

# Marco-o1 v2: Towards Widening The Distillation Bottleneck for Reasoning Models

Huifeng Yin<sup>1,2\*</sup> Yu Zhao<sup>1\*</sup> Minghao Wu<sup>1,3</sup> Xuanfan Ni<sup>1</sup> Bo Zeng<sup>1</sup> Hao Wang<sup>1</sup> Tianqi Shi<sup>1</sup>  
Liangying Shao<sup>1</sup> Chenyang Lyu<sup>1</sup> Longyue Wang<sup>1†</sup> Weihua Luo<sup>1†</sup> Kaifu Zhang<sup>1</sup>

<sup>1</sup>Alibaba International Digital Commerce

<sup>2</sup>Tsinghua University

<sup>3</sup>Monash University

{fengli.zy, wanglongyue.wly, weihua.luowh, kaifu.zkf}@alibaba-inc.com

## Abstract

Large Reasoning Models (LRMs) such as OpenAI o1 and DeepSeek-R1 have shown remarkable reasoning capabilities by scaling test-time compute and generating long Chain-of-Thought (CoT). Distillation post-training on LRMs-generated data is a straightforward yet effective method to enhance the reasoning abilities of smaller models, but faces a critical bottleneck: we found that distilled long CoT data poses *learning difficulty* for small models and leads to the *inheritance of biases* (i.e. over-thinking) when using Supervised Fine-tuning (SFT) and Reinforcement Learning (RL) methods. To alleviate this bottleneck, we propose constructing tree-based CoT data from scratch via Monte Carlo Tree Search (MCTS). We then exploit a set of CoT-aware approaches, including *Thoughts Length Balance*, *Fine-grained DPO*, and *Joint Post-training Objective*, to enhance SFT and RL on the constructed data. We conduct evaluation on various benchmarks such as math (GSM8K, MATH, AIME), instruction-following (Multi-IF) and planning (Blocksworld), results demonstrate our approaches substantially improve the reasoning performance of distilled models compared to standard distilled models via reducing the hallucinations in long-time thinking. The project homepage is <https://github.com/AIDC-AI/Marco-o1>.

## 1 Introduction

Recent advancements in large reasoning models (LRMs), such as OpenAI o1 (OpenAI, 2024), QwQ (Qwen Team, 2024) and DeepSeek-R1 (Guo et al., 2025), have led to significant progress in handling complex tasks spanning mathematics, coding, and even open-ended queries (Zhong et al., 2024b; Huang et al., 2024b; Zhao et al., 2024a). The success is largely attributed to “scaling test-time

compute” by extending the length of the reasoning process. Given that most state-of-the-art LRMs are computationally expensive, recent efforts attempt to distill their reasoning capabilities into smaller lightweight models, demonstrating competitive performances (Qin et al., 2024). For instance, Guo et al. (2025) explored *direct distillation*, where they fine-tuned smaller dense models (e.g. Qwen2.5 7B) using reasoning patterns generated by DeepSeek-R1 671B model, outperforming GPT-4 on math benchmarks (e.g. AIME: 9.3% vs. 55.5%).

However, we observed that these distilled models often exhibit hallucinations during long-time thinking, such as content repetition and over-reflection, leading to no final answer being produced (as shown in Table 1). We refer to this phenomenon as **formalistic long-time thinking**, where smaller models mechanically replicate the reasoning patterns of large models without internalizing the reasoning logic. Recent research shows that LRMs face both over-thinking and under-thinking issues (Chen et al., 2024; Wang et al., 2025b), while smaller models struggle to learn general reasoning (Fu et al., 2023a). Accordingly, the root cause may be that distillation methods introduce *bias inheritance* and *learning difficulties* in smaller models. A natural research question arise: *How can long CoT reasoning be effectively transferred to smaller models through data construction, SFT and RL methods?*

To tackle this challenge, we explore improvements in reasoning distillation from both data and methodological perspectives. First, we propose a fundamental framework for constructing tree-based CoT data, which generates pre-defined thought nodes using general LLMs (rather than LRMs) and heuristically expands these nodes into a tree structure via the Monte Carlo Tree Search (MCTS) algorithm (Browne et al., 2012). The constructed data is not only more effective compared to directly distilled data, but also inherently more flex-

\*Equal contribution.

