

# *rSIM*: Incentivizing Reasoning Capabilities of LLMs via Reinforced Strategy Injection

Anonymous ACL submission

## Abstract

Large language models (LLMs) are post-trained through reinforcement learning (RL) to evolve into Reasoning Language Models (RLMs), where the hallmark of this advanced reasoning is “aha” moments when they start to perform *strategies*, such as self-reflection and deep thinking, within chain of thoughts (CoTs). Motivated by this, this paper proposes a novel reinforced strategy injection mechanism (*rSIM*), that enables any LLM to become an RLM by employing a small planner to guide the LLM’s CoT through the adaptive injection of reasoning strategies. To achieve this, the planner (leader agent) is jointly trained with an LLM (follower agent) using multi-agent RL (MARL), based on a leader-follower framework and straightforward rule-based rewards. Experimental results show that *rSIM* enables Qwen2.5-0.5B to become an RLM and significantly outperform Qwen2.5-14B. Moreover, the planner is generalizable: it only needs to be trained once and can be applied as a plug-in to substantially improve the reasoning capabilities of existing LLMs. In addition, the planner supports continual learning across various tasks, allowing its planning abilities to gradually improve and generalize to a wider range of problems.

## 1 Introduction

Large language models (LLMs) have been enhanced with advanced reasoning capabilities, evolving into Reasoning Language Models (RLMs) (Besta et al., 2025) that solve problems through step-by-step reasoning, commonly referred to as chain-of-thought (CoT) (Wei et al., 2022). A key advancement in RLMs is their ability to integrate reasoning *strategies*, such as self-reflection, decomposition, and deliberative thinking, into the CoT process, contributing to improved problem-solving accuracy.

Existing literature (Trung et al., 2024; Havrilla et al., 2024), especially the recent Group Relative

Policy Optimization (GRPO) (Guo et al., 2025), primarily post-trains LLMs to evolve into RLMs purely through reinforcement learning (RL) algorithms. A hallmark of this evolution is the emergence of the “aha moment” where LLMs start to perform *strategies* such as self-reflection within CoTs. We empirically found that the better performance of evolved RLMs correlates with the higher appearance frequency and more appropriate positions of these strategies compared to those in LLMs. However, our findings also indicate that when LLMs, particularly smaller ones, inherently lack the capacity to perform basic reasoning strategies, RL-based post-training is unable to transform them into capable RLMs.

Therefore, this paper proposes a reinforced strategy injection mechanism (*rSIM*), which enables any LLM, in particular those as small as 0.5B, to evolve into an RLM with minimal or even no additional training. To achieve this, *rSIM* introduces only an auxiliary planner that, at each reasoning step of an LLM, adaptively selects an appropriate strategy from a predefined set such as *self-reflection*, *decomposition*, and others, and injects it into the chain of thought (CoT) to guide the next step generation. Specifically, the *rSIM* offers the following four key contributions:

- Injecting reasoning strategies adaptively into the CoT process of any LLM, including small and large ones, via a planner enables the LLM to directly gain the advanced reasoning ability like RLMs.
- Training the planner and LLM jointly as two agents under multi-agent RL (MARL) with the leader-follower algorithm (Gerstgrasser and Parkes, 2023) reinforces the planner’s ability to inject strategies.
- Planners are pluggable, meaning that a planner trained on one task can be directly integrated

with another LLM to enhance its reasoning ability on similar tasks.

- Planners support continual learning, as a planner in our *rSIM* can be continuously trained across tasks to enhance its planning ability over a broader range of problems.

Our results across seven datasets covering mathematics, multi-task reasoning, and code generation verify the benefits of *rSIM*. First, even small LLMs such as Qwen2.5-0.5B, when jointly trained with a planner (Qwen2.5-0.5B), can evolve into an RLM, achieving accuracy on par with Qwen2.5-14B on MATH (Hendrycks et al., 2021). Second, using the trained planner as a plugin enables another LLM to outperform larger models by a significant margin without any additional training. Third, a planner trained on mathematics can be continuously fine-tuned on coding tasks such as CodeAlpaca-20k (Chaudhary, 2023) to further guide an LLM like Qwen2.5-0.5B, achieving 17% higher accuracy on code generation.

## 2 Related Work

**Data distillation.** Knowledge distilled from large language models (LLMs) can be transferred to smaller models to enhance their performance (Xu et al., 2024). In particular, (Guo et al., 2025) verifies that larger LLMs can transfer their step-by-step reasoning abilities to smaller models by distilling chain-of-thought (CoT) (Wei et al., 2022) samples, where LLM-generated reasoning traces serve as additional fine-tuning data. Fine-tuning with teacher-generated CoT outputs (Magister et al., 2022; Yu et al., 2024; Dai et al., 2024), rationalizations (Li et al., 2023a), specialized reasoning skills (Liao et al., 2025), CoT and Program of Thought (PoT) (Chenglin et al., 2024), or even incorrect CoT samples (Huang et al., 2022; Hosseini et al., 2024) can significantly improve the reasoning abilities of smaller models. In this paper, using the planner from our DPR framework as a plugin to improve the reasoning of smaller models can be viewed as transferring a human-level planning chain to them.

**Reinforcement learning.** Reinforcement learning (RL) (Thrun and Littman, 2000) has been widely applied to decision-making tasks, as demonstrated by AlphaGo (Silver et al., 2016) and AlphaZero (Silver et al., 2017). RLHF (Ouyang et al., 2022) first leveraged PPO (Schul-

man et al., 2017) to align models with human preferences. ReFT (Trung et al., 2024) pioneered the use of RL as a fine-tuning paradigm to enhance LLM reasoning performance. Building on this progress, DeepSeek-R1-Zero (Guo et al., 2025) made a breakthrough by demonstrating that self-verification, reflection, and the ability to generate long CoTs in LLMs can be incentivized purely through RL, specifically using GRPO (Shao et al., 2024). This enables smaller base LLMs, such as 3B and 7B models, to be trained directly to achieve reasoning performance comparable to stronger models. However, as noted by (Havrilla et al., 2024), models are inherently constrained in their ability to explore CoT solutions beyond their existing capabilities. Consequently, weaker base models as small as 0.5B fail to benefit from RL training and lag behind in reasoning performance.

**Multi-agent LLMs.** Building on the development of using a single LLM as a planning or decision-making agent, multi-agent LLMs (Li et al., 2024) based on MARL (Lowe et al., 2017), where multiple language models collaborate, have achieved significant progress in complex problem-solving (Guo et al., 2024; Li et al., 2023b; Zhang et al., 2024). As highlighted by (Li et al., 2024), this structure has been successfully applied to practical tasks (Hong et al., 2023; Qian et al., 2023; Mandi et al., 2024). To enhance performance, SOCRATIC (Shridhar et al., 2023) trains a combination of two small distilled models to perform CoT reasoning in LLMs. In contrast to our DPR framework, SOCRATIC still relies on distilling the abilities of large models into smaller ones. Meanwhile, ReAct (Yao et al., 2023) enables LLMs to generate both reasoning traces and task-specific actions in an interleaved manner. Additionally, CORY (Ma et al., 2024) fine-tunes LLMs as two autonomous agents, a pioneer and an observer, which leads to superior performance compared to standard PPO. However, our work is the first to decouple planning from the reasoning process and to build a two-agent system under the MARL framework that enables any LLM to benefit from the trained planner agent.

