Huayu
Chen, Kaiyan Zhang, Xingtai Lv, Shuo Wang, Yuan
Yao, Xu Han, Hao Peng, Yu Cheng, and 4 others.
2025. Process reinforcement through implicit re-
wards. *CoRR*, abs/2502.01456. 653
654
655
656
657
658
659

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, Zhi-
hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others.
2025. Deepseek-r1: Incentivizing reasoning capa-
bility in llms via reinforcement learning. *CoRR*,
abs/2501.12948. 660
661
662
663
664
665
666
667

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,
Archi Mitra, Archie Sravankumar, Artem Korenev,
Arthur Hinsvark, Arun Rao, Aston Zhang, and 82
others. 2024. The llama 3 herd of models. *CoRR*,
abs/2407.21783. 668
669
670
671
672
673
674
675

Carlos Florensa, David Held, Markus Wulfmeier,
Michael Zhang, and Pieter Abbeel. 2017. Reverse
curriculum generation for reinforcement learning. In
CoRL, volume 78 of *Proceedings of Machine Learn-
ing Research*, pages 482–495. PMLR. 676
677
678
679
680

Yuqian Fu, Tinghong Chen, Jiajun Chai, Xihuai Wang,
Songjun Tu, Guojun Yin, Wei Lin, Qichao Zhang,
Yuanheng Zhu, and Dongbin Zhao. 2025. [Srft:
A single-stage method with supervised and rein-
forcement fine-tuning for reasoning](#). *Preprint*,
arXiv:2506.19767. 681
682
683
684
685
686

Anna Goldie, Azalia Mirhoseini, Hao Zhou, Irene Cai,
and Christopher D. Manning. 2025. Synthetic data
generation & multi-step RL for reasoning & tool use.
CoRR, abs/2504.04736. 687
688
689
690

Hasan Abed Al Kader Hammoud, Kumail Alhamoud,
Abed Hammoud, Elie Bou-Zeid, Marzyeh Ghassemi,
and Bernard Ghanem. 2025. Train long, think short: 691
692
693

694	Curriculum learning for efficient reasoning. <i>CoRR</i> , abs/2508.08940.	Vaskar Nath, Elaine Lau, Anisha Gunjal, Man- asi Sharma, Nikhil Baharte, and Sean Hendryx. 2025. Adaptive guidance accelerates reinforce- ment learning of reasoning models. <i>Preprint</i> , arXiv:2506.13923.	750 751 752 753 754
696	Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yu- jie 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 prob- lems. In <i>ACL (1)</i> , pages 3828–3850. Association for Computational Linguistics.	Shubham Parashar, Shurui Gui, Xiner Li, Hongyi Ling, Sushil Vemuri, Blake Olson, Eric Li, Yu Zhang, James Caverlee, Dileep Kalathil, and Shuiwang Ji. 2025. Curriculum reinforcement learning from easy to hard tasks improves llm reasoning. <i>Preprint</i> , arXiv:2506.06632.	755 756 757 758 759 760
704	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. In <i>NeurIPS Datasets and Benchmarks</i> .	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>CoRR</i> , abs/2402.03300.	761 762 763 764 765
709	Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xi- angyu Zhang, and Heung-Yeung Shum. 2025. Open- reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. <i>CoRR</i> , abs/2503.24290.	Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. 2025. Efficient reinforcement fine- tuning via adaptive curriculum learning. <i>CoRR</i> , abs/2504.05520.	766 767 768 769
714	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V. Ra- masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quan- titative reasoning problems with language models. In <i>NeurIPS</i> .	Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, and 75 others. 2025. Kimi k1.5: Scaling reinforcement learning with llms. <i>CoRR</i> , abs/2501.12599.	770 771 772 773 774 775 776 777
721	HuanYu Liu, Jia Li, Hao Zhu, Kechi Zhang, Yihong Dong, and Ge Li. 2025a. SATURN: sat-based rein- forcement learning to unleash language model rea- soning. <i>CoRR</i> , abs/2505.16368.	Xuezhi Wang and Denny Zhou. 2024. Chain-of-thought reasoning without prompting. In <i>NeurIPS</i> .	778 779
725	Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. 2025b. Prorl: Prolonged reinforcement learning expands rea- soning boundaries in large language models. <i>CoRR</i> , abs/2505.24864.	Jinyang Wu, Chonghua Liao, Mingkuan Feng, Shuai Zhang, Zhengqi Wen, Pengpeng Shao, Huazhe Xu, and Jianhua Tao. 2025. Thought-augmented policy optimization: Bridging external guidance and inter- nal capabilities. <i>CoRR</i> , abs/2505.15692.	780 781 782 783 784
730	Yihao Liu, Shuo Cheng Li, Lang Cao, Yuhang Xie, Mengyu Zhou, Haoyu Dong, Xiaojun Ma, Shi Han, and Dongmei Zhang. 2025c. Superrl: Reinforcement learning with supervision to boost language model reasoning. <i>CoRR</i> , abs/2506.01096.	Zhiheng Xi, Wenxiang Chen, Boyang Hong, Senjie Jin, Rui Zheng, Wei He, Yiwen Ding, Shichun Liu, Xin Guo, Junzhe Wang, Honglin Guo, Wei Shen, Xiaoran Fan, Yuhao Zhou, Shihan Dou, Xiao Wang, Xinbo Zhang, Peng Sun, Tao Gui, and 2 others. 2024. Training large language models for reasoning through reverse curriculum reinforcement learning. In <i>ICML</i> . OpenReview.net.	785 786 787 788 789 790 791 792
735	Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. https://github.com/ agentica-project/rllm . Notion Blog.	Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhi- rong Wu, and Chong Luo. 2025. Logic-rl: Unleash- ing LLM reasoning with rule-based reinforcement learning. <i>CoRR</i> , abs/2502.14768.	793 794 795 796 797
741	Lu Ma, Hao Liang, Meiyi Qiang, Lexiang Tang, Xi- aochen Ma, Zhen Hao Wong, Junbo Niu, Chengyu Shen, Runming He, Bin Cui, and Wentao Zhang. 2025. Learning what reinforcement learning can't: Interleaved online fine-tuning for hardest questions. <i>Preprint</i> , arXiv:2506.07527.	Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. 2025. Learning to reason under off-policy guidance. <i>CoRR</i> , abs/2504.14945.	798 799 800 801
747	Sanmit Narvekar, Jivko Sinapov, Matteo Leonetti, and Peter Stone. 2016. Source task creation for curricu- lum learning. In <i>AAMAS</i> , pages 566–574. ACM.	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi- heng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian	802 803 804