†Corresponding author.

| Error Type  | Thoughts   |
|---|--|
| <b>Content Repetition</b><br>(Math)               | <p>Five men and nine women stand equally spaced around a circle in random order. The probability that every man stands diametrically opposite a woman is <math>\frac{m}{n}</math>, where <math>m</math> and <math>n</math> are relatively prime positive integers. Find <math>m + n</math>.</p> <p>So I've got this probability problem here. There are five men and nine women standing equally spaced around a circle in random order. I need to find the probability that ...</p> <p>...</p> <p>The total number of circular permutations of 14 distinct people is <math>(14 - 1)! = 13!</math>, because in circular arrangements, one position is fixed to handle the rotational symmetry. So, total number of possible arrangements is <math>13!</math>.</p> <p>...</p> <p>Wait, but in circular arrangements, positions are considered up to rotation, so fixing one position is standard. Alternatively, perhaps I need to consider that the <i>positions</i> are <b>considered</b> up to <b>consider</b> that the <i>positions</i> are <b>considered</b> up to that the <i>positions</i> are <b>considered</b> up to <b>consider</b> that the <i>positions</i> are <b>considered</b> up to ... (no answer)</p> |
| <b>Over-Reflection</b><br>(Instruction-Following) | <p>How to tell others that your major is computer science, without saying the word computer or science? Your entire response should not contain the word computer and science.</p> <p>So I have this problem: I need to explain what my major is without saying "computer" or "science". Hmm, my major is in computer science, but I have to find another way to describe it.</p> <p>...</p> <p>Maybe I can describe it in terms of what I do. I learn how to design and develop software, understand how information is processed and stored, and how to solve problems using logical methods. I also study the theory behind it, like algorithms and data structures.</p> <p>...</p> <p>I should also consider ... Alternatively, ... I think another angle is ... Wait, perhaps ... I could say that ... Or, I might ... I think another way ... Alternatively, ... I should also consider ... (no answer)</p>  |
| <b>Instruction Failure</b><br>(Translation)       | <p>Translate the following sentence into English: 您好，很高兴为您效劳，请问有什么可以帮您。</p> <p>Translation: <i>Okay</i>, I'm happy to help you. Do you have any questions for me?</p> <p><b>Wait, that's not correct.</b> The original sentence is in Chinese, and the translation should be accurate. Let me think again.</p> <p>The original sentence is: 您，很高兴为您效劳，请问有什么可以帮您。</p> <p>Breaking it down:<br/>"您好" is "Hello" or "Hi". "很高兴为您效劳" translates to "I am very pleased to serve you." ...</p> <p>Wait, the original sentence is more like a greeting and offer to help. ...</p> <p>Alternatively, "Hello, I'm happy to help you. How can I assist you?" So perhaps the correct translation is: "Hello, I'm here to help you. How can I assist you?"</p> <p>Wait, ...</p> <p>Alternatively, ... (no answer)</p>  |

Table 1: The illustration of **formalistic long-time thinking** generated by distilled reasoning models across different tasks. Error tokens in thoughts are highlighted in *blue* and *red* colors. Notably, due to excessively long thoughts, there are *no final answers* in above cases. The Quantitative Analysis is detailed in Section 2.4.

ible, allowing the extraction of different types of reasoning paths as training data. Secondly, regarding commonly-used SFT and direct preference optimization (DPO) (Rafailov et al., 2023) as post-training framework, we empirically investigate a set of CoT-aware methods on the effects of formalistic long-time thinking. Specifically, this includes: 1) *Thoughts Length Balance*, where we extract CoT data of varying lengths; 2) *Fine-grained DPO*, where we employ conservative DPO (cDPO) (Mitchell, 2023) and mask-based DPO to better leverage the fine-grained information in long CoT; 3) *Joint Post-training Objective*, where we combine the DPO loss with SFT loss to mitigate the over-optimization observed in DPO (Fernando et al.,

2024; Wang et al., 2024b).

