

# Enhancing Reasoning Abilities of Small LLMs with Cognitive Alignment

Anonymous ACL submission

## Abstract

The reasoning capabilities of large language reasoning models (LRMs), such as OpenAI’s o1 and DeepSeek-R1, have seen substantial advancements through deep thinking. However, these enhancements come with significant resource demands, underscoring the need for training effective small reasoning models. A critical challenge is that small models possess different reasoning capacities and cognitive trajectories compared with their larger counterparts. Hence, directly distilling chain-of-thought (CoT) results from large LRMs to smaller ones can sometimes be ineffective and often requires a substantial amount of annotated data. In this paper, we first introduce a novel Critique-Rethink-Verify (CRV) system, designed for training smaller yet powerful LRMs. Our CRV system consists of multiple LLM agents, each specializing in unique abilities: (i) critiquing the CoT qualities according to the cognitive capabilities of smaller models, (ii) rethinking and refining these CoTs based on the critiques, and (iii) verifying the correctness of the refined results. Based on the CRV system, we further propose the Cognitive Preference Optimization (CogPO) algorithm to continuously enhance the reasoning abilities of smaller models by aligning their reasoning processes with their cognitive capacities. Comprehensive evaluations on challenging reasoning benchmarks demonstrate the efficacy of our CRV+CogPO framework, which outperforms other methods by a large margin.<sup>1</sup>

## 1 Introduction

The remarkable progress in language reasoning models (LRMs) has revolutionized NLP (Zhao et al., 2023). Recently, leading models such as OpenAI’s o1<sup>2</sup> and DeepSeek-R1 (DeepSeek-AI,

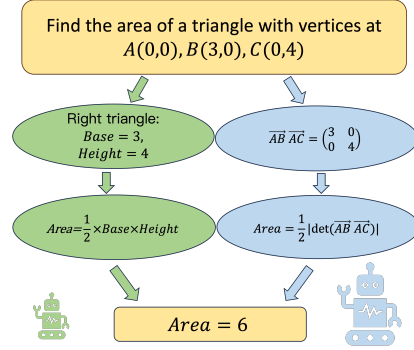


Figure 1: A motivation example. Large models (right) apply vector-based algebraic abstraction to solve the problem, while small models (left) employ simple formulaic geometric decomposition. This trajectory mismatch underscores the inefficacy of direct CoT distillation across models with substantial capacity gaps.

2025) have leveraged slow thinking to solve complex tasks. Despite their impressive capabilities, the scale of these models results in substantial computational demands. Consequently, there is a growing need to train reasoning models with fewer parameters.

A straightforward approach to address this challenge is the direct distillation of Chain-of-Thought (CoT) outputs (Wei et al., 2022a) or other deep thoughts (such as Tree-of-Thought (Yao et al., 2023b)) from larger LRMs to smaller ones. This technique is widely applied to improve the capacities of smaller LRMs (Hsieh et al., 2023; Shridhar et al., 2022; Li et al., 2023; Yue et al., 2024). However, smaller models<sup>3</sup> inherently possess different reasoning capacities and cognitive trajectories when solving problems compared to their larger counterparts, as illustrated in Figure 1. Similar findings have also been presented in (Li et al., 2022; Zhang et al., 2024; Hu et al., 2024; Li et al., 2024). This phenomenon indicates that direct distillation of CoTs from larger models can sometimes

<sup>1</sup>Source codes, datasets and models will be released upon paper acceptance.

<sup>2</sup><https://openai.com/o1/>

<sup>3</sup>In this work, we regard smaller LLMs as decoder-only language models typically with fewer than 10B parameters.

be ineffective due to the large capacity gap. Thus, a natural question arises: *How can we improve the reasoning abilities of smaller LRMs to align with their own cognitive capacity?*

In this paper, we introduce “Critique-Rethink-Verify” (CRV), a novel system to enhance the reasoning capabilities of smaller models. CRV leverages multiple LLM agents, each with specialized functions and working in synergy. These functions include (i) critiquing the CoT by considering the cognitive limits of smaller LRMs, (ii) rethinking and refining these CoTs, integrating the feedback received from the previous critiques, and (iii) verifying the accuracy and validity of the refined reasoning paths. Extending the direct preference optimization (DPO) technique (Rafailov et al., 2023), we further propose the cognitive preference optimization (CogPO) algorithm to align the reasoning process with the cognitive capacities of smaller LRMs on the basis of CRV system. Ultimately, the reasoning performance of smaller models can be improved effectively.

In the experiments, the effectiveness of our approach is evaluated on several challenging reasoning benchmarks that are difficult for models with limited parameter sizes, such as AIME 2024, MATH-500 (Lightman et al., 2023), GPQA-Diamond (Rein et al., 2023), and LiveCodeBench. The results indicate that the small LRMs trained using the CRV+CogPO framework achieve outstanding reasoning performance. In summary, we make the following major contributions:

- We present the CRV system for training small yet powerful LRMs, leveraging multiple LLM agents, each specializing in unique tasks.
- We propose the CogPO algorithm that continuously enhances the reasoning abilities of small models by aligning their reasoning processes with their cognitive capacities.
- Evaluations on challenging benchmarks demonstrate that the CRV+CogPO framework significantly improves the reasoning performance of small models, outperforming other popular training methods.

## 2 Related Work

### 2.1 Prompting LLMs to Reason

Prompting strategies to improve reasoning in LLMs have become a critical focus. Initial studies showed

that LLMs could perform basic reasoning tasks using meticulously crafted prompts, such as linguistic analysis (Chen et al., 2021) and commonsense inference (Latcinnik and Berant, 2020; Schwartz et al., 2020). To name a few, Chain-of-Thought (CoT) (Wei et al., 2022b) prompting explicitly guides LLMs through step-by-step reasoning, enabling them to decompose complex problems into manageable intermediate reasoning steps. Tree-of-Thought (ToT) (Yao et al., 2023a) prompting introduces a hierarchical structure to reasoning trajectories, allowing models to explore multiple solution paths. Furthermore, self-refine (Shinn et al., 2023; Madaan et al., 2023) prompting incorporates verification checkpoints, where models validate intermediate results before advancing.

### 2.2 Reasoning LLMs

With the advancement of LLMs, model capabilities have steadily improved (Chen and Varoquaux, 2024; Bansal et al., 2024). Models with approximately 7B to 14B parameters show remarkable performance, and their fine-tuning costs have become increasingly feasible. This has led to the emergence of specialized small models tailored for mathematical and code-related reasoning tasks such as Qwen-Math<sup>4</sup>, Qwen-Coder<sup>5</sup>, and Macro-ol (Zhao et al., 2024).

Recent studies (Shridhar et al., 2023; Yan et al., 2023; Liang et al., 2024; Yuan et al., 2024) have investigated fine-tuning methods to enhance the reasoning abilities of smaller models. By utilizing intermediate reasoning steps, LLMs can iteratively refine their outputs (Jiang et al., 2024; Wang et al., 2024; Chen et al., 2025). This methodology facilitates the development of small reasoning models, particularly following the release of stronger reasoning models such as DeepSeek-R1 (DeepSeek-AI, 2025) and QwQ-32B<sup>6</sup>.