805 Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Ji-
806 axi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and
807 22 others. 2024. Qwen2.5 technical report. *CoRR*,
808 abs/2412.15115.

809 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,
810 Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu,
811 Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole
812 Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang,
813 Mofan Zhang, Wang Zhang, Hang Zhu, and 16 others.
814 2025. DAPO: an open-source LLM reinforcement
815 learning system at scale. *CoRR*, abs/2503.14476.

816 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai
817 Wang, Yang Yue, Shiji Song, and Gao Huang. 2025.
818 Does reinforcement learning really incentivize reason-
819 ing capacity in llms beyond the base model?
820 *CoRR*, abs/2504.13837.

821 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D.
822 Goodman. 2022. Star: Bootstrapping reasoning with
823 reasoning. In *NeurIPS*.

824 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Ke-
825 qing He, Zejun Ma, and Junxian He. 2025. Simpler1-
826 zoo: Investigating and taming zero reinforcement
827 learning for open base models in the wild. *CoRR*,
828 abs/2503.18892.

829 Rosie Zhao, Alexandru Meterez, Sham M. Kakade, Cen-
830 giz Pehlevan, Samy Jelassi, and Eran Malach. 2025.
831 Echo chamber: RL post-training amplifies behaviors
832 learned in pretraining. *CoRR*, abs/2504.07912.

Appendix

A The Pseudocode of EvoCoT Algorithm

Algorithm 1 presents the complete algorithmic workflow of **EvoCoT**. We introduce an additional `train_steps` argument to enforce that **EvoCoT** uses the same total number of training steps as all baselines, enabling a fair comparison and controlled ablation. For example, in the ablation study shown in Table 6, we set `train_steps = 500` for all variants.

B The Prompt Templates of EvoCoT

This appendix provides the prompt templates used for **EvoCoT** Stage 1: Answer-Guided CoT Generation and Stage 2: Step-Wise Curriculum Learning. Figure 5 shows the Qwen2.5 prompt template. For other models, Stage 1 templates remain the same, while Stage 2 templates follow the special token concatenation scheme in (Zeng et al., 2025). All experiments’ evaluations also use the same Stage 2 template.

C The Training and Evaluation Details of EvoCoT

Table 5: **EvoCoT** Training Hyperparameters

Parameter	Value	Parameter	Value
Advantage estimator	GRPO	Learning rate	1×10^{-6}
Train batch size	32	Mini-batch size	32
Prompt length (max)	3000	Response length (max)	5192
Samples per problem	8	Temperature	1.0
KL loss enabled	Yes	KL loss coefficient	0.0001
Shuffle dataset	No	Micro batch size	1

This appendix provides additional details on the framework and hyperparameters used for training and evaluation of **EvoCoT**. We use the Ver1 framework for training the models, which provides an efficient RL pipeline. The full list of training hyperparameters is shown in Table 5. For evaluation, we use the Qwen2.5-7B-Math framework⁴ to evaluate LLMs’ performance across various benchmarks.

All other evaluation parameters not explicitly mentioned follow the default settings of frameworks. The specific implementation code is provided in the supplementary materials.

⁴<https://github.com/QwenLM/Qwen2.5-Math>

D EvoCoT Hyperparameters Ablation Studies

In this section, we discuss the effect of the `delta_step` hyperparameter in **EvoCoT**.

In our preliminary experiments, we explicitly examined how many CoT steps should be removed at each iteration. As shown in Appendix A, Algorithm 1, line 25, we consider: for `k` in range(`n`, `-1`, `-delta_step`). Based on empirical sampling, the initial CoT-length distribution of Qwen2.5-7B is shown in Figure 4.