We validated our approaches on five exam-oriented and open-ended benchmarks, covering three different difficulty levels of math (GSM8K, MATH and AIME) (Cobbe et al., 2021a; Lightman et al., 2023), instruction-following in eight languages (Multi-IF) (He et al., 2024), and real-world planning tasks (Blocksworld) (Valmeekam et al., 2022). Experimental results show that the proposed method consistently and orthogonally improve reasoning performance over the standard distilled models. The improvements come from the reduced hallucinations during long-time thinking, particularly content repetition, which leads to fewer “no answer” phenomena and better overall accuracy.

The **main contributions** of this work are:

- Our study reveals the side effect of standard distillation on transferring long CoT reasoning, which results in sub-optimal training of smaller models when using the distilled data (in Section 2.4).
- We propose a novel approach to construct CoT trees from scratch, which not only scales up the solution space but also more closely mimics human-like reasoning patterns. To the best of our knowledge, it is the first attempt of its kind (in Section 2).
- We investigate a set of effective approaches to widen the distillation bottleneck, demonstrating that they are orthogonal and complementary to each other, and robustly applicable to different reasoning tasks and languages (in Section 3&4).

## 2 Tree-Based CoT Data Construction

We propose a flexible and customizable tree-based CoT data construction method that generates high-quality CoT data from scratch. In this section, we introduce the overall framework in Section 2.1, followed by the thought nodes (Section 2.2), reasoning patterns (Section 2.3), and how to extract CoT data for post-training (Section 2.4).

### 2.1 Overall Framework

We introduce the tree-based CoT data construction process as shown in Figure 1. This tree structure not only constrains the search space to prevent unbounded expansions but also guides the model to produce reasoning steps (nodes) systematically. For instance, we specify that each node in the search tree corresponds to a particular action role (e.g., *thinking*, *reflection*), and each edge represents a transition to the next step. By constraining the transitions among these nodes, we ensure the search is both tractable and coherent. With the structure in place, we use MCTS to explore the search tree. During each step:

- **Node Selection.** We select a thought node to expand based on MCTS principles, such as upper confidence bound (UCB). If  $\text{Child}(n)$  denotes the set of child nodes of node  $n$ , then UCB balances exploration and exploitation via a score:

$$UCB(n_i) = \frac{v(n_i)}{n_{\text{visits}}(n_i)} + C \sqrt{\frac{\ln(n_{\text{visits}}(n_{\text{parent}}))}{n_{\text{visits}}(n_i)}},$$

where  $v(n_i)$  is an estimated value (or reward) of node  $n_i$ ,  $n_{\text{visits}}(\cdot)$  denotes the visit count, and  $C$  is the exploration constant.

| Thought Node | Prompt   |
|--------------|--|
| Thinking     | (continuation generation)  |
| Sub-Task     | Firstly, I need to break down this task.                             |
| Reflection   | Let’s check the result. Wait! something is wrong, let’s think again. |
| Hypothesis   | I propose the following hypothesis:                                  |
| Double-Check | Now, I need to check whether all the requirements are met.           |
| Reclarify    | To ensure clarity, let me restate the question or issue at hand:     |
| Answer       | The answer is:   |

Table 2: The pre-defined Thought Node. For a selected node, its corresponding prompt is continuously fed to LLMs for MCTS expansion.

- **Expansion.** We expand the selected node by prompting an LLM by adding a thought prompt (detailed in Table 2) that specifies the required action role. The LLM then generates the textual content for that node.
- **Rollout.** If the expansion reaches an *answer* node, we compute a reward based on correctness determined by rules and backpropagate this reward up the tree.

### 2.2 Thought Node

**Definition** A Thought Node corresponds to a distinct step or action within the CoT reasoning process. As outlined in Table 2, each node has a dedicated role and prefix prompt that guides the language model to generate specific content or revise previously generated reasoning. This structured design facilitates modular expansion and systematic backtracking within the MCTS framework. Notably, *Thinking* is treated as a special node that does not require any prefix prompt; instead, it admits unconditioned continuation generation to foster open-ended exploration of partial solutions. By combining multiple node types into a coherent tree, we can more effectively elicit and refine multi-step reasoning from the model.

**Multi-Model Coordination and Reflection** We adopt multi-model coordination to further diversify and correct the generated reasoning paths: 1) For nodes such as *Thinking*, we use Qwen2.5-72B-Instruct to generate logical steps or partial solutions; 2) For *Reflection* nodes, we switch to a different model, e.g., Llama3.1-70B-Instruct, to perform self-checks and corrections.