### 2.3 Alignment Training

To effectively train LLMs, a reinforcement learning stage is typically employed after the supervised fine-tuning (SFT) phase, which serves to improve the model’s alignment towards certain objectives. Reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) has shown effectiveness in aligning LLMs with human feedback. A potential drawback of RLHF is the ex-

<sup>4</sup><https://qwenlm.github.io/blog/qwen2.5-math/>

<sup>5</sup><https://qwenlm.github.io/blog/qwen2.5-coder-family/>

<sup>6</sup><https://qwenlm.github.io/blog/qwq-32b/>

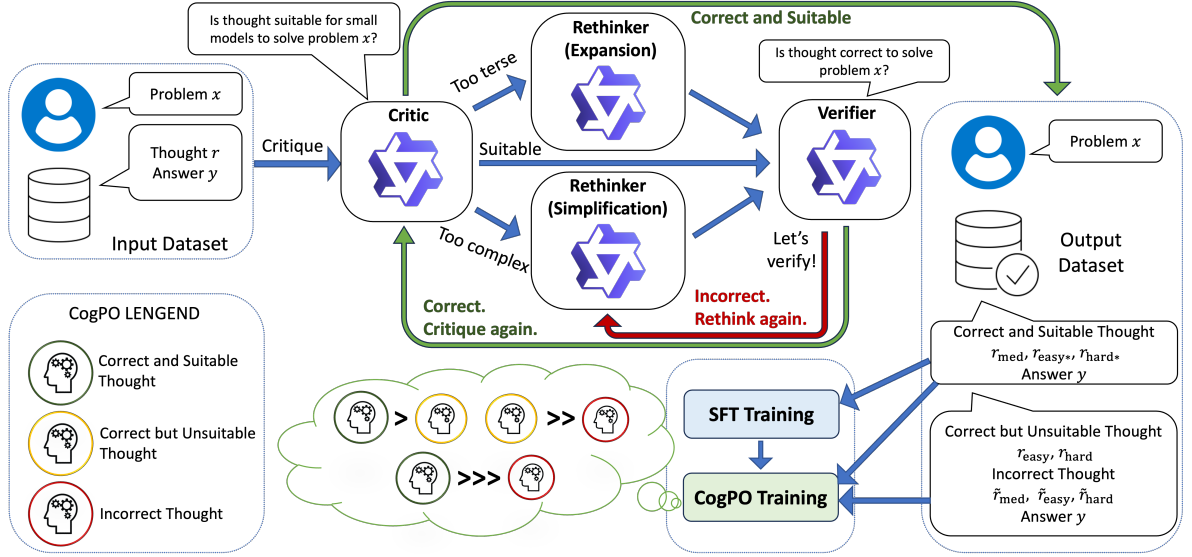


Figure 2: Overview of our CRV+CogPO framework, consisting of two synergistic phases: (1) SFT training with cognitively aligned data generated by CRV system, and (2) CogPO: dynamic  $\beta$  adjustment preference optimization training using cognitive reasoning pairs with different quality gaps. **Disclaimer:** We use the Qwen logo as our backbones; however, any LLMs with sufficient capabilities can serve as the agents as well.

licit need for a reward model and the unstable RL training process. Direct preference optimization (DPO) (Rafailov et al., 2023) trains LLMs based on chosen and rejected responses. Since the introduction of DPO, several approaches have been proposed to enhance its efficacy and efficiency. For example, CPO (Xu et al., 2024) extends DPO to avoid generating adequate but not perfect machine translations. SimPO (Meng et al., 2024) simplifies DPO by eliminating the reference model. KTO (Ethayarajh et al., 2024) and NCA (Chen et al., 2024) develop novel optimization goals that leverage unpaired data for model alignment. Furthermore, SPPO (Wu et al., 2024b) employs on-policy sampling to generate preference data, outperforming off-policy DPO methods. In our work, we extend DPO to align reasoning abilities with the cognitive limits of small LLMs.

### 3 Proposed Approach

#### 3.1 Overall Framework

Our framework consists of two synergistic phases: (1) SFT with cognitively aligned data generated by CRV system, and (2) CogPO with dynamic  $\beta$  adjustment. As illustrated in Figure 2, the CRV system first refines data tailored to the cognitive capacity of smaller LLMs for SFT training, and CogPO further aligns reasoning preferences through suitability-aware optimization using pairs with different quality gaps. This design ensures that the model initially acquires capacity-matched

reasoning patterns, followed by the refinement of its decision boundaries through gap-sensitive learning.<sup>7</sup>

#### 3.2 The CRV System

The CRV system employs LLM agents to construct the SFT dataset aligned with the cognitive limits of smaller models to be trained. The input to CRV system is an initial training set  $\mathcal{D}_{\text{SFT}} = \{(x, y, r_{\text{orig}})\}$ , where the three elements denote the problem, the correct answer, and the original reasoning process generated by any large LLMs (e.g., DeepSeek-R1), which has been validated as correct. The following provides descriptions of each agent in the CRV system.

##### 3.2.1 Critic

An LLM agent first evaluates the appropriateness of reasoning processes for the target small LLM (denoted as  $\pi_{\text{base}}$ ). For each  $(x, y, r_{\text{orig}}) \in \mathcal{D}_{\text{SFT}}$ , the Critic assesses  $r_{\text{orig}}$  using the criteria of *Cognitive Matching Degree*, where the Critic checks whether the complexity and difficulty of  $r_{\text{orig}}$  aligns with the cognitive capacity of  $\pi_{\text{base}}$ . Specifically, the Critic classifies the reasoning processes into three subsets: i)  $\mathcal{D}_{\text{easy}} : (x, y, r_{\text{easy}})$ , cases where the reasoning process is overly terse, making it difficult for  $\pi_{\text{base}}$  to follow; ii)  $\mathcal{D}_{\text{med}} : (x, y, r_{\text{med}})$ ,

<sup>7</sup>The decision boundary refers to the model’s ability to judge whether the produced CoT is correct and aligns with its own cognitive capabilities, enabling it to successfully solve problems following its CoT.

Level/Model Size	1.5B	7B	32B
Easy	195	80	19
Medium	296	389	354
Hard	9	31	127

Table 1: Complexity distributions of CoTs generated by different sizes of DeepSeek-R1-Distill-Qwen models.

cases with appropriate steps that enable successful problem solving; and iii)  $\mathcal{D}_{\text{hard}} : (x, y, r_{\text{hard}})$ , cases with overly redundant or excessively complex reasoning steps that exceed the comprehension of  $\pi_{\text{base}}$ , making it extremely prone to fail to guide  $\pi_{\text{base}}$  in solving  $x$ .

**Remarks.** An intuitive approach would be to use  $\pi_{\text{base}}$  itself as the Critic. However, due to its small parameter size (e.g., 7B), certain CoTs exceed  $\pi_{\text{base}}$ ’s comprehension, rendering it incapable of reliable complexity classification. Thus, we leverage the same LLM for the Rethinker (denoted as  $\pi_{\text{large}}$ ) to serve as the Critic, forcing it to “think” from the perspective of the small model  $\pi_{\text{base}}$ . A detailed analysis of the Critic choices is provided in the Experiments 4.3 and Appendix A.5.

**Hypothesis Verification.** To further verify that the complexity levels of CoTs are closely related to the cognitive capacities of reasoning models, we conduct an experiment in which we evaluate DeepSeek-R1-Distill-Qwen-1.5B/7B/32B on MATH500, collecting each model’s outputs. We employ the Critic to rate the level of model’s CoT outputs; each CoT is evaluated three times, and the final rating is determined by majority vote. For each model, we quantify the distribution of these CoTs across different complexity levels in Table 1. As shown, DeepSeek-R1-Distill-Qwen-1.5B yields the largest number of simple CoTs, while DeepSeek-R1-Distill-Qwen-32B generates the greatest number of difficult CoTs.

These findings demonstrate that the complexity of CoTs escalate as the model size increases, suggesting that larger models possess higher reasoning and cognitive capacities. Consequently, overly terse or complex CoTs may not be suitable for training models with lower cognitive abilities. It is therefore essential to use CoTs that align with the model’s cognitive trajectory to improve its reasoning capabilities, a strategy akin to “teaching according to the student’s ability.”

### 3.2.2 Rethinker

An LLM agent  $\pi_{\text{large}}$  is tasked with rewriting reasoning processes to achieve cognitive alignment.

For each  $(x, y, r_{\text{easy}}) \in \mathcal{D}_{\text{easy}}$ , the Rethinker expands  $r_{\text{easy}}$  by adding necessary steps for easier understanding, i.e.,  $r_{\text{easy}^*} = \pi_{\text{large}}(x, y, r_{\text{easy}})$ . Similarly, for each  $(x, y, r_{\text{hard}}) \in \mathcal{D}_{\text{hard}}$ , the Rethinker simplifies  $r_{\text{hard}}$  by removing redundancies or using simpler methods to solve the problem grounded in the correct answer:  $r_{\text{hard}^*} = \pi_{\text{large}}(x, y, r_{\text{hard}})$ . Cases of the rewriting process of the Rethinker are shown in Tables 11 and 12.