## 3 Preliminary and Motivation

### 3.1 Reasoning Language Models via Reinforcement Learning

Given a question  $q$ , the reasoning language model (RLM), parameterized by  $\theta$ , generates a sequence

of thoughts, denoted as  $\mathbf{o} = [z_1, z_2, \dots, z_n]$ , where each  $z_i$  with  $i \in [1, \dots, n]$  is a textual description of the thought at the  $i$ -th reasoning step. The predicted solution  $\tilde{y}$  is extracted from  $z_n$  and compared with the ground truth  $y$ . In the reinforcement learning (RL) framework, represented as  $[\mathcal{S}, \mathcal{A}, \pi_\theta, \mathcal{R}]$ , the state  $s \in \mathcal{S}$  corresponds to the tokens generated so far, while the action  $a \in \mathcal{A}$  is the next token sampled from the policy  $a \sim \pi_\theta(a|s)$ . The corresponding reward is denoted as  $r \in \mathcal{R}$ . To optimize the policy model  $\theta$ , the RL algorithm aims to maximize the expected cumulative reward, which is formulated as follows:

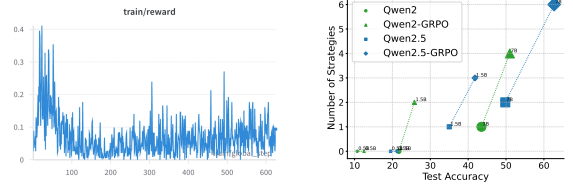
$$J(\theta) = \mathbb{E}_{s, a \sim \pi_{\theta_{old}}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)} A^{\pi_{\theta_{old}}}(s, a) \right]$$

where  $\theta_{old}$  represents the most recent policy model. The advantage function  $A(\cdot)$  estimates how much better an action is compared to the expected return. To balance bias and variance, we employ Generalized Advantage Estimation (GAE), formulated as  $A^{\pi_{\theta_{old}}}(s_t, a_t) = \sum_{l=0}^{T-t} (\gamma\lambda)^l \delta_{t+l}$  where  $\delta_t = r_t + \gamma V^{\pi_{\theta_{old}}}(s_{t+1}) - V^{\pi_{\theta_{old}}}(s_t)$ , where  $t$  denotes the token index, and  $\gamma \in [0, 1]$  is the discount factor. In general, to mitigate the discrepancy between  $\theta$  and  $\theta_{old}$ , methods such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) and Group Relative Policy Optimization (GRPO) (Shao et al., 2024) incorporate a KL penalty term.

As demonstrated by GRPO, the policy model  $\theta$  is trained using rule-based rewards, where the rewards encode human-defined rules to guide the model in adhering to these rules while improving reasoning accuracy.

### 3.2 LLMs Without Inherent Reasoning Strategies Show Limited Improvement

Using a pure RL algorithm like GRPO, a base model is optimized to achieve an “aha” moment — learning to exhibit more reasoning strategies, such as self-reflection and deep thinking, during training, leading to a sudden increase of model response length and ultimately to high accuracy. We argue that (1) when the base LLM is capable of performing reasoning strategies, reinforcement learning (RL) algorithms can optimize it to apply more strategies effectively during reasoning, and (2) there exists a positive correlation between the number of strategies and the accuracy.



(a) GRPO-based training failure on Qwen2.5-0.5B (b) Accuracy vs. number of strategies across LLMs

Figure 1: (a) Shows the absence of the “aha” moment in Qwen2.5-0.5B on the MATH dataset (Hendrycks et al., 2021). (b) Compares performance of LLMs before and after training with GRPO. By detecting key words, we count the average number of reasoning strategies, as presented in Figure 2, that are used in answering each question.

To prove this, we first define a set of the most commonly used *strategies*: *self-reflection*, *decomposition*, *deliberative thinking*, *validation*, *summarization*, *prioritization*, *continuation*, *sub-planning*, and *termination*. Following subsection A.2, we then count the number of times these strategies appear in the CoTs generated by the base LLMs and post-trained models.

As supporting evidence, we train LLMs of sizes 0.5B, 1.5B, and 7B from both Qwen2 and Qwen2.5 on the MATH dataset using the GRPO. First, as shown in Figure 1, when Qwen2.5-0.5B is used as the base model, the total reward initially increases to around 0.3 but then abruptly drops to 0 under GRPO training. Second, from Figure 1b, we observe that base models such as Qwen2-0.5B and Qwen2.5-0.5B, which lack inherent reasoning strategies (i.e., strategy count is 0), cannot be trained with RL to gain reasoning intelligence and show only limited improvement in accuracy. In contrast, models such as 1.5B and 7B, which demonstrate inherent reasoning strategies (with strategy counts greater than 0), can be further optimized through RL to apply more strategies during reasoning. More importantly, we observe that as the number of strategies increases, model accuracy also improves.

## 4 Methodology

Motivated by our observations in Subsection 3.2, we introduce reinforced strategy injection mechanism (*rSIM*), which allows a planner, implemented as a small-sized LLM, to guide another LLM by adaptively providing strategy instruction during CoTs. Through rejection, the planner can incorporate rich human-crafted knowledge and



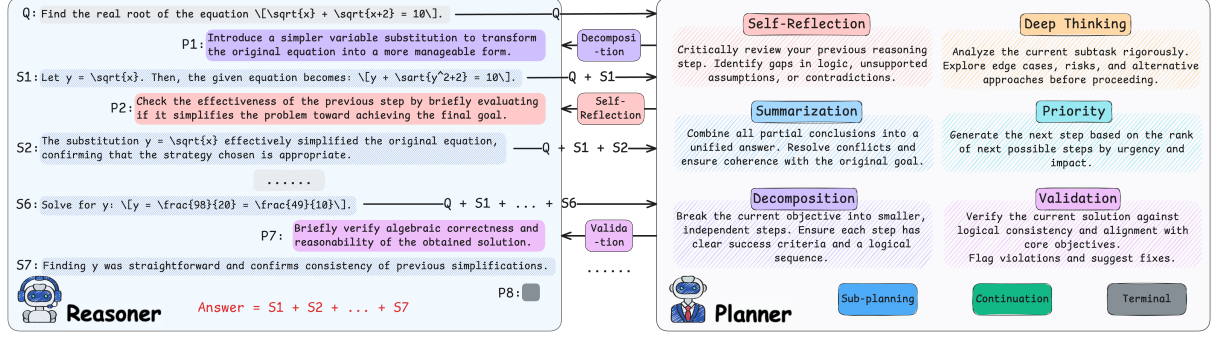


Figure 2: Illustration of the cooperative pipeline between the *rSIM* planner (leader) and the LLM (reasoner/follower). The planner receives the question and the current reasoning steps, and selects one of nine strategies to inject into the reasoning process to guide the reasoner in generating the next step. This demo is based on a question from the MATH dataset.

prior information, presented as those *strategies*, to help the LLM effectively and directly gain advanced reasoning abilities and evolve toward becoming an RLM.

#### 4.1 Training Objective of the Two-Agent Framework

The planner is designed to guide the reasoner at each reasoning step by selecting one strategy from a predefined set of nine human-designed strategies (Figure 2), which encode core reasoning strategies that LLMs should follow. These strategies, expressed as prompts, help bridge the gap in weak LLMs that inherently lack such capabilities, enabling meaningful improvement during post-training with reinforcement learning. In this framework, the base LLM serves as the reasoner, responsible for generating the next reasoning step based on the planner’s selected strategy. This collaboration is naturally modeled as a leader-follower paradigm (Gerstgrasser and Parkes, 2023) within a multi-agent RL system. The planner, acting as the leader, takes an action  $\mathbf{a}^p \sim \pi_\phi^p(\mathbf{a}^p | \mathbf{s})$ , where  $\pi_\phi^p$  is the policy parameterized by  $\phi$ ,  $\mathbf{a}^p$  is a strategy sampled from the action space  $\mathcal{A}^p$ , and  $\mathbf{s}^p \in \mathcal{S}^p$  includes the question and previous reasoning steps. The reasoner, as the follower, then takes an action  $\mathbf{a} \sim \pi_\theta(\mathbf{a} | \mathbf{s}, \mathbf{a}^p)$ , where  $\mathbf{a}$  is the next token and  $\mathbf{s}$  denotes all tokens generated thus far. Eventually, we define the planner’s reward as  $R^p = R_{acc} + R_{terminal} + R_{penalty}$  and the reasoner’s reward as  $R = R_{acc} + R_{format} + R_{follow}$ , where  $R_{terminal} = 1$  if the final plan is the ‘Terminal’ strategy, and  $-1$  otherwise;  $R_{penalty} = -$  (ratio of the most frequently selected strategy);  $R_{follow}$  is the ratio of reasoning steps that follow the given plan; and  $R_{acc}$  and

$R_{format}$  are the accuracy and format rewards as defined in GRPO (Shao et al., 2024).