This separation enhances the reliability of reflec-

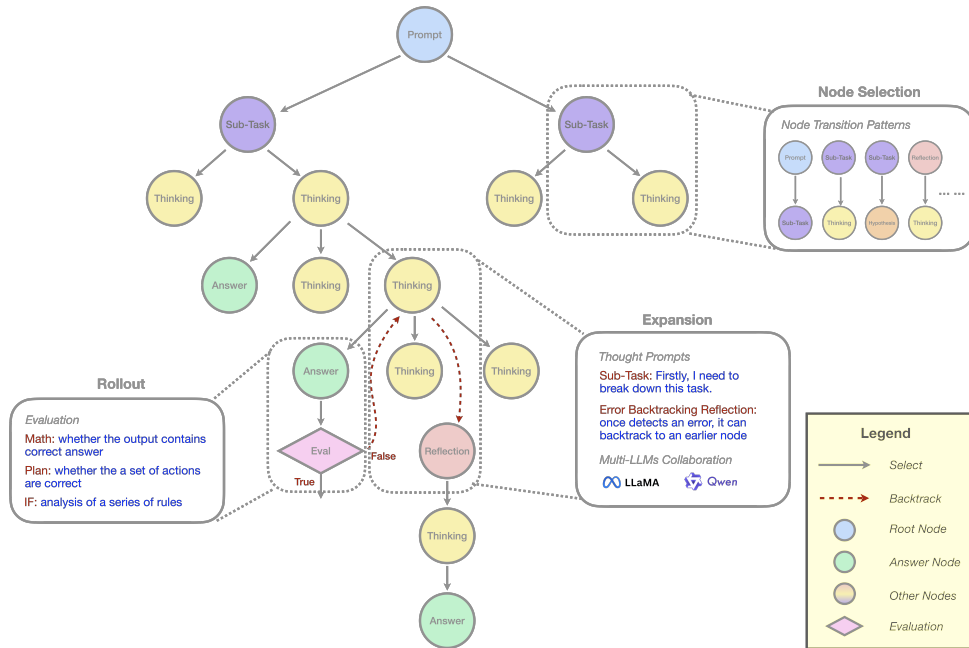


Figure 1: MCTS-based CoT data generation framework. Starting from an initial prompt (root node), the system proceeds through predefined nodes (e.g., *Sub-Task*, *Thinking*, *Reflection*) according to a customizable node transfer matrix. Each node is expanded by prompting either Qwen or Llama, allowing multi-model collaboration. If a wrong answer is detected, we perform error backtracking (pink arrows) to a prior node and trigger *Reflection* in another model, enhancing the overall correctness and diversity of the final reasoning path.

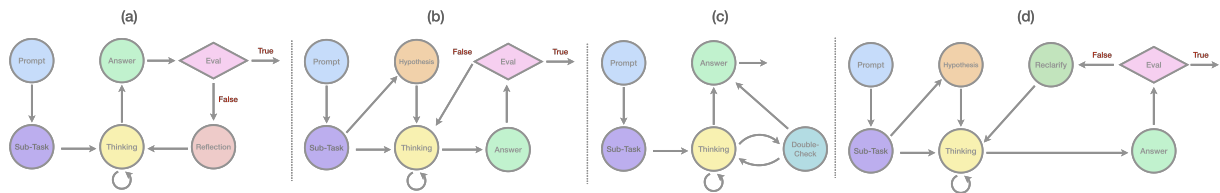


Figure 2: Representative node transition patterns in our search tree. Each sub-figure (a–d) illustrates a distinct sequence of transitions (e.g., *Sub-Task*, *Thinking*, *Reflection*, *Double Check*, *Hypothesis*) toward arriving at an *Answer* node. These variations allow the search to adaptively expand or backtrack based on correctness checks, thereby generating rich and context-specific chain-of-thought data.

tion. When the same model that made a mistake also attempts to correct itself, it may fall back on the same erroneous distributional patterns. In our pipeline, if a *Reflection* node detects an error, it can backtrack to a specific earlier node (also configurable in the MCTS design) and request a re-generation of the *Thinking* steps. By alternating between models, we reduce the risk of repeated mistakes and improve the diversity of exploration.