### 3.2.3 Verifier

Finally, we leverage the LLM agent  $\pi_{\text{base}}$  to validate the correctness of  $r_{\text{med}}$ ,  $r_{\text{easy}^*}$ , and  $r_{\text{hard}^*}$  in order to preserve the high quality of the dataset. It predicts whether  $\pi_{\text{base}}$  can derive the correct answer  $y$  from the rewritten thoughts  $r_{\text{easy}^*}$  or  $r_{\text{hard}^*}$ . Note that  $r_{\text{med}}$  has already been validated as correct in the original dataset, and we send  $r_{\text{med}}$  to the Verifier to further ensure data quality,

After verification, incorrect cases are sent back to the Rethinker to be continuously rewritten until they pass verification. In the implementation, cases that fail to pass verification after three iterations are discarded. For the cases that pass verification, we invoke the Critic to make the judgment again (please refer to Figure 2 for the algorithmic flow).

The final SFT dataset is composed of verified medium-level data:  $\mathcal{D}_{\text{SFT}^*} = \mathcal{D}_{\text{med}} \cup \mathcal{D}_{\text{easy}^*} \cup \mathcal{D}_{\text{hard}^*}$ , where  $\mathcal{D}_{\text{med}}$  denotes the verified medium-level data, and  $\mathcal{D}_{\text{easy}^*}$  and  $\mathcal{D}_{\text{hard}^*}$  represent the rewritten versions of  $\mathcal{D}_{\text{easy}}$  and  $\mathcal{D}_{\text{hard}}$  that have passed verification and have been re-rated as medium by the Critic, respectively.  $\mathcal{D}_{\text{SFT}^*}$  serves as the SFT training set in the CRV stage. Prompt templates used in CRV system are provided in Appendix C.

## 3.3 Cognitive Preference Optimization

The CogPO algorithm aligns CoT processes of smaller LLMs with their inherent cognitive capacities, following the SFT training using CRV system.

### 3.3.1 Preliminaries

Briefly speaking, the CogPO algorithm is extended from DPO (Rafailov et al., 2023) and its variants. Let  $y_w$  and  $y_l$  be the chosen and rejected responses for an instruction  $x$  (not restricted to reasoning problems addressed in this work), respectively. We further denote  $\pi_{\theta}$  as the model to be optimized after SFT and  $\pi_{\text{ref}}$  as the reference model. DPO seeks to maximize the following margin:  $M_{\beta}(x, y_w, y_l) = \beta \cdot \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$



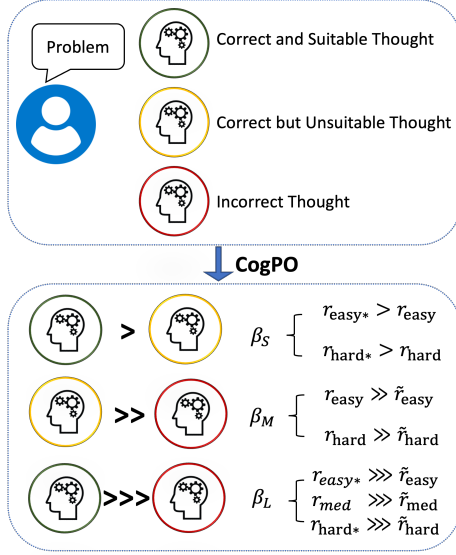


Figure 3: An illustration of CogPO, showing the different preference gaps between CoT pairs and the corresponding mini-tasks.

where  $\beta$  is a temperature hyperparameter. Based on  $M_\beta(x, y_w, y_l)$ , the DPO loss is defined as:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \log \sigma(M_\beta(x, y_w, y_l)). \quad (1)$$

The settings of  $\beta$  are critical to the performance of DPO.  $\beta$ -DPO (Wu et al., 2024a) further adjusts  $\beta$  according to  $M_\beta(x, y_w, y_l)$ , either at the instance level or batch level, allowing the model to adapt  $\beta$  based on the reward differential of the input data.

### 3.3.2 Algorithmic Description

As noted, DPO and  $\beta$ -DPO do not require any prior knowledge of how the model learns the user’s preferences. We suggest that this type of prior knowledge is critical for training better smaller reasoning models, as the cognitive trajectories of large and small models often differ (Li et al., 2022; Zhang et al., 2024; Hu et al., 2024), which may not be directly reflected in the reward differential. We propose CogPO to align reasoning preferences by encoding more prior knowledge and continuously training on a series of *mini-tasks*.

We leverage the Rethinker in CRV to also generate incorrect reasoning processes when asked to rewrite the original thought  $r_{\text{orig}}$  (prompt template is provided in Appendix C). The incorrect thoughts are denoted as  $\tilde{r}_{\text{med}}$ ,  $\tilde{r}_{\text{easy}}$ , and  $\tilde{r}_{\text{hard}}$ , based on their origins from  $\mathcal{D}_{\text{med}}$ ,  $\mathcal{D}_{\text{easy}}$ , and  $\mathcal{D}_{\text{hard}}$ . These thoughts contain factual errors or invalid reasoning steps, which can mislead  $\pi_{\text{base}}$ , rendering it impossible to solve  $x$ . Thus, we categorize the properties of all the thoughts we have collected

into the following three types: i)  $r_{\text{med}}$ ,  $r_{\text{easy}^*}$ , and  $r_{\text{hard}^*}$ : medium-level reasoning processes that are both correct and cognitively suitable for  $\pi_{\text{base}}$ ; ii)  $r_{\text{easy}}$  and  $r_{\text{hard}}$ : easy or hard thoughts that are correct but unsuitable for  $\pi_{\text{base}}$ ; iii)  $\tilde{r}_{\text{med}}$ ,  $\tilde{r}_{\text{easy}}$ , and  $\tilde{r}_{\text{hard}}$ : incorrect reasoning processes with logical flaws or invalid reasoning steps (regardless of the difficulty levels). To define the mini-tasks used for CogPO training, we consider the preference gaps in these three types of CoT pairs as follows:

**1. Small Gap Mini-task:** The pairs are  $(r_{\text{easy}^*}, r_{\text{easy}})$  and  $(r_{\text{hard}^*}, r_{\text{hard}})$ . Both are correct but differ in complexity (suitable vs. unsuitable for  $\pi_{\text{base}}$ ). We treat  $r_{\text{easy}^*}$  and  $r_{\text{hard}^*}$  as chosen reasoning processes ( $r_w$ ), and  $r_{\text{easy}}$  and  $r_{\text{hard}}$  as rejected ( $r_l$ ).

**2. Medium Gap Mini-task:** The pairs are  $(r_{\text{easy}}, \tilde{r}_{\text{easy}})$  and  $(r_{\text{hard}}, \tilde{r}_{\text{hard}})$ . The former are correct but unsuitable, while the latter are completely incorrect. As correctness is more important than suitability for our model, the preference gap of this mini-task should be higher than that in the previous case. For this mini-task,  $r_{\text{easy}}$  and  $r_{\text{hard}}$  are treated as  $r_w$ , while  $\tilde{r}_{\text{easy}}$  and  $\tilde{r}_{\text{hard}}$  are treated as  $r_l$ .

**3. Large Gap Mini-task:** The pairs are  $(r_{\text{med}}, \tilde{r}_{\text{med}})$ ,  $(r_{\text{easy}^*}, \tilde{r}_{\text{easy}})$ , and  $(r_{\text{hard}^*}, \tilde{r}_{\text{hard}})$ . Intuitively, the preference gaps should be the largest between suitable and correct thoughts and incorrect ones. Here,  $r_{\text{med}}$ ,  $r_{\text{easy}^*}$ , and  $r_{\text{hard}^*}$  are treated as  $r_w$ , while  $\tilde{r}_{\text{med}}$ ,  $\tilde{r}_{\text{easy}}$ , and  $\tilde{r}_{\text{hard}}$  are treated as  $r_l$ .

Following our modeling framework, each training instance  $(x, r_w, r_l)$  receives its specific  $\beta$  value, as illustrated in Figure 3. The CogPO objective function aggregates these preferences:

$$\mathcal{L}_{\text{CogPO}} = -\mathbb{E}_{(x, r_w, r_l) \sim \mathcal{D}} \log \sigma(M_{\beta_{\text{CogPO}}}(x, r_w, r_l)), \quad (2)$$

where  $\beta_{\text{CogPO}}$  is selected from  $\{\beta_S, \beta_M, \beta_L\}$ , depending on the specific types of mini-tasks (with  $\beta_S < \beta_M < \beta_L$ , corresponding to the three gaps). Overall, our CogPO algorithm enables granular preference learning: strong regularization ( $\beta_L$ ) for validity discrimination, moderate guidance ( $\beta_M$ ) for suitability alignment, and subtle refinement ( $\beta_S$ ) for reasoning style adaptation. This design provides more control over the alignment process, leading to further improvements on the basis of SFT (using CRV system).