As presented in Figure 2, to answer each question, the two agents interact with each other for  $n$  rounds, where  $n$  corresponds to the number of reasoning steps required to produce the final answer, leading to  $\mathbf{o}^{dpr} = [\mathbf{p}_1, \mathbf{z}_1, \mathbf{p}_2, \mathbf{z}_2, \dots, \mathbf{p}_n, \mathbf{z}_n]$ , where  $\mathbf{p}_n$  is plan selected by the planner for the  $n$ -th reasoning step. In the context of reinforcement learning, we define two advantage functions:  $A^{\pi_\phi^p}(\mathbf{s}_t, \mathbf{a}_t^p)$  for the planner and  $A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t, \mathbf{a}^p)$  for the reasoner. Here,  $t$  denotes the index of the generated token. For the planner, we assume that all tokens within a single reasoning step share the same advantage, which is set equal to the plan-level reward. Therefore, we have the following objective  $\mathcal{J}_{\mathbf{o}^{dpr}}$ :

$$\frac{1}{|\mathbf{o}^{dpr}|} \sum_{t=1}^{|\mathbf{o}^{dpr}|} \left[ \lambda \cdot \left( \frac{\pi_\phi^p(\mathbf{a}_t^p | \mathbf{s}_t)}{\pi_{\phi_{old}}^p(\mathbf{a}_t^p | \mathbf{s}_t)} \right) \cdot A^{\pi_\phi^p}(\mathbf{s}_t, \mathbf{a}_t^p) + (1 - \lambda) \cdot \left( \frac{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t, \mathbf{a}_t^p)}{\pi_{\theta_{old}}(\mathbf{a}_t | \mathbf{s}_t, \mathbf{a}_t^p)} \right) \cdot A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t, \mathbf{a}_t^p) \right]$$

where  $\mathbf{o}^{dpr} \sim (\pi_{\phi_{old}}^p, \pi_{\theta_{old}})$ , and we set  $A^{\pi_\phi^p}(\mathbf{s}_t, \mathbf{a}_t^p) = A^{\pi_\phi^p}(\mathbf{s}_i, \mathbf{a}_i^p)$  if the  $t$ -th token belongs to the reasoning step  $\mathbf{z}_i$ . This means that the reward of the  $i$ -th plan is assigned to all tokens within the reasoning step generated based on that plan.

We follow the GRPO to define the final objective as  $\mathbb{E} \left\{ \mathbf{o}_j^{dpr} \right\}_{j=1}^G \sim (\pi_{\phi_{old}}^p, \pi_{\theta_{old}}) \left[ \frac{1}{G} \sum_{j=1}^G \mathcal{J}_{\mathbf{o}_j^{dpr}} \right]$ , where  $A^{\pi_\phi^p}(\mathbf{s}_t, \mathbf{a}_t^p) = \frac{R_j^p - \text{mean}(\{R_j^p\}_{j=1}^G)}{\text{std}(\{R_j^p\}_{j=1}^G)}$ .

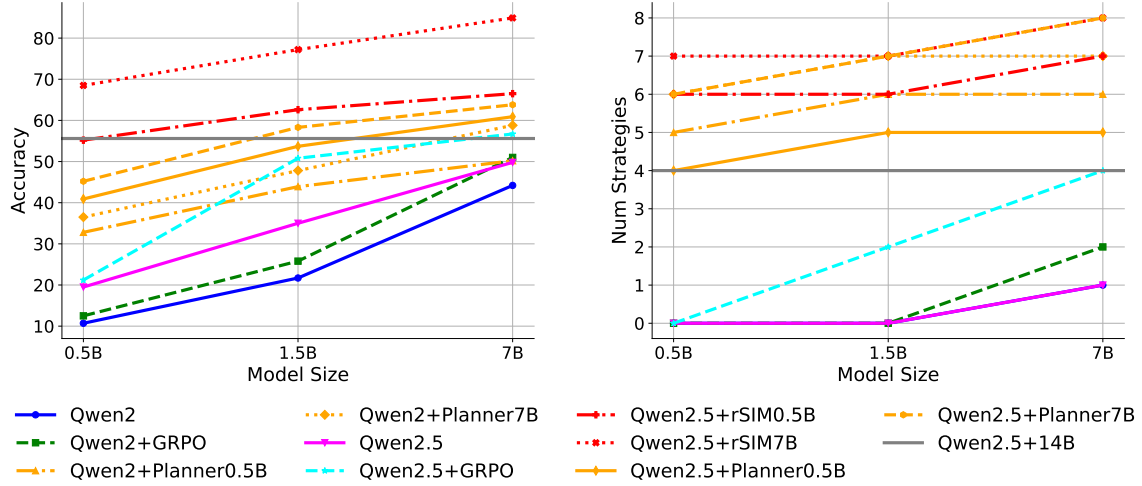


Figure 3: Accuracy of Qwen2.5 models (0.5B, 1.5B, 7B) on the MATH dataset under various settings. “+GRPO” indicates training with GRPO (Shao et al., 2024). “+rSIM” denotes joint training with a Qwen2.5 planner, with the number indicating planner size. “+Planner” refers to using trained planners (0.5B, 7B) in plugin mode, derived from the MATH dataset (see Figure 5). Qwen2.5-14B is shown as a baseline with gray horizontal lines.

## 4.2 A Two-Stage Training Scheme

We observe that simultaneously and equally optimizing the policies of both agents often leads to several issues: (1) conflicting policy updates, where the gradients of the two agents may “pull” in opposing directions; (2) credit assignment ambiguity, as it becomes unclear whether success or failure is due to the leader’s plan or the follower’s execution; and (3) competing exploration versus exploitation, where joint exploration may result in catastrophic miscoordination.

Therefore, to enable effective training, we propose a two-stage scheme in which the first stage prioritizes the policy optimization of the planner, while the second stage shifts focus to optimizing the policy of the reasoner, i.e., the base model. Specifically, we ensure a stable training process by adjusting the weighting parameter  $\lambda$  across the two stages. In the first stage, we set  $\lambda = 0.7$  to emphasize planner optimization, and in the second stage, we reduce it to  $\lambda = 0.3$  to prioritize the optimization of the reasoner. To keep the main paper concise, further implementation details are provided in the appendix.

## 5 Experiments

**Datasets.** Experiments are performed on seven datasets from the HuggingFace website: 1). MATH (Hendrycks et al., 2021), GSM8K (Cobbe et al., 2021), AIME2024 for mathematics, 2). MMLU-Pro (Wang et al., 2024), TheoremQA (Chen et al., 2023) for multi-task reasoning,

3) CodeAlpaca-20k (Chaudhary, 2023) and HumanEval for code generation.

**Training Settings.** Our experiments uses a range of LLMs, including Qwen2 models at 0.5B, 1.5B, and 7B scales, as well as Qwen2.5 models at 0.5B, 1.5B, 7B, and 14B scales. In addition, we incorporate Open-o1 and Deepseek-R1 (Guo et al., 2025) as base reasoners, which are paired with the trained planner used as a plugin. We use a batch size of 16 with gradient accumulation set to 4. For the GRPO-related hyperparameters, we configure the temperature at 0.9, the maximum prompt and completion lengths at 1024 each, G at 16, and the KL coefficient (beta) at 0.04. The learning rates for the 0.5B, 1.5B, 3B, 7B, and 14B models are  $2e-5$ ,  $1e-5$ ,  $8e-6$ ,  $5e-6$ , and  $2e-6$ , respectively. These settings remain consistent for both the planner and the reasoner. We employ the AdamW optimizer with a cosine learning rate scheduler. During evaluation, we set the planner’s temperature to 0 and the reasoner’s temperature to 0.3.