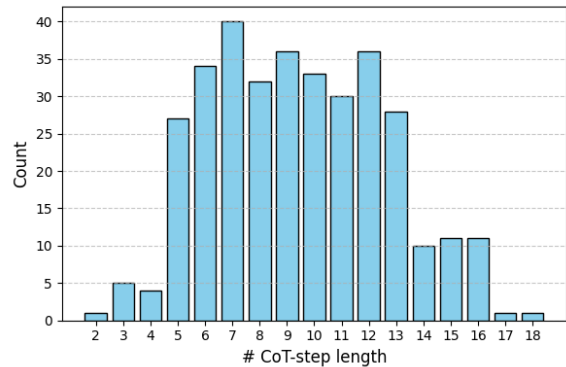


Figure 4: Distribution of CoT steps for Qwen2.5-7B on GSM8K.

All variants use the same total number of 500 RL training steps. Even if the CoT is not fully shortened, training stops at step 500. The results on Qwen2.5-7B are shown in Table 6.

delta_step	GSM8K	MATH	AIME 24	AMC 23	Minerva Math	Olympiad Bench	Avg.
1	91.4	76.5	20.0	60.0	37.1	35.9	53.5
2	90.8	77.1	16.7	52.5	38.6	38.7	52.4
3	91.6	77.2	13.3	55.0	31.6	36.3	50.8

Table 6: **EvoCoT** performance under different `delta_step` values on Qwen2.5-7B.

A larger shortening step makes each stage much harder and disrupts the smooth progression of curriculum learning. Removing one CoT step per iteration achieves the highest average score 53.5. Increasing the step size to 2 reduces the average score to 52.4, and a step size of 3 further decreases it to 50.8. Removing multiple steps at once expands the exploration space too abruptly and limits the model’s ability to adapt during RL training. Based on this observation, for weaker models such as Qwen2.5-7B, removing only one CoT step at each iteration is sufficient.

894 In future work, if applying EvoCoT to stronger
895 flagship models (e.g., DeepSeek-R1), the step size
896 can be dynamically adjusted based on the reward.
897 If the model answers correctly at each step, the step
898 size can be increased, and the number of CoT steps
899 removed in each iteration can be set to 2, 3, or even
900 more until the rollouts in a batch reach a balance
901 between positive and negative samples.

902 **E LLMs Usage**

903 In preparing this manuscript, we use LLMs to aid
904 and polish the writing. Specifically, LLMs improve
905 clarity, grammar, and phrasing, ensuring the text
906 is concise and readable. The use of LLMs **does**
907 **not** influence the technical contributions or the in-
908 terpretation of experimental findings. All content
909 polished by LLMs is carefully checked by the au-
910 thors.

Algorithm 1 EvoCoT: Self-Evolving Curriculum Learning

```
1 def EvoCoT(LLM, D, T, train_steps, delta_step=1):
2     """
3     LLM: initial language model
4     D: dataset of (Q, A) pairs
5     T: number of self-evolving iterations
6     train_steps: maximum number of training steps
7     delta_step: number of CoT steps to remove each shortening (default = 1)
8     """
9     train_cnt = 0
10    for t in range(T):
11        # Stage 1: Answer-guided CoT generation
12        C_set = []
13        for (Q, A) in D:
14            # Generate reasoning trajectory conditioned on answer
15             $\hat{C}$  = LLM.generate(Q, A)
16
17            # Verify: can the LLM derive A from answer-guided CoT?
18             $\hat{A}$  = LLM.generate(Q,  $\hat{C}$ )
19            if  $\hat{A}$  == A:
20                C_split = split_steps( $\hat{C}$ ) # split into step-wise CoT
21                C_set.append((Q, C_split, A))
22
23        # Stage 2: Step-wise curriculum training
24        for (Q, C_split, A) in C_set:
25            n = len(C_split)
26            # progressively shorten CoT from full length to 0
27            for k in range(n - 1, -1, -delta_step):
28                 $C_k$  = partial_CoT(Q, C_split, k) # keep first k steps
29                 $\tilde{C}, \tilde{A}$  = LLM.roll_out(Q,  $C_k$ ) # roll out CoT & answer
30                # Train LLM with reward based on answer correctness
31                LLM.train( $\tilde{C} + \tilde{A}$ , reward=( $\tilde{A}$  == A))
32
33            train_cnt += 1
34            if train_cnt >= train_steps:
35                return
```

Stage 1: Answer-Guided CoT Generation

Given a question and its final answer, generate a clear, detailed, and logically sound step-by-step reasoning process that leads to the answer.

Each step should be separated by two newline characters `\n\n` for clarity.

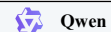
You must not contradict, challenge, or reevaluate **Correct Answer** under any circumstances.

Question: {QUESTION}

Correct Answer: {ANSWER}

Now write the reasoning. Ensure that your reasoning matches the correct answer.

Stage 2: Step-Wise Curriculum Learning.



```
<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
{input}
Please reason step by step, and put your final answer
within \boxed{}.<|im_end|>
<|im_start|>assistant
{output}
```

Figure 5: Qwen2.5 Prompt format used for EvoCoT