**Remarks.** CogPO can be naturally combined with  $\beta$ -DPO (Wu et al., 2024a). We can redefine the  $\beta$  values  $\{\beta_S, \beta_M, \beta_L\}$  as follows:  $\beta_i^* = \beta_i + \alpha \cdot (M_i - M_0) \cdot \beta_i$  where  $\beta_i$  is chosen from  $\{\beta_S, \beta_M, \beta_L\}$  based on the corresponding gap type,  $M_i$  is the

Dataset/Model	Zero-shot	SFT	CRV+SFT	DPO	$\beta$ -DPO	SimPO	CogPO
AIME2024	10.0	20.0	<b>26.7</b>	23.3	23.3	<u>26.7</u>	<b>30.0</b>
MATH-500	73.6	80.0	<b>84.0</b>	83.4	83.8	<u>84.2</u>	<b>84.4</b>
GSM8K	89.5	92.3	<b>92.7</b>	92.6	<u>93.0</u>	92.6	<b>93.3</b>
GPQA Diamond	33.3	37.4	<b>40.9</b>	40.0	37.4	<b>40.9</b>	<b>40.9</b>
LiveCodeBench V2	30.7	31.3	<b>34.4</b>	34.4	35.8	<u>36.2</u>	<b>36.6</b>
MMLU	71.9	76.1	<b>76.5</b>	76.1	76.0	<b>76.5</b>	<b>76.5</b>
OlympiadBench (math-en)	40.1	43.6	<b>45.8</b>	45.7	<u>46.5</u>	46.0	<b>46.6</b>

Table 2: Performance comparison of various training methods. The LLM backbone is Qwen2.5-7B-Instruct, and the training set is Bespoke-Stratos-17k. Results are shown for zero-shot (without further training), SFT, CRV+SFT, DPO,  $\beta$ -DPO, SimPO, and CogPO. DPO,  $\beta$ -DPO, SimPO, and CogPO are conducted on the same model checkpoints of CRV+SFT, using the same preference pair dataset. The metrics represent scores for these tasks, with the best results for each dataset in each group marked in bold and the second-best underlined.

Dataset/Model	LLaMA-O1	Macro-o1	Bespoke-Stratos-7B	Ours	OpenThinker-7B	Ours
Training Set Size	332K	60K	17K	17K	114K	114K
AIME2024	3.3	6.7	20.0	<b>30.0</b>	31.3	<b>43.3</b>
MATH500	28.6	38.4	82.0	<b>84.4</b>	83.0	<b>88.4</b>
GPQA Diamond	26.3	31.8	37.8	<b>40.9</b>	42.4	<b>42.9</b>
LiveCodeBench V2	1.6	24.9	36.1	<b>36.6</b>	39.9	<b>46.4</b>

Table 3: Comparison between our model and other small reasoning models in the open-source community. Specifically, we train two versions using our approach on Bespoke-Stratos-17k and OpenThoughts-114k, respectively, where the two training sets are the same with Bespoke-Stratos-7B and OpenThinker-7B, respectively.

instance-level reward differential, and  $M_0$  is a pre-defined threshold as in (Wu et al., 2024a).<sup>8</sup>

## 4 Experiments

To evaluate the effectiveness of the CRV framework and the CogPO algorithm, we conduct a series of experiments on challenging reasoning benchmarks. Due to space limitation, datasets and experimental settings are shown in the Appendix A.1 and A.2.

### 4.1 Main Experimental Results and Ablations

We choose Bespoke-Stratos-17k as the training set. Table 2 presents the results of our CRV framework and the CogPO algorithm on various reasoning benchmarks. CRV+SFT surpasses direct SFT on all benchmarks. Building on CRV+SFT, CogPO further enhances the model’s reasoning capability, surpasses other preference-optimization algorithms, and ultimately achieving the most outstanding performance, demonstrating its ability to align the model’s reasoning processes with its cognitive capacities. These results reveal that our CRV+CogPO framework effectively enhances the reasoning capabilities of smaller models, outperforming other traditional methods by a large margin.

<sup>8</sup>In our experiment, the combination does not yield substantial improvements, as prior knowledge is more important for our task. Hence, we stick to the usage of  $\mathcal{L}_{\text{CogPO}}$ .

### 4.2 Comparison Against Other Models

We compare our trained 7B model with other models released in the open-source community. We consider two reasoning LLMs available before the launch of DeepSeek-R1, namely Macro-o1 (Zhao et al., 2024) and LLaMA-O1<sup>9</sup>. We also compare other models trained on datasets distilled from DeepSeek-R1, including Bespoke-Stratos-7B<sup>10</sup> and OpenThinker-7B<sup>11</sup>. Using our CRV+CogPO framework, we also train two models on the Bespoke-Stratos-17k and OpenThoughts-114k training sets, respectively. Thus, it is fair to compare our method against those of Bespoke-Stratos-7B and OpenThinker-7B. The results, along with the sizes of the training sets, are shown in Table 3. It can be observed that employing DeepSeek-R1-generated CoT data yields superior results. At the algorithmic level, both Bespoke-Stratos-7B and our model are trained on the 17K CoTs from DeepSeek-R1. Under identical data conditions, our model significantly outperforms Bespoke-Stratos-7B across all benchmarks and achieves performance comparable to OpenThinker-7B, which is trained on 114K CoTs from DeepSeek-R1. Moreover, when trained on the same dataset

<sup>9</sup><https://huggingface.co/SimpleBerry/LLaMA-O1-Supervised-1129>

<sup>10</sup><https://huggingface.co/bespokelabs/Bespoke-Stratos-7B>

<sup>11</sup><https://huggingface.co/open-thoughts/OpenThinker-7B>

Model Backbone (The Critic)	AIME2024	MATH-500	GPQA-D	GSM8K	LCB V2	OlympiadBench
Qwen2.5-7B-Instruct	13.3	80.2	<b>40.9</b>	92.3	30.5	43.9
Qwen2.5-32B-Instruct	23.3	82.2	39.9	92.6	33.3	45.1
Qwen2.5-72B-Instruct	20.0	81.8	36.4	<b>92.7</b>	30.5	42.0
DeepSeek-R1-Distill-Qwen-32B	<b>26.7</b>	<b>84.0</b>	<b>40.9</b>	<b>92.7</b>	<b>34.4</b>	<b>45.8</b>

Table 4: Comparison using different backbones as the Critic. All the results are produced using CRV+SFT without CogPO on Bespoke-Stratos-17k.

Dataset/Model	Easy	Medium	Hard
AIME2024	13.3	23.3	16.7
MATH500	75.4	82.8	78.2
GPQA-D	34.3	37.4	33.3
LCB V2	31.9	36.2	32.5

Table 5: Experimental results on training data of different complexity levels.

as OpenThinker-7B, our model substantially surpasses OpenThinker-7B on all benchmarks. These findings demonstrate that, given the same data, our CRV + CogPO training framework exhibits superior performance, confirming its effectiveness.

### 4.3 Study on Choices of the Critic

In the previous section, we claimed that using the small target LLM  $\pi_{\text{base}}$  as the Critic does not necessarily produce satisfactory results due to its limited parameter size. In contrast, larger LLMs  $\pi_{\text{large}}$  can “think like small models” better. The results of using different backbones as the Critic are shown in Table 4, with the backbones for the Rethinker and the Verifier unchanged. From the results, we can see that they confirm our findings, as larger models consistently perform better than the 7B model in almost all tasks. Among the three large agents, DeepSeek-R1-Distill-Qwen-32B exhibits the best performance based on majority voting across all testing sets. A detailed and in-depth analysis of the selection of the Critic is provided in Appendix A.5.

### 4.4 Training with CoT Datasets of Different Complexity Levels

To further investigate whether medium-level data are indeed the most suitable for base model, we conduct experiments on the OpenThoughts-114K dataset. We used the Critic to rate all CoTs in the dataset, then randomly sampled 10K CoTs from each of the derived easy, medium, and hard subsets to construct three training sets. We then perform SFT with Qwen2.5-7B-Instruct on these three training sets under identical configurations. The results are shown in Table 5, indicating that when the num-

Dataset/Model	SFT	w. C	w. CR	w. CRV
AIME2024	20.0	23.3	26.7	26.7
MATH500	80.0	83.4	83.2	84.0
GPQA-D	37.4	38.4	39.9	40.9
LCB V2	31.3	34.3	34.1	34.4

Table 6: Ablation results on the CRV system.

ber of training data is the same, the model trained on the medium subset achieves the highest scores, fully supporting our hypothesis. The CoTs in the easy and hard sets are either too terse or overly complex, preventing the base model from effectively comprehending all CoTs in those sets. In contrast, the medium subset data align with the model’s cognitive capabilities and thus yield the best results.