**Baselines.** We compare rSI with the recent GRPO method (Shao et al., 2024). In addition, we include the Plan-and-Solve (PS+) prompting approach (Wang et al., 2023) and Planner Prompting as two baselines. Specifically, “w/ prompting” in the experiments refers to Planner Prompting, where an LLM is directly prompted to act as the planner and provide the reasoning strategy for each step of generation.

**Metrics.** The primary evaluation metric is accuracy, defined as the percentage of correct answers. Additionally, to assess the planner’s effectiveness

Table 1: Performance of Llama series models with *rSIM* on four datasets such as TheoremQA, along with comparisons to the Plan-and-Solve (PS+) Prompting method (Wang et al., 2023) and Planner Prompting baseline. When the reasoner (in the Models column) collaborates with a planner, the planner’s size is indicated instead of using ‘No’. When a model name such as Qwen2.5 is specified in the ‘Planner’ column, we perform cross-model evaluation, where Llama serves as the reasoner and Qwen2.5 as the planner. When the ‘Planner’ column specifies ‘plug-in’, the planner is used off-the-shelf without training to guide the reasoner. Moreover, the row containing ‘w/ prompting’ indicates that we directly prompt an LLM to act as the planner. The ‘w/’ specifies the method used to enhance reasoning. The ‘-’ indicates that training did not converge, while the ‘×’ denotes missing results due to untrainable models.

Llama3.2 Models	Planner	MATH		GSM8K		MMLU-Pro		TheoremQA	
		Score	#Strategy	Score	#Strategy	Score	#Strategy	Score	#Strategy
1B w/ ZeroCoT	No	30.6	0	44.4	0	21.2	0	13.7	0
1B w/ PS+ [4]	No	28.2	0	43.7	0	19.4	0	12.2	0
1B w/ prompting	3B	27.4	7	42.6	3	16.8	8	6.6	6
1B w/ prompting	70B	33.3	5	46.9	3	22	6	14.3	5
3B w/ ZeroCoT	No	48	0	77.7	0	30.1	0	20.7	0
3B w/ PS+ [4]	No	47.5	0	77.7	0	30	0	18.6	0
3B w/ prompting	3B	46.4	7	77.1	5	28.5	7	19.9	8
3B w/ prompting	70B	55.5	7	81.8	4	31.8	5	22.8	6
1B w/ GRPO	No	-	0	-	0	×	0	×	0
1B w/ <i>rSIM</i>	1B	57	3	83.9	1	30.8	4	20.9	3
1B w/ <i>rSIM</i>	3B	61.5	4	86.3	3	33	6	25	6
1B w/ <i>rSIM</i>	Qwen2.5-1.5B	59.1	4	84.4	1	31.8	4	24.2	5
Llama3.3 70B w/ ZeroCoT	No	77	0	90.5	0	68.9	0	32.3	0
Llama3.3 70B w/ PS [4]	No	78.3	0	90.9	0	70	0	32	0
Llama3.3 70B w/ Prompting	3B	79.1	7	90.5	4	68.9	7	31.9	8
Llama3.3 70B w/ Prompting	70B	84	6	92.9	4	71.5	4	38.6	5
Llama3.3 70B	1B plug-in	83.2	3	91.7	1	71.8	6	39	5
Llama3.3 70B	3B plug-in	86.3	4	92.1	2	72.7	5	41.8	6
Llama3.3 70B	Qwen2.5-1.5B	83.7	4	92	2	72.3	6	40.7	6

as a guidance mechanism, we measure the average number of strategies applied per problem. For the baseline models, we include the candidate strategies defined in our *rSIM* as part of the prompt, allowing the models to select them when appropriate.

## 5.1 Main Results

With *rSIM*, any LLM, especially smaller ones, can be trained to convergence, achieving dramatic improvements in reasoning performance, as shown by the training curve in Figure 5 and the high problem-solving accuracy in Figure 3. Specifically, for the base Qwen2.5 models in sizes of 0.5B, 1.5B, and 7B, jointly training with a planner, either Qwen2.5-0.5B or Qwen2.5-7B, under the results in significant accuracy gains over the base models, even surpassing stronger base models. For example, all models trained with *rSIM* (tagged as “+*rSIM*” in Figure 3) outperform their GRPO-trained counterparts (“Qwen2.5+GRPO”). Moreover, the base Qwen2.5-0.5+*rSIM*0.5B model, which employs a 0.5B planner, notably outperforms the Qwen2.5-7B model. Overall, all Qwen2.5+*rSIM*7B models across different base

sizes achieve a new state-of-the-art (SOTA) performance. Moreover, as shown in Table 3, the planner introduced by our *rSIM* does not significantly increase the token cost when facilitating reasoning in a single LLM.

By comparing the left and right sub-figures of Figure 3, we observe that after training under *rSIM*, every base model learns to execute a variety of human-level reasoning strategies as defined by the nine options in Figure 2. In particular, even weaker base models, such as Qwen2.5-0.5B, which initially showed no discernible strategy, can be guided by the planner to employ six or seven strategies during reasoning. These results also reveal a positive correlation between accuracy and the number of strategies used.

More importantly, as shown in Table 1, applying *rSIM* to other types of LLMs, such as a series of Llama models, still yields consistent performance improvements due to the strategy injection provided by the planner. For example, when Llama3.2-1B is used as the reasoner and paired with a Llama3.2-1B planner, the system achieves accuracy close to that of Llama3.3-70B. This suggests that the planner in our DPR framework is

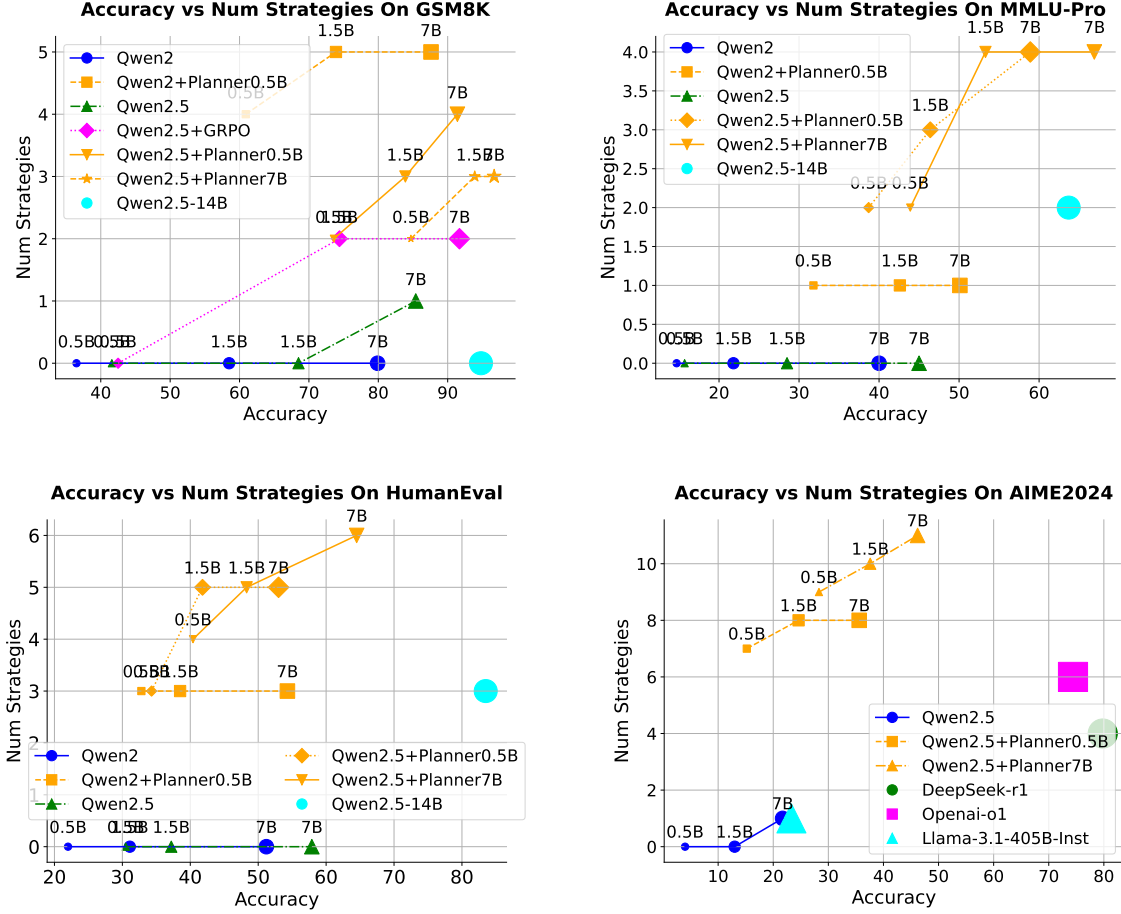


Figure 4: Illustration of accuracy for models of various sizes across different datasets, using the trained *rSIM* planner as a plugin. The terms with “+Planner” indicate that the base model collaborates with a trained planner during reasoning. The trained planners are derived from the MATH dataset, as shown in Figure 5.

transferable across different LLMs, although the degree of improvement may depend on the base model used for planning.