### 4.5 Study on Effectiveness of Critic, Rethinker and Verifier

To further explore the collaborative mechanism within the CRV system and the individual roles and contributions of each module, we conduct extensive ablation experiments on the Bespoke-Stratos-17k dataset. Table 6 presents our ablation results. The “SFT” row reports results from directly performing SFT on the original dataset without any CRV intervention; the “w. C” row shows performance when only the Critic is applied before SFT, using only the traces rated as medium by the Critic for SFT; the “w. CR” row indicates results when both the Critic and the Rethinker participate prior to SFT, utilizing the medium-rated traces and the refined easy/hard traces that have not yet been verified; the “w. CRV” row reflects outcomes when the Critic, Rethinker, and Verifier are all applied.

As the Critic, Rethinker, and Verifier participate sequentially, the model’s reasoning ability exhibits a progressively improving trend, which clearly illustrates the role of each component. Notably, “w. CR” experiences a performance drop on MATH500 and LCB V2, indicating that omitting the Verifier after the Rethinker’s refinement could impair the model’s reasoning ability. Therefore, each component of the CRV system plays an indispensable

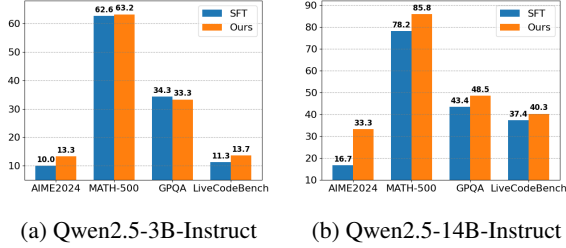


Figure 4: Experimental results of different sizes of Qwen2.5 models on AIME2024, MATH500, GPQA Diamond, and LiveCodeBench V2.

role. To achieve optimal performance, we recommend processing the data using the complete CRV system.

#### 4.6 Study on Model Scales

To study the effectiveness of different parameter sizes on the student models, we further report the performance of Qwen2.5-3B-Instruct and Qwen2.5-14B-Instruct. The experimental settings are identical to those of Qwen2.5-7B-Instruct. The results are presented in Figure 4. We observe that our method is also effective across different model scales. An interesting observation is that the improvement is more significant in Qwen2.5-14B-Instruct compared to Qwen2.5-3B-Instruct. This is because, even when we leverage the CRV system to rewrite the CoTs, the large capacity gap between the teacher and student models makes it more challenging for Qwen2.5-3B-Instruct to capture the CoTs through SFT. This finding is also consistent with the recently discovered “distillation scaling law” (Busbridge et al., 2025).

#### 4.7 Study on Other Model Backbones

To evaluate the universality of the proposed approach, we perform additional experiments on multiple backbones beyond the Qwen2.5 series on Bespoke-Stratos-17k dataset. Figure 5 demonstrates that, for both LLaMA and Mistral series, our approach achieves notable performance gains over the direct SFT baseline across diverse mathematical and coding tasks. These results indicate that the CRV+CogPO framework enables seamless adaptation to other backbones, demonstrating the universality of our approach on various LLM backbones, which also shows the potential of our work to produce stronger models based on other LLMs.

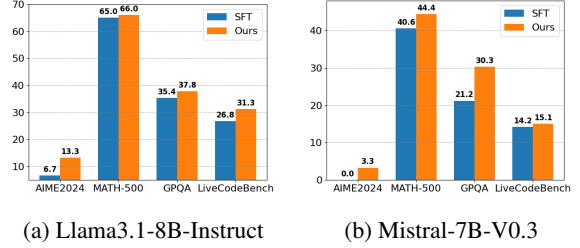


Figure 5: Experimental results of different model series (Llama3.1-8B-Instruct, Mistral-7B-V0.3) other than Qwen2.5 on AIME2024, MATH500, GPQA Diamond, and LiveCodeBench V2.

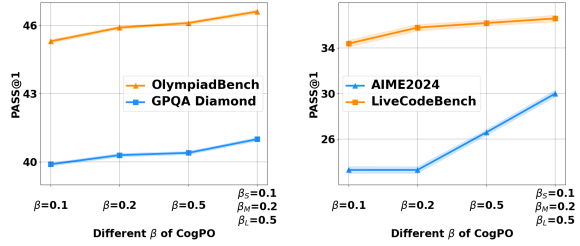


Figure 6: Impact of different  $\beta$  on AIME2024, GPQA Diamond, LiveCodeBench V2 and OlympiadBench.

#### 4.8 Hyper-parameter Analysis

To evaluate the impact of  $\beta$  values in CogPO, we perform a series of experiments with varying  $\beta$  values to assess the algorithm’s effectiveness. As shown in Figure 6, the highest performance is attained when assigning tailored  $\beta$  values to samples based on their specific gaps, which is a core principle of the CogPO algorithm.

#### 4.9 Case Studies

Due to space limitations, case studies are shown in the appendix. They clearly show how our approach can effectively expand or simplify the reasoning processes based on the Critic.

### 5 Conclusion and Future Work

In this paper, we present the CRV framework where we leverage the strengths of LLM agents to critique, refine, and verify CoT outputs for optimizing CoT training sets. The CogPO algorithm further aligns model outputs with their inherent cognitive capacities, improving the performance on several challenging reasoning tasks. In the future, we will i) train and release stronger small models using larger CoT datasets; ii) improve the effectiveness of the CRV framework, especially for much smaller models; and iii) investigate our approach for other domain-specific applications, such as medical diagnosis and legal reasoning.



## Limitations

While our proposed framework shows promising results in enhancing the reasoning capabilities of smaller LLMs, several limitations still remain. The CRV framework relies heavily on the contributions of larger models refining the CoT output. This dependency may create challenges in situations where access to larger models is restricted, or these larger models generate incorrect results. In addition, although our framework is designed for smaller LLMs, there remains a ceiling on their performance. By nature, smaller models inherently have reduced capacity to encode complex information and handle nuanced reasoning tasks, which may limit their effectiveness in certain scenarios.

## Ethical Considerations

Our work is fully methodological; hence, there are no direct ethical issues. However, smaller models trained on data distilled from larger ones might inherit or exacerbate biased outputs, which can still influence outcomes. We suggest that continuous evaluation of trained LLMs based on ethical guidelines is indispensable.

## References

- Hritik Bansal, Arian Hosseini, Rishabh Agarwal, Vinh Q. Tran, and Mehran Kazemi. 2024. [Smaller, weaker, yet better: Training LLM reasoners via compute-optimal sampling](#). *CoRR*, abs/2408.16737.
- Dan Busbridge, Amitis Shidani, Floris Weers, Jason Ramapuram, Etai Littwin, and Russ Webb. 2025. [Distillation scaling laws](#). *CoRR*, abs/2502.08606.
- Huayu Chen, Guande He, Hang Su, and Jun Zhu. 2024. [Noise contrastive alignment of language models with explicit rewards](#). *CoRR*, abs/2402.05369.
- Kaiyuan Chen, Jin Wang, and Xuejie Zhang. 2025. [Learning to reason via self-iterative process feedback for small language models](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3027–3042. Association for Computational Linguistics.
- Lihu Chen and Gaël Varoquaux. 2024. [What is the role of small models in the LLM era: A survey](#). *CoRR*, abs/2409.06857.
- Zeming Chen, Qiyue Gao, and Lawrence S. Moss. 2021. [NeuralLog: Natural language inference with joint neural and logical reasoning](#). In *Proceedings of \*SEM 2021: The Tenth Joint Conference on Lexical and Computational Semantics*, pages 78–88, Online. Association for Computational Linguistics.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168.
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *CoRR*, abs/2501.12948.
- 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, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Alonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. [The llama 3 herd of models](#). *CoRR*, abs/2407.21783.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. [KTO: model alignment as prospect theoretic optimization](#). *CoRR*, abs/2402.01306.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [Olympiadbench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pages 3828–3850. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language](#)