In addition, we present cross-model evaluation results in Table 1, where the trained planner (Qwen2.5-1.5B) is used to assist Llama3.2-1B and Llama3.3-70B by injecting step-wise strategies. The reasoning performance of the Llama models shows an improvement, indicating that the *rSIM* planner possesses cross-model generalization capability.

## 5.2 Evaluation of the Pluggable Planner

As shown in Figures 3 and 4, which report the performance of base LLMs on five datasets, the trained planner can be used as a plugin to directly introduce human-level reasoning strategies **without any additional post-training**. In the MATH dataset (Figure 3), Qwen2.5+Planner models achieve significantly higher accuracy than both the base Qwen2.5 models and those post-

trained with GRPO. A similar trend is observed across all four datasets in Figure 4: once integrated with the pre-trained planner, base models show substantial accuracy gains and often outperform larger counterparts. For instance, Qwen2.5-1.5B+Planner consistently surpasses the 7B models, and Qwen2.5-7B+Planner achieves SOTA accuracy, even outperforming more powerful models such as Llama-3.1-405B-Inst in the final sub-figure of Figure 3. More importantly, as shown by the number of strategies used during problem-solving across all datasets, the trained planner enables any LLM to incorporate effective reasoning strategies into the reasoning process. This effect is especially evident in the challenging AIME2024 dataset, where the planner enables LLMs to apply human-level strategies more than eight times per problem, resulting in a significant improvement in accuracy.

These results demonstrate that once planners are jointly trained with even weak reasoners, they



Reasoner	Planner (MATH)	MATH	GSM8K	HumanEval	MMLU_Pro	AIME2024
0.5B	0.5B	-0.8	+0	+17	+0.4	+0
	7B	+0	+1.1	+22.3	+0.9	+0
7B	0.5B	+1.2	+1.8	+24.4	+0	+0
	7B	+0	+0.6	+22.8	+0	+0

Table 2: Performance gain of the off-the-shelf planner after continued training on the coding task CodeAlpaca-20k under the *rSIM*. With the enhanced planner used as a plugin, the reasoner models (listed in the “reasoner” column: Qwen2.5-0.5B and Qwen2.5-7B) address problems across five datasets.

can be directly reused to enhance the performance of LLMs on other tasks without any additional post-training. For instance, planners such as Qwen2.5-0.5B and Qwen2.5-7B, trained with the simple Qwen2.5-0.5B reasoner on the MATH dataset, work adaptively with different LLMs across various datasets. This opens a new direction where, instead of fine-tuning LLMs for every task, we can integrate high-intelligence modules such as the planner in *rSIM* to directly boost their reasoning capabilities. In doing so, we shift the focus from optimizing a coupled planning and reasoning process to developing an advanced planner that can collaborate with any LLM.

Similarly, when using Llama series models as base LLMs, the trained planner can serve as a plug-in module to directly inject human-level reasoning strategies. As shown in Table 1, without any additional post-training, Llama3.2-1B, when guided by strategies provided by the Llama3.2-1B or Llama3.2-3B planner, achieves consistent and significant accuracy improvements across all corresponding datasets.

### 5.3 Evaluating the Continuous Learning Capability of the Planner

Table 2 shows that the planner can be continuously trained to achieve improved performance on all tasks while preserving its effectiveness on previously seen tasks. Specifically, the off-the-shelf planner, originally optimized on the MATH dataset as shown in Figure 5, is further trained with Qwen2.5-0.5B (serving as the reasoner) using the CodeAlpaca-20k dataset under the *rSIM* to obtain a coding-enhanced planner. This enhanced planner is then used as a plugin with the reasoner models (0.5B and 7B of Qwen2.5) to address questions from five diverse datasets. Table 2 reports the performance gains relative to the off-the-shelf MATH planner. Overall, aside from a small drop of 0.8 in accuracy on MATH, the performance of

the enhanced planner is either maintained or improved across all tasks. In particular, the improvement on HumanEval is substantial, with gains ranging from 17% to 24.4%. This demonstrates that continued training allows the planner to better guide coding-related reasoning using human-level strategies. Such a significant advantage highlights the practical value of the planner, as it can be continuously trained with a small reasoner on mixed-task datasets to improve its **capabilities in a cumulative way**. It is worth noting that the coding-enhanced planner does not show any accuracy gain on AIME2024.

## 6 Ablation Study

In Table 1 and Table 5, we compare the planner of *rSIM* with baseline prompting methods, including PS+ prompting (Wang et al., 2023) and direct planner prompting. From our results, we observe that these prompting-based approaches generally require a powerful language model, such as Qwen2.5 14B, to yield even limited improvements. When the model is smaller, such as LLaMA3.2 3B or Qwen2.5 0.5B, accuracy often decreases, demonstrating the limited reliability and scalability of prompting for plan generation. In addition, as shown in Table 4 and Figure 6, we examine the impact of different strategies on reasoning performance across various datasets. The results show that self-reflection is consistently important, while other strategies’ effectiveness varies across tasks.

## 7 Concluding Remarks

In this paper, we proposed the reinforced strategy injection mechanism (*rSIM*), which enables any large language model (LLM), including small ones like Qwen2.5-0.5B, to become an advanced reasoning language model (RLM). *rSIM* trains a planner with a reasoner jointly via multi-agent reinforcement learning using a leader-follower algorithm. The planner learns to select the best strategy from a set of nine human-designed options to guide the reasoner step by step. Experiments showed that *rSIM* significantly improves reasoning accuracy across multiple tasks. Importantly, the planner can be used as a plugin without additional training, enabling any LLM to gain advanced reasoning capabilities immediately. The planner also supports continual learning, enabling its reasoning guidance to improve continuously across diverse tasks.



## References

- Maciej Besta, Julia Barth, Eric Schreiber, Ales Kubicek, Afonso Catarino, Robert Gerstenberger, Piotr Nyczyk, Patrick Iff, Yueling Li, Sam Houlston, and 1 others. 2025. Reasoning language models: A blueprint. *arXiv preprint arXiv:2501.11223*.
- Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023. Theoremqa: A theorem-driven question answering dataset. In *Proc. Conference on Empirical Methods in Natural Language Processing*.
- Li Chenglin, Qianglong Chen, Liangyue Li, Caiyu Wang, Feng Tao, Yicheng Li, Zulong Chen, and Yin Zhang. 2024. Mixed distillation helps smaller language models reason better. In *Findings of the Association for Computational Linguistics*, pages 1673–1690.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Chengwei Dai, Kun Li, Wei Zhou, and Songlin Hu. 2024. Improve student’s reasoning generalizability through cascading decomposed cots distillation. In *Conference on Empirical Methods in Natural Language Processing*.
- Matthias Gerstgrasser and David C Parkes. 2023. Oracles & followers: Stackelberg equilibria in deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 11213–11236. PMLR.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. In *Proc International Joint Conference on Artificial Intelligence*, pages 8048–8057.
- Alex Havrilla, Yuqing Du, Sharath Chandra Rapparthi, Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. 2024. Teaching large language models to reason with reinforcement learning. In *International Conference on Machine Learning*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1 others. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6.
- Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. 2024. V-star: Training verifiers for self-taught reasoners. In *Conference on Language Modeling*.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *Advances in neural information processing systems*.
- Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. 2024. More agents is all you need. *Transactions on Machine Learning Research*.
- Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023a. Symbolic chain-of-thought distillation: Small models can also "think" step-by-step. In *The Association for Computational Linguistics*.
- Yuan Li, Yixuan Zhang, and Lichao Sun. 2023b. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. *arXiv preprint arXiv:2310.06500*.
- Huanxuan Liao, Shizhu He, Yao Xu, Yuanzhe Zhang, Kang Liu, and Jun Zhao. 2025. Neural-symbolic collaborative distillation: Advancing small language models for complex reasoning tasks. In *Association for the Advancement of Artificial Intelligence*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *International Conference on Learning Representations*.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30.
- Hao Ma, Tianyi Hu, Zhiqiang Pu, Liu Boyin, Xiaolin Ai, Yanyan Liang, and Min Chen. 2024. Coevolving with the other you: Fine-tuning llm with sequential cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 37:15497–15525.

678	Lucie Charlotte Magister, Jonathan Mallinson, Jakub	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu,	731
679	Adamek, Eric Malmi, and Aliaksei Severyn. 2022.	Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.	732
680	Teaching small language models to reason. In <i>The</i>	2023. Plan-and-solve prompting: Improving zero-	733
681	<i>Association for Computational Linguistics</i> .	shot chain-of-thought reasoning by large language	734
		models. In <i>Proc. Annual Meeting of the Associa-</i>	735
682	Zhao Mandi, Shreeya Jain, and Shuran Song. 2024.	<i>tion for Computational Linguistics</i> , volume 1, pages	736
683	Roco: Dialectic multi-robot collaboration with large	2609–2634.	737
684	language models. In <i>IEEE International Conference</i>		
685	<i>on Robotics and Automation</i> , pages 286–299.		
		Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,	738
686	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Abhranil Chandra, Shiguang Guo, Weiming Ren,	739
687	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others.	740
688	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	2024. Mmlu-pro: A more robust and challenging	741
689	others. 2022. Training language models to follow in-	multi-task language understanding benchmark. In	742
690	structions with human feedback. <i>Advances in neural</i>	<i>Conference on Neural Information Processing Sys-</i>	743
691	<i>information processing systems</i> , 35:27730–27744.	<i>tems Datasets and Benchmarks Track</i> .	744
692	Chen Qian, Xin Cong, Cheng Yang, Weize Chen,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	745
693	Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	746
694	Sun. 2023. Communicative agents for software de-	and 1 others. 2022. Chain-of-thought prompting	747
695	velopment. <i>arXiv preprint arXiv:2307.07924</i> , 6(3).	elicits reasoning in large language models. <i>Ad-</i>	748
		<i>vances in neural information processing systems</i> ,	749
		35:24824–24837.	750
696	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec	Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen,	751
697	Radford, and Oleg Klimov. 2017. Proximal	Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao,	752
698	policy optimization algorithms. <i>arXiv preprint</i>	and Tianyi Zhou. 2024. A survey on knowledge dis-	753
699	<i>arXiv:1707.06347</i> .	stillation of large language models. <i>arXiv preprint</i>	754
		<i>arXiv:2402.13116</i> .	755
700	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,		
701	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	756
702	Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-	Shafraan, Karthik Narasimhan, and Yuan Cao. 2023.	757
703	math: Pushing the limits of mathematical reason-	React: Synergizing reasoning and acting in language	758
704	ing in open language models. <i>arXiv preprint</i>	models. In <i>International Conference on Learning</i>	759
705	<i>arXiv:2402.03300</i> .	<i>Representations (ICLR)</i> .	760
706	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya	Ping Yu, Jing Xu, Jason Weston, and Ilia Kulikov. 2024.	761
707	Sachan. 2023. Distilling reasoning capabilities into	Distilling system 2 into system 1. In <i>Workshop on</i>	762
708	smaller language models. In <i>Findings of the Asso-</i>	<i>Neural Information Processing Systems</i> .	763
709	<i>ciation for Computational Linguistics</i> , pages 7059–		
710	7073.	Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfis-	764
		ter, Rui Zhang, and Sercan Arik. 2024. Chain	765
711	David Silver, Aja Huang, Chris J Maddison, Arthur	of agents: Large language models collaborating on	766
712	Guez, Laurent Sifre, George Van Den Driessche, Ju-	long-context tasks. <i>Advances in Neural Information</i>	767
713	lian Schrittwieser, Ioannis Antonoglou, Veda Pan-	<i>Processing Systems</i> , 37:132208–132237.	768
714	neershelvam, Marc Lanctot, and 1 others. 2016.		
715	Mastering the game of go with deep neural networks		
716	and tree search. <i>nature</i> , 529(7587):484–489.		
717	David Silver, Thomas Hubert, Julian Schrittwieser,		
718	Ioannis Antonoglou, Matthew Lai, Arthur Guez,		
719	Marc Lanctot, Laurent Sifre, Dharshan Kumaran,		
720	Thore Graepel, and 1 others. 2017. Mastering chess		
721	and shogi by self-play with a general reinforcement		
722	learning algorithm. <i>nature</i> .		
723	Sebastian Thrun and Michael L Littman. 2000. Rein-		
724	forcement learning: An introduction. <i>AI Magazine</i> ,		
725	21(1):103–103.		
726	Luong Trung, Xinbo Zhang, Zhanming Jie, Peng Sun,		
727	Xiaoran Jin, and Hang Li. 2024. Reft: Reasoning		
728	with reinforced fine-tuning. In <i>Proc Annual Meet-</i>		
729	<i>ing of the Association for Computational Linguistics</i> ,		
730	pages 7601–7614.		

## A Implementation Details

This section presents the design, training, and evaluation details of our proposed *rSIM*.

### A.1 Finetuning the Reasoner LLM toward Generating Step-wise Reasoning

We first ensure that the reasoner LLM generates the reasoning process in a step-by-step format. To accomplish this, we adopt the mechanism proposed by (Lightman et al., 2023), finetuning the model to use ‘\n\n’ to clearly separate each reasoning step in the problem-solving process. Specifically, we construct a prompt consisting of five chain-of-thought examples formatted such that each step is separated by ‘\n\n’. Using this prompt, we let the reasoner generate reasoning processes randomly on the MATH dataset. To avoid introducing dataset-specific solution information into the model, we select only those generated samples that conform to the desired format but contain incorrect solutions. Ultimately, we obtain 1,000 such samples and fine-tune the reasoner on these samples for one epoch. This process ensures that the reasoning generated by the model consistently follows the intended step-by-step formatting.

### A.2 Counting the Number of Strategies

To determine how many strategies are used by the LLM during reasoning as shown in Fig. 1(B), Fig. 4, and Fig. 5 we adopt a keyword-matching approach based on strategy names and their synonyms. Specifically, we construct a list of candidate keywords for each of the seven main strategies: *self-reflection*, *decomposition*, *deep thinking*, *validation*, *summarization*, *prioritization*, and *sub-planning*, as follows:

- Self-Reflection: review, revisit, reflect, reevaluate, rethink, reexamine, reassess, reconsider, analyze, assess, validate, critique, inspect, examine, audit, diagnose, cross-check.
- Decomposition: decompose, break down, divide, split, separate, segment, partition, dissect, analyze, unfold, unwrap, reduce, map out, organize, structure.
- Deep Thinking: contemplate, deliberate, reflect, ponder, mull over, reason, deduce, infer, evaluate, scrutinize, meditate, analyze, consider, investigate, explore.