672	<a href="#">understanding</a> . In <i>9th International Conference on Learning Representations</i> . OpenReview.net.	
673		
674	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh,	
675	Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay	
676	Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. <a href="#">Dis-</a>	
677	<a href="#">tilling step-by-step! outperforming larger language</a>	
678	<a href="#">models with less training data and smaller model</a>	
679	<a href="#">sizes</a> . In <i>Findings of the Association for Computa-</i>	
680	<i>tational Linguistics: ACL 2023</i> , pages 8003–8017.	
681	Association for Computational Linguistics.	
682	Linmei Hu, Hongyu He, Duokang Wang, Ziwang Zhao,	
683	Yingxia Shao, and Liqiang Nie. 2024. <a href="#">LLM vs small</a>	
684	<a href="#">model? large language model based text augmenta-</a>	
685	<a href="#">tion enhanced personality detection model</a> . In <i>Thirty-</i>	
686	<i>Eighth AAAI Conference on Artificial Intelligence</i> ,	
687	pages 18234–18242. AAAI Press.	
688	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia	
689	Yan, Tianjun Zhang, Sida Wang, Armando Solar-	
690	Lezama, Koushik Sen, and Ion Stoica. 2024. <a href="#">Live-</a>	
691	<a href="#">codebench: Holistic and contamination free eval-</a>	
692	<a href="#">uation of large language models for code</a> . <i>CoRR</i> ,	
693	abs/2403.07974.	
694	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	
695	sch, Chris Bamford, Devendra Singh Chaplot, Diego	
696	de Las Casas, Florian Bressand, Gianna Lengyel,	
697	Guillaume Lample, Lucile Saulnier, L��lio Ren-	
698	nard Lavaud, Marie-Anne Lachaux, Pierre Stock,	
699	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-	
700	th��e Lacroix, and William El Sayed. 2023. <a href="#">Mistral</a>	
701	<a href="#">7b</a> . <i>CoRR</i> , abs/2310.06825.	
702	Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu,	
703	Yu Zhang, Zhenguo Li, and James T. Kwok. 2024.	
704	<a href="#">Forward-backward reasoning in large language mod-</a>	
705	<a href="#">els for mathematical verification</a> . In <i>Findings of</i>	
706	<i>the Association for Computational Linguistics, ACL</i>	
707	<i>2024</i> , pages 6647–6661. Association for Computa-	
708	tional Linguistics.	
709	Veronica Latcinnik and Jonathan Berant. 2020. <a href="#">Explain-</a>	
710	<a href="#">ing question answering models through text genera-</a>	
711	<a href="#">tion</a> . <i>arXiv preprint arXiv:2004.05569</i> .	
712	Chenglin Li, Qianglong Chen, Caiyu Wang, and	
713	Yin Zhang. 2023. <a href="#">Mixed distillation helps</a>	
714	<a href="#">smaller language model better reasoning</a> . <i>CoRR</i> ,	
715	abs/2312.10730.	
716	Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen,	
717	Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian,	
718	Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng	
719	Yan. 2022. <a href="#">Explanations from large language models</a>	
720	<a href="#">make small reasoners better</a> . <i>CoRR</i> , abs/2210.06726.	
721	Zhiming Li, Yushi Cao, Xiufeng Xu, Junzhe Jiang,	
722	Xu Liu, Yon Shin Teo, Shang-Wei Lin, and Yang	
723	Liu. 2024. <a href="#">LLMs for relational reasoning: How far</a>	
724	<a href="#">are we?</a> In <i>LLM4CODE@ICSE</i> , pages 119–126.	
	Zhenwen Liang, Ye Liu, Tong Niu, Xiangliang Zhang,	725
	Yingbo Zhou, and Semih Yavuz. 2024. <a href="#">Improv-</a>	726
	<a href="#">ing LLM reasoning through scaling inference com-</a>	727
	<a href="#">putation with collaborative verification</a> . <i>CoRR</i> ,	728
	abs/2410.05318.	729
	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri	730
	Edwards, Bowen Baker, Teddy Lee, Jan Leike, John	731
	Schulman, Ilya Sutskever, and Karl Cobbe. 2023.	732
	<a href="#">Let’s verify step by step</a> . <i>CoRR</i> , abs/2305.20050.	733
	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	734
	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	735
	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	736
	Sean Welleck, Bodhisattwa Prasad Majumder,	737
	Shashank Gupta, Amir Yazdanbakhsh, and Peter	738
	Clark. 2023. <a href="#">Self-refine: Iterative refinement with</a>	739
	<a href="#">self-feedback</a> . <i>CoRR</i> , abs/2303.17651.	740
	Yu Meng, Mengzhou Xia, and Danqi Chen. 2024.	741
	<a href="#">Simplo: Simple preference optimization with a</a>	742
	<a href="#">reference-free reward</a> . <i>CoRR</i> , abs/2405.14734.	743
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	744
	Carroll L. Wainwright, Pamela Mishkin, Chong	745
	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray,	746
	John Schulman, Jacob Hilton, Fraser Kelton, Luke	747
	Miller, Maddie Simens, Amanda Askell, Peter Welin-	748
	der, Paul F. Christiano, Jan Leike, and Ryan Lowe.	749
	2022. <a href="#">Training language models to follow instruc-</a>	750
	<a href="#">tions with human feedback</a> . In <i>Advances in Neural</i>	751
	<i>Information Processing Systems 35: Annual Con-</i>	752
	<i>ference on Neural Information Processing Systems</i>	753
	<i>2022</i> .	754
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	755
	pher D. Manning, Stefano Ermon, and Chelsea Finn.	756
	2023. <a href="#">Direct preference optimization: Your language</a>	757
	<a href="#">model is secretly a reward model</a> . In <i>Advances in</i>	758
	<i>Neural Information Processing Systems 36: Annual</i>	759
	<i>Conference on Neural Information Processing Sys-</i>	760
	<i>tems 2023</i> .	761
	David Rein, Betty Li Hou, Asa Cooper Stickland,	762
	Jackson Petty, Richard Yuanzhe Pang, Julien Di-	763
	rani, Julian Michael, and Samuel R. Bowman. 2023.	764
	<a href="#">GPQA: A graduate-level google-proof q&amp;a bench-</a>	765
	<a href="#">mark</a> . <i>CoRR</i> , abs/2311.12022.	766
	Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023.	767
	<a href="#">Reflexion: an autonomous agent with dynamic mem-</a>	768
	<a href="#">ory and self-reflection</a> . <i>CoRR</i> , abs/2303.11366.	769
	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya	770
	Sachan. 2022. <a href="#">Distilling multi-step reasoning ca-</a>	771
	<a href="#">pabilities of large language models into smaller</a>	772
	<a href="#">models via semantic decompositions</a> . <i>CoRR</i> ,	773
	abs/2212.00193.	774
	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya	775
	Sachan. 2023. <a href="#">Distilling reasoning capabilities into</a>	776
	<a href="#">smaller language models</a> . In <i>Findings of the Associa-</i>	777
	<i>tion for Computational Linguistics: ACL 2023</i> , pages	778
	7059–7073. Association for Computational Linguis-	779
	tics.	780