- Validation: validate, verify, confirm, check, test, justify, prove, cross-check, ensure, affirm, support, substantiate, corroborate, authenticate, evaluate.
- Summarization: summarize, recap, restate, paraphrase, condense, outline, highlight, abstract, generalize, simplify, extract, distill, encapsulate, conclude, report.
- Prioritization: prioritize, rank, order, select, choose, emphasize, highlight, focus on, weigh, assess, sort, filter, arrange, allocate, favor.
- Sub-planning: plan, outline, design, strategize, organize, arrange, map out, formulate, structure, prepare, coordinate, blueprint, set up.

Therefore, once a reasoning step is completed, indicated by the generation of ‘\n\n’, we identify keywords within that sequence. If a match is found, we increment the count for the corresponding strategy by one.

### A.3 Structure of the Planner

The planner is a simple decoder-only LLM equipped with an action head a linear layer that outputs a 9-dimensional vector corresponding to the number of available strategies. The planner takes as input the question and the current reasoning steps. The hidden state of the final token in the sequence is fed into the action head to produce a strategy priority vector. The strategy with the highest score is then selected to guide the reasoner’s next reasoning step. For example, when using Qwen2.5-0.5 as the planner, the action head is implemented as a fully connected layer with shape  $896 \times 9$ . In this case, the hidden vector of the final token in the input sequence serves as the input to the action head for strategy selection.

### A.4 Interactive Sampling Mechanism

Different from GRPO (Shao et al., 2024), which generates  $G$  samples for each question by directly forwarding the question through the LLM  $G$  times, our multi-agent framework adopts an interactive sampling mechanism in which the planner agent and the reasoner agent of *rSIM* interact throughout the reasoning process to generate each sample. We follow the notation introduced in Subsections 3.1 and 4.1, and thus present the detailed procedure in Algorithm Table 1. It is important to note that when the planner collaborates with an LLM

---

**Algorithm 1: Interactive Sampling Mechanism**

---

**Input:** Question  $q$ , the reasoner policy  $\pi_\theta$  with its hyperparameters, such as the temperature, the planner policy  $\pi_\phi^p$  with its hyperparameters, and the number of generations  $G$ .

**Output:** Generated  $G$  samples.

```
1 Sample  $G$  first strategies from the planner
    $\{a_{1,g}^p\}_{g=1}^G \sim \pi_\phi^p(\cdot|q)$ 
2 Begin parallel sampling:
3 for any  $a_1^p \in \{a_{1,g}^p\}_{g=1}^G$  do
4    $a_{trace}^p \leftarrow a_1^p$ ,  $\mathbf{o}^{dpr} = []$ ,  $\mathbf{a}^p = [a_1^p]$ ,
    $n \leftarrow 1$ 
5   while  $a_{trace}^p$  is not Terminal do
6     Autoregressively Decode with
        $\mathbf{a} \sim \pi_\theta(\cdot|q, \mathbf{o}^{dpr}, a_{trace}^p)$  until
       ‘\n\n’ is generated
7     Collect the sequence of decoded  $\mathbf{a}$ 
       as the reasoning step  $\mathbf{p}_n, \mathbf{s}_n$ 
8     Append  $\mathbf{o}^{dpr} \leftarrow [\mathbf{o}^{dpr}, \mathbf{p}_n, \mathbf{s}_n]$ 
9     Select strategy action
        $a^p \sim \pi_\phi^p(\cdot|q, \mathbf{o}^{dpr})$ 
10    Set  $a_{trace}^p \leftarrow a^p$ , Append
        $\mathbf{a}^p \leftarrow [\mathbf{a}^p, a_{trace}^p]$ 
11    Set  $n \leftarrow n + 1$ 
12  end
13 end
14 Return  $\{\mathbf{o}_{1...G}^{dpr}\}, \{\mathbf{a}_{1...G}^p\}$ 
```

---

---

**Algorithm 2: Policy Optimization of Multi-Agent Framework**

---

**Input:**  $q$ , reasoner policy  $\pi_\theta$ , planner policy  $\pi_\phi^p$ , dataset  $\mathcal{D}$ ,  $N$ .

**Output:** Optimized  $\theta, \phi$ .

```
1 for  $e = 1, \dots, E$  do
2   Set reference models  $\pi_{\theta_{ref}} \leftarrow \pi_\theta$ ,
    $\pi_{\phi_{ref}} \leftarrow \pi_\phi$ 
3    $\lambda \leftarrow 0.7$ 
4   for step  $1, \dots, M$  do
5     Sample a batch samples  $\mathcal{D}_b$  from  $\mathcal{D}$ 
6     Update old policy models
        $\pi_{\theta_{old}} \leftarrow \pi_\theta$ ,  $\pi_{\phi_{old}} \leftarrow \pi_\phi$ 
7     Perform Interactive Sampling Mechanism to generate
        $\{\mathbf{o}_{1...G}^{dpr}\}, \{\mathbf{a}_{1...G}^p\}$  for each
       question in  $\mathcal{D}_b$ 
8     Compute rule-based rewards
        $\{R_j^p\}_{j=1}^G$  and  $\{R_j\}_{j=1}^G$  for each
        $\mathbf{o}_j^{dpr}, \mathbf{a}_j^p$ 
9     Compute advantage  $A^{\pi_\phi^p}(\mathbf{s}_t, \mathbf{a}_t^p)$  for
       the  $t$ -th token of each  $j \in G$ 
       through group relative advantage
       estimation.
10     $\lambda \leftarrow 0.3$  iff.  $\text{step} < N$ 
11    Update policy models  $\pi_\theta, \pi_\phi^p$  by
       maximizing the training objective
       (Equation 1)
12  end
13 end
```

---

for problem solving, the step-by-step reasoning is generated following the same procedure.

### A.5 Training and Evaluation Details

Algorithm Table 2 presents how to jointly optimize the policies of the two agents. During the *Interactive Sampling Mechanism*, the temperatures for the planner and reasoner are set to be 0.9. We set the maximum prompt length to 2,048 tokens and the maximum generation length to 1,024 tokens for all models. **It is worth noting that in the training objective (Equation 2), we include only  $\{\mathbf{o}_{1...G}^{dpr}\}$  for simplicity, although the full formulation should involve both  $\{\mathbf{o}_{1...G}^{dpr}\}, \{\mathbf{a}_{1...G}^p\}$ . We should note that for tokens generated guided by the action of the planner, this action reward is assigned to each to-**

**ken’s reward.**

**Evaluation setup.** We set the maximum prompt length to 2,048 tokens and the maximum generation length to 1,024 tokens for all models. During evaluation, we adopt the zero-shot setting with a temperature of 0, and report pass@1 (accuracy) using stochastic decoding. When using the planner as a plugin for evaluation, we strictly follow the *Interactive Sampling Mechanism* described in Algorithm Table 1.

### B Additional Results

As shown in Fig. 5, we train the Qwen2.5-0.5B model, initially lacking high-level reasoning capabilities, using *rSIM* with planners implemented by both Qwen2.5-0.5B and Qwen2.5-7B models. With the two-stage training scheme, total re-



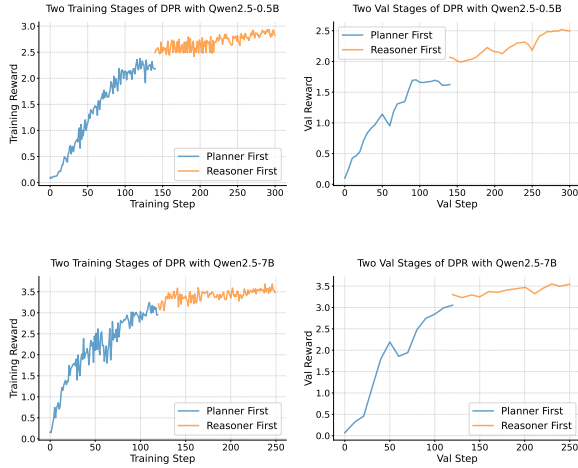


Figure 5: Training and evaluation (Eval) curves of the *rSIM* on the MATH dataset, using either the Qwen2.5-0.5B or Qwen2.5-7B model as the planner, paired with the Qwen2.5-0.5B model as the reasoner.