781	Vered Shwartz, Peter West, Ronan Le Bras, Chandra	2023b. <a href="#">Tree of thoughts: Deliberate problem solving</a>	838
782	Bhagavatula, and Yejin Choi. 2020. <a href="#">Unsupervised</a>	<a href="#">with large language models</a> . In <i>Advances in Neural</i>	839
783	<a href="#">commonsense question answering with self-talk</a> . In	<i>Information Processing Systems 36: Annual Confer-</i>	840
784	<i>Proceedings of the 2020 Conference on Empirical</i>	<i>ence on Neural Information Processing Systems 2023,</i>	841
785	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<i>NeurIPS 2023</i> .	842
786	pages 4615–4629, Online. Association for Computa-		
787	tional Linguistics.		
788	WeiQi Wang, Tianqing Fang, Chunyang Li, Haochen	Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding,	843
789	Shi, Wenxuan Ding, Baixuan Xu, Zhaowei Wang,	Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen,	844
790	Jiaxin Bai, Xin Liu, Cheng Jiayang, Chunkit Chan,	Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen	845
791	and Yangqiu Song. 2024. <a href="#">CANDLE: iterative con-</a>	Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun.	846
792	<a href="#">ceptualization and instantiation distillation from large</a>	2024. <a href="#">Advancing LLM reasoning generalists with</a>	847
793	<a href="#">language models for commonsense reasoning</a> . In	<a href="#">preference trees</a> . <i>CoRR</i> , abs/2404.02078.	848
794	<i>Proceedings of the 62nd Annual Meeting of the Asso-</i>		
795	<i>ciation for Computational Linguistics</i> , pages 2351–	Yuanhao Yue, Chengyu Wang, Jun Huang, and Peng	849
796	2374. Association for Computational Linguistics.	Wang. 2024. <a href="#">Distilling instruction-following abilities</a>	850
797	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<a href="#">of large language models with task-aware curriculum</a>	851
798	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	<a href="#">planning</a> . In <i>Findings of the Association for Com-</i>	852
799	and Denny Zhou. 2022a. <a href="#">Chain-of-thought prompt-</a>	<i>putational Linguistics: EMNLP 2024</i> , pages 6030–	853
800	<a href="#">ing elicits reasoning in large language models</a> . In	6054. Association for Computational Linguistics.	854
801	<i>Advances in Neural Information Processing Systems</i>		
802	<i>35: Annual Conference on Neural Information Pro-</i>	Biao Zhang, Zhongtao Liu, Colin Cherry, and Orhan	855
803	<i>cessing Systems 2022</i> .	Firat. 2024. <a href="#">When scaling meets LLM finetuning:</a>	856
804	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<a href="#">The effect of data, model and finetuning method</a> . In	857
805	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	<i>The Twelfth International Conference on Learning</i>	858
806	and Denny Zhou. 2022b. <a href="#">Chain-of-thought prompt-</a>	<i>Representations</i> . OpenReview.net.	859
807	<a href="#">ing elicits reasoning in large language models</a> . In		
808	<i>NeurIPS</i> .	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	860
809	Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu,	Xiaolei Wang, Yupeng Hou, Yingqian Min, Be-	861
810	Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan	ichen Zhang, Junjie Zhang, Zican Dong, Yifan Du,	862
811	He. 2024a. <a href="#"><math>\beta</math>-dpo: Direct preference optimization</a>	Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao	863
812	<a href="#">with dynamic <math>\beta</math></a> . <i>CoRR</i> , abs/2407.08639.	Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang	864
813	Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yim-	Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen.	865
814	ing Yang, and Quanquan Gu. 2024b. <a href="#">Self-play pref-</a>	2023. <a href="#">A survey of large language models</a> . <i>CoRR</i> ,	866
815	<a href="#">erence optimization for language model alignment</a> .	abs/2303.18223.	867
816	<i>CoRR</i> , abs/2405.00675.		
817	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan,	Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi	868
818	Lingfeng Shen, Benjamin Van Durme, Kenton Mur-	Shi, Chenyang Lyu, Longyue Wang, Weihua Luo,	869
819	ray, and Young Jin Kim. 2024. <a href="#">Contrastive prefer-</a>	and Kaifu Zhang. 2024. <a href="#">Marco-o1: Towards open</a>	870
820	<a href="#">ence optimization: Pushing the boundaries of LLM</a>	<a href="#">reasoning models for open-ended solutions</a> . <i>CoRR</i> ,	871
821	<a href="#">performance in machine translation</a> . In <i>Forty-first In-</i>	abs/2411.14405.	872
822	<i>ternational Conference on Machine Learning</i> . Open-		
823	Review.net.		
824	Junbing Yan, Chengyu Wang, Taolin Zhang, Xiaofeng		
825	He, Jun Huang, and Wei Zhang. 2023. <a href="#">From complex</a>		
826	<a href="#">to simple: Unraveling the cognitive tree for reasoning</a>		
827	<a href="#">with small language models</a> . In <i>Findings of the Asso-</i>		
828	<i>ciation for Computational Linguistics: EMNLP 2023</i> ,		
829	pages 12413–12425. Association for Computational		
830	Linguistics.		
831	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
832	Thomas L. Griffiths, Yuan Cao, and Karthik		
833	Narasimhan. 2023a. <a href="#">Tree of thoughts: Deliberate</a>		
834	<a href="#">problem solving with large language models</a> . <i>CoRR</i> ,		
835	abs/2305.10601.		
836	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
837	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.		



Dataset	Size
AIME2024	30
MATH-500	500
GSM8K	1319
GPQA Diamond	198
LiveCodeBench V2	511
MMLU	14042
OlympiadBench (math-en)	674

Table 7: Testing set statistics.

## A Supplementary Experiments

### A.1 Datasets

In our experiments, we evaluate our work on several benchmarks, including AIME2024<sup>12</sup>, MATH-500 (Lightman et al., 2023), GSM8K (Cobbe et al., 2021), GPQA Diamond (Rein et al., 2023), LiveCodeBench V2 (Jain et al., 2024), MMLU (Hendrycks et al., 2021), and OlympiadBench (math-en) (He et al., 2024). The sizes of our testing sets are summarized in Table 7.

For our training set  $\mathcal{D}_{\text{SFT}^*}$ , we leverage Bespoke-Stratos-17k<sup>13</sup>, which contains 17K tuples of questions, reasoning processes, and answers directly distilled from DeepSeek-R1 (DeepSeek-AI, 2025). We also utilize two released CoT datasets to conduct supplementary experiments. The first one is Sky-T1-data-17k<sup>14</sup>, which is distilled from QwQ-32B-Preview, whose reasoning abilities are reported to be weaker than those of DeepSeek-R1. The second one is OpenThoughts-114k<sup>15</sup>, which is distilled from DeepSeek-R1 and verified using a data curation recipe. We have chosen not to use some previously released CoT datasets (e.g., OpenLongCoT-SFT<sup>16</sup>) due to their significantly weaker reasoning abilities, while some benchmarks (e.g., AIME2024, OlympiadBench) are extremely challenging.

### A.2 Experimental Details

In our work, we utilize Qwen2.5-7B-Instruct as the default model backbone and extend our evaluation to Llama3.1-8B-Instruct (Dubey et al., 2024) and

<sup>12</sup><https://maa.org/math-competitions/american-invitational-mathematics-examination-aime>

<sup>13</sup><https://huggingface.co/datasets/bespokelabs/Bespoke-Stratos-17k>

<sup>14</sup><https://github.com/NovaSky-AI/SkyThought>

<sup>15</sup><https://huggingface.co/datasets/open-thoughts/OpenThoughts-114k>

<sup>16</sup><https://huggingface.co/datasets/SimpleBerry/OpenLongCoT-SFT>

Hyperparameter	Value
<i>CRV Stage</i>	
Batch size	96
Learning rate	1e-5
Learning epoch	3.0
<i>CogPO Stage</i>	
Batch size	96
Learning rate	5e-7
Learning epoch	1.0
<i>SFT (Baseline)</i>	
Batch size	96
Learning rate	1e-5
Learning epoch	3.0
<i>DPO (Baseline)</i>	
Batch size	96
Learning rate	5e-7
Learning epoch	1.0
$\beta$	0.1
<i>SimPO (Baseline)</i>	
Batch size	96
Learning rate	5e-7
Learning epoch	1.0
$\beta$	2.0
$\gamma$	0.3

Table 8: Training hyperparameters.

Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), along with other sizes of Qwen2.5 models, to validate the generalizability of our algorithm across diverse model architectures and sizes. We first establish a baseline by assessing the model’s zero-shot capabilities. Subsequent experiments leverage this result to quantify the performance improvements attributable to CRV and CogPO. During the CRV phase, the same generation hyperparameters are applied to the Critic, Rethinker, and Verifier for inference: temperature  $T = 0.7$ , top\_p = 0.9, and top\_k = 50. The default backbone is DeepSeek-R1-Distill-Qwen-32B, while we test other backbone choices in the experiments. For CogPO training, the default  $\beta$  settings are:  $\beta_S = 0.1$ ,  $\beta_M = 0.2$ , and  $\beta_L = 0.5$ . Training details for model training and baselines are shown in Table 8.

On the Bespoke-Stratos-17k dataset, for the 3B model, we use a single node with 8 A800 GPUs (80GB), with an approximate training time of 4 hours. For the 7B model, we use a single node with 8 A800 GPUs (80GB), with a training time of about 5 hours. For the 14B model, we use 4 nodes, each with 8 A800 GPUs (80G), resulting in a training time of approximately 14 hours.