Table 3: Evaluation of the average generation token cost across different methods on MATH and TheoremQA. We report both the average and standard deviation (mean  $\pm$  std) of the total tokens used per question, including tokens used for prompting the LLMs and those generated by the models. This evaluation is conducted on two challenging datasets: MATH and TheoremQA.

Methods	Planner	MATH	TheoremQA
0.5B w/ ZeroCoT	No	221.6 $\pm$ 172.5	250.3 $\pm$ 110.2
14B w/ ZeroCoT	No	261.8 $\pm$ 192.2	308.7 $\pm$ 137.5
0.5B w/ PS+ [4]	No	327.5 $\pm$ 176.7	367.5 $\pm$ 153.6
0.5B w/ Prompting	14B	815.2 $\pm$ 356.7	993.4 $\pm$ 390.5
0.5B w/ <i>rSIM</i>	7B	780.3 $\pm$ 210.9	800.7 $\pm$ 230.2
7B w/ ZeroCoT	No	246.9 $\pm$ 189.5	291.2 $\pm$ 160.8
7B w/ PS [4]	No	357.2 $\pm$ 200.9	390.6 $\pm$ 190
7B w/ Prompting	14B	934.5 $\pm$ 390.2	1103 $\pm$ 487.5
7B w/ <i>rSIM</i>	7B	900.6 $\pm$ 350.8	970.2 $\pm$ 427.1

wards steadily increase until convergence, stabilizing around 2.8 for the 0.5B planner and 3.5 for the 7B planner. Crucially, the planner and reasoner policy models effectively learn to generate and follow step-wise strategies, respectively, resulting in improved reasoning performance. Additionally, the 7B planner achieves higher rewards more rapidly than the 0.5B planner. Consistent training and evaluation trends confirm effective policy optimization. Therefore, we conclude that **by decoupling planning from reasoning, we can introduce human priors in the form of integrated reasoning strategies into any LLM that lacks them, significantly enhancing reasoning intelligence through our multi-agent RL framework.**

## B.1 Additional Ablation Study

Table 4 presents the importance of different strategies in improving the reasoning accuracy of Qwen2.5-7B planner in different datasets. Due to time constraints, we do not re-adjust the strategy set for fine-tuning the LLMs across various tasks, nor do we directly compute statistical metrics for each strategy. Instead, we evaluate the importance of each strategy by iteratively masking out one strategy at a time during the evaluation phase. This allows us to observe how the performance is affected when a specific strategy is not considered by the planner. Specifically, when a strategy is masked out, the planner selects the alternative strategy with the second-highest score.

## C Limitations

The limitations of our proposed *rSIM* can be categorized into three main aspects: (1) the human-defined action space of the planner is not continuously optimizable and (2) the planner exhibits an imbalanced preference over strategies.

First, the planner’s action space is defined solely based on human understanding of the task. As a result, its effectiveness, diversity, and generalization depend heavily on expert priors and cannot be enhanced by the planner itself during reasoning. Without broader exploration during training particularly in our multi-agent reinforcement learning setting the *rSIM*’s performance may significantly degrade when necessary strategies are absent from the planner’s action space. Furthermore, since each strategy is encoded as a short prompt describing the guidance, the quality and robustness of these prompts directly affect the reasoner’s reasoning generation. Misleading, ambiguous, or overly narrow prompts can harm the reasoning process. Unfortunately, such issues are common in human-crafted prompts, including biases and lack of applicability across tasks. More importantly, we currently lack a mechanism to dynamically update or expand the planner’s action space. Although this challenge is shared by many reinforcement learning-based approaches where any change to the action space invalidates the underlying policy it is particularly critical in our context. Enhancing the strategy space is necessary to improve the generality of LLM-based reasoning frameworks, which is essential for their practical deployment at scale.

Second, we count the total number of strategies

Table 4: Evaluation of the importance of different strategies on reasoning performance when Qwen2.5-0.5 serves as the reasoner and Qwen2.5-7B acts as the planner during inference. In each row labeled with a strategy name, we remove the corresponding strategy from the planner’s option set to evaluate its impact on reasoning performance.

Dataset/Strategy	full	self-rection	deep thinking	decomposition	summarization	validation	prioritization	sub-planning
MATH	45.2	36.9	41.3	42.8	43.7	44.3	44.6	44.8
HumanEval	40.2	32.3	39	35.4	38.4	39	39.6	35.4
MMLU Pro	43.9	31.3	38.7	34.5	39.6	41.8	42.7	44.1
TheoremQA	38.7	30	35.3	32.8	38.4	38.6	37.1	36.9

Table 5: Performance of our *rSIM* on datasets such as TheoremQA, along with comparisons to the Plan-and-Solve (PS) Prompting method (Wang et al., 2023) and Planner Prompting baseline. We use Qwen2.5 models in sizes 0.5B, 7B, and 14B. The format and experimental setting of this table are consistent with those in Table 1.

Methods	Planner	MATH		MMLU-Pro		TheoremQA	
		Score	#Strategy	Score	#Strategy	Score	#Strategy
0.5B w/ ZeroCoT	No	19.5	0	15.7	0	9.5	0
0.5B w/ PS+ [4]	No	17.2	0	13	0	8	0
0.5B w/ Prompting	7B	21.6	7	15.9	8	9.6	8
0.5B w/ Prompting	14B	26.3	6	18.6	6	10	5
14B w/ ZeroCoT	No	55.6	0	51.2	0	43	0
14B w/ PS+ [4]	No	57.1	0	52.3	0	43.4	0
14B w/ Prompting	7B	56.8	6	53.5	6	43.8	7
14B w/ Prompting	14B	60	5	57	5	47.4	5
0.5B w/ <i>rSIM</i>	0.5B	40.9	4	38.7	2	34.1	3
0.5B w/ <i>rSIM</i>	7B	45.2	6	43.9	2	38.7	4
7B w/ ZeroCoT	No	49.8	0	40	0	36	0
7B w/ PS [4]	No	49.6	0	39.2	0	34.9	0
7B w/ Prompting	7B	51	6	40.2	8	36.4	7
7B w/ Prompting	14B	56.8	5	46.7	6	41.8	5
7B w/ <i>rSIM</i>	0.5B	60.9	5	58.9	4	48.3	6
7B w/ <i>rSIM</i>	7B	63.8	8	66.9	4	53	5

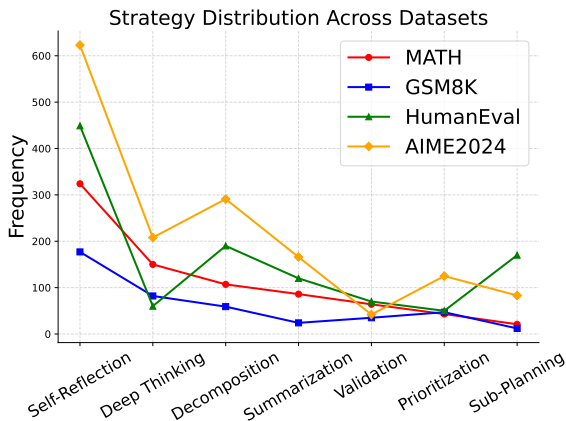


Figure 6: Illustration of how many strategies are used in solving 100 questions from MATH, GSM8K, HumanEval, and AIME2024.

used by both the planner (Qwen2.5-0.5B) and the reasoner (Qwen2.5-0.5B) during problem-solving across different datasets. As shown in Fig. 6, the planner consistently emphasizes the use of *Self-Reflection* across all four datasets, while placing less focus on *Validation*, *Prioritization*, and *Sub-Planning*. Notably, in the MATH dataset, the second most frequently used strategy is *Decomposition*, which is also true for the coding task HumanEval. In addition, for the complex mathematical problems in AIME2024, *Decomposition* again shows a high usage rate, indicating that breaking down the problem during reasoning is crucial for arriving at the correct solution. Therefore, since strategies are not used in a balanced manner during the problem-solving of questions across different tasks, the planner may be unable to comprehensively explore diverse solution paths to reach a reliable answer.