Dataset/Model	Zero-shot	SFT	Ours
AIME2024	10.0	16.7	<b>20.0</b>
MATH-500	73.6	73.2	<b>77.0</b>
GPQA Diamond	33.3	28.8	<b>36.9</b>
LiveCodeBench V2	30.7	20.9	<b>33.3</b>

Table 9: Performance comparison using Sky-T1-data-17k as the training set.

Dataset/Model	Zero-shot	SFT	Ours
AIME2024	10.0	31.3	<b>43.3</b>
MATH-500	73.6	83.0	<b>88.4</b>
GPQA Diamond	33.3	42.4	<b>42.9</b>
LiveCodeBench V2	30.7	39.9	<b>46.4</b>

Table 10: Performance comparison using OpenThoughts-114k as the training set.

### A.3 Results on Weaker CoT Dataset

To demonstrate that our approach is truly superior to vanilla SFT over CoT datasets, we conduct an experiment on the Sky-T1 dataset, which is relatively weaker than Bespoke-Stratos-17k due to the choice of the teacher model (i.e., QwQ-32B-Preview) and the data curation pipeline. The results are presented in Table 9. As shown, in some cases, the SFT baseline cannot even beat the zero-shot performance. This negative finding is also consistent with their own blog regarding the model size and data quality<sup>17</sup>. Nonetheless, by comparing our method with the SFT baseline, we can observe clear improvement, which demonstrates that our approach has the efficacy to enhance the reasoning abilities of small models in various scenarios.

### A.4 Results on Larger CoT Dataset

We further evaluate the performance of our method using OpenThoughts-114k as the training set, which is much larger than other training sets. This dataset is distilled from DeepSeek-R1 and goes through several quality verification steps. The results are presented in Table 10. It can be seen that our method ultimately exhibits exceptionally strong reasoning performance, significantly surpassing SFT on all benchmarks. This underscores the scalability and generalizability of our CRV+CogPO framework to larger datasets.

### A.5 Design Choice of the Critic

An initial, straightforward approach is to employ  $\pi_{\text{base}}$  as the Critic. However, owing to the small

model’s limited reasoning capability, it consistently faces difficulties in distinguishing the difficulty levels of CoTs effectively within our datasets. Note that for CoTs rated as “easy” or “hard”, the CoT is either overly concise (omitting necessary steps) or excessively complex, rendering it unintelligible to the small model and preventing it from following the chain to arrive at the correct answer. Under these circumstances, it is clearly unreasonable to require the small model to classify the CoT difficulty that it cannot comprehend effectively.

Another intuitive CoT evaluation approach is to input the problem and its corresponding CoT into the small model and then verify whether the model can arrive at the correct answer. However, applying this method directly would only partition CoT processes into “correct” or “incorrect” categories. For incorrect CoTs, this binary classification fails to distinguish the root cause of errors (i.e., whether the CoT is overly simplified or overly complex), which is critical for determining appropriate refinement strategies (e.g., expansion for overly simplified processes vs. simplification for overly complex ones).

Consequently, we utilize the larger and stronger LLM used in both the Rethinker and the Verifier (referred to as  $\pi_{\text{large}}$ ) to act as the Critic. This involves guiding the large model to simulate the cognitive approach of the smaller model,  $\pi_{\text{base}}$ . The prompt template of the Critic is shown in Table 13. This setting is akin to educational practices, where professors, instead of students, customarily curate academic content across a spectrum of difficulty levels due to their broader knowledge base. As shown in Table 4, the experiments clearly demonstrate the superior evaluative proficiency of the large model, confirming its advantage in categorizing CoT complexity from the perspective of the smaller model efficiently.

## B Case Studies

Case studies are presented in Table 11 and Table 12.

## C Prompt Templates

Prompt templates of the Critic, the Rethinker and the Verifier in our CRV system are shown in Table 13.

<sup>17</sup><https://novasky-ai.github.io/posts/sky-t1/>

<b>Problem</b>	Find the inverse of matrix $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$
<b>Answer</b>	$A^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$
<b>Original reasoning process</b> (correct but unsuitable)	Calculate determinant $\det(A) = 3$ , thus $A^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$
<b>Extended reasoning process</b> (correct and suitable)	Compute determinant: $2 \times 2 - 1 \times 1 = 3$ Construct adjugate: $\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$ Normalize: $A^{-1} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$
<b>Incorrect reasoning process</b>	Swap diagonal elements: $\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$

Table 11: Case study of how the reasoning process is extended.

<b>Problem</b>	Find the area of a triangle with vertices at $(0, 0)$ , $(3, 0)$ , and $(0, 4)$
<b>Answer</b>	6
<b>Original reasoning process</b> (correct but unsuitable)	Vector Representation: $\vec{AB} = \begin{pmatrix} 3 \\ 0 \end{pmatrix}, \vec{AC} = \begin{pmatrix} 0 \\ 4 \end{pmatrix}$ Determinant Method: $\text{Area} = \frac{1}{2} \left\  \begin{vmatrix} 3 & 0 \\ 0 & 4 \end{vmatrix} \right\  = \frac{1}{2}(12) = 6$
<b>Simplified reasoning process</b> (correct and suitable)	Recognize right-angled triangle: $\text{Base} = 3, \text{Height} = 4$ Apply elementary formula: $\text{Area} = \frac{1}{2} \times \text{Base} \times \text{Height} = \frac{1}{2} \times 3 \times 4 = 6$
<b>Incorrect reasoning process</b>	$\text{Area} = \frac{1}{2}(\text{Sum of sides}) = \frac{1}{2}(3 + 4 + 5) = 6$

Table 12: Case study of how the reasoning process is simplified.

Role	Prompt Template
<b>Critic</b>	<p>You are a highly capable evaluator.</p> <p>Your task is to assess the given reasoning process from the perspective of a small language model (e.g., 7B). Specifically, determine whether the reasoning process provides sufficient detail for a small model to solve the problem, or whether it is too terse (i.e., lacking critical details) or too complex (i.e., containing unnecessary or confusing steps).</p> <p>Complexity Definitions (from the perspective of a small model):</p> <ul style="list-style-type: none"> <li>- Easy: The reasoning process is overly terse; it omits essential details that a small model needs to solve the problem.</li> <li>- Medium: The reasoning process is appropriately balanced, offering enough detailed guidance.</li> <li>- Hard: The reasoning process is overly complex, with extraneous or convoluted steps that could hinder a small model to follow it.</li> </ul> <p>Output Format:</p> <p>You must output exactly one word: easy, medium, or hard.</p>
<b>Rethinker</b> (easy)	<p>You are a helpful assistant who is highly skilled at extending reasoning processes.</p> <p>Given a problem ,its correct answer and its terse reasoning process, your task is to extend the reasoning process by adding necessary details and intermediate steps so that a small language model (e.g., a 7B model) can follow the extended reasoning process to solve the problem.</p> <p>You should add necessary steps and details based on the correct answer.</p> <p>You must output ONLY the extended reasoning process with no additional explanation or commentary.</p>
<b>Rethinker</b> (hard)	<p>You are a helpful assistant who is highly skilled at simplifying reasoning processes.</p> <p>Given a problem, its correct answer and its overly complex reasoning process, your task is to simplify the reasoning process so that a small language model (e.g., a 7B model) can reliably follow the steps to solve the problem.</p> <p>You should remove redundancies or use simpler method on the basis of correct answer.</p> <p>You must output ONLY the simplified reasoning process with no additional explanation or commentary.</p>
<b>Verifier</b>	<p>You are a highly capable Verifier.</p> <p>Your task is to assess a given reasoning process based on a problem and its correct answer.</p> <p>Specifically, determine whether the reasoning process is sufficient and accurate for you to reach the correct answer.</p> <p>If the reasoning process correctly guides you to derive the the correct answer, output YES.</p> <p>If the reasoning process fails to guide you to the correct answer, output NO.</p> <p>You must output exactly one word: YES or NO.</p>
<b>Rethinker</b> (incorrect thought)	<p>You are an assistant who is skilled at converting correct reasoning processes to incorrect reasoning processes.</p> <p>Given a problem ,its answer and its correct reasoning process, your task is to corrupt the correct reasoning process by introducing logical fallacies and misleading steps, so that a small language model (e.g., a 7B model) cannot follow the incorrect reasoning process to solve the problem.</p> <p>You must output ONLY the incorrect reasoning process with no additional explanation or commentary.</p>

Table 13: Prompt templates for the CRV+CogPO framework.