### 2.3 Reasoning Pattern

As illustrated in Figure 2, we design a set of *customizable search tree structures* to reflect the diverse ways in which humans reason about different tasks. Each tree is configured to capture a variety of reasoning modes. For instance, in Figure 2(a), we demonstrate a sequence of nodes to solve a ques-

tion: we first break down the task via a *Sub-Task* node, then perform a general *Thinking* step, and finally provide an *Answer*. We evaluate correctness through rule-based checks: if correct, we output the result; otherwise, we prompt the model to reflect on potential mistakes (entering the *Reflection* node). Here, the model revisits or re-checks its chain of thought, then either formulates new reasoning or proposes a revised answer. This feedback loop repeats until the solution is correct or a pre-set search limit is reached. Some tasks, however, also require formulating assumptions or provisional conclusions, so in Figure 2(b), we incorporate a *Hypothesis* node immediately after the *Sub-Task* node for tasks that benefit from explicitly positing assumptions or preliminary formulas early on. For example, “Find the sum of all ordered pairs  $(x, y)$

of positive integers such that  $x + y = 5$ ." After breaking down the task (*Sub-Task*), the model proposes a hypothesis that  $x$  can range from 1 to 4 (with  $y = 5 - x$ ) and enumerates each pair, obtaining a total sum of 20. If the answer is verified correct, the model outputs 20.

In practice, we randomly sample from these different reasoning-flow templates (e.g., Figure 2(a-d)) to ensure we capture diverse, human-like modes of thought. By introducing node transition patterns, our framework produces richer and more flexible CoT data and fosters more robust reasoning in downstream tasks.

## 2.4 CoT Data Extraction for Post-Training

Once MCTS completes its exploration, we have a large set of candidate paths. At this point, we must extract final CoT data for SFT or DPO.

- **CoT Data for SFT.** We select successful paths that lead to the correct final answer. Depending on the data volume requirements, one can: 1) Pick the highest-reward path according to MCTS; 2) Pick the long or short path that yields the correct answer, if specific chain lengths are desired.
- **CoT Data for DPO.** Constructing DPO data requires both positive and negative examples for each prompt: 1) The positive example is the CoT path that correctly solves the problem, like the SFT data; 2) The negative example is a flawed path (an incorrect final answer) that shares a minimal prefix with the positive path to mitigate excessive overlapping tokens, which can degrade DPO performance.

We find that many existing QA prompts (especially those frequently seen during Qwen or Llama training) are too easy for the models, producing few or no negative paths. As a result, fewer DPO pairs are generated. One can overcome this limitation by using more challenging questions or those that the models have not encountered extensively.

**Quantitative Analysis: Formalistic Long-time Thinking** We explore the impact of the CoT length distribution on model performance. We experiment with different sampling strategies for DPO datasets, where responses to the same question are ranked and categorized based on their length (Their distributions are shown in Figure 3). As shown in Table 3, the selection of shorter CoT paths leads to a noticeable reduction in ineffective outputs. In addition, we note that shorter reasoning paths tend to mitigate the issue of "formalistic

| Datasets | Long   | Middle        | Short         |
|----------|--------|---------------|---------------|
| GSM8K    | 5.38%  | <b>5.08%</b>  | <b>5.08%</b>  |
| MATH     | 28.40% | 20.60%        | <b>16.20%</b> |
| AIME     | 51.66% | <b>50.00%</b> | 60.00%        |
| Plan.    | 6.40%  | 6.40%         | <b>6.20%</b>  |
| IF (Zh)  | 32.30% | 4.23%         | <b>1.9%</b>   |
| IF (En)  | 22.36% | <b>3.80%</b>  | 4.00%         |
| IF (Ot.) | 18.69% | 1.70%         | <b>1.59%</b>  |

Table 3: Effects of thoughts length(long, medium, and short CoT paths) on model performance across different datasets.

long-time thinking", thus improving the quality of the reasoning output.

## 3 CoT-Aware Post-Training

Section 2.4 identifies that DPO training is prone to causing the formalistic long-time thinking. In this section, we propose three methods to address this problem: Thoughts Length Balance (Section 3.1), Fine-grained DPO (Section 3.2), Joint Post-training Objective (Section 3.3).

### 3.1 Thoughts Length Balance

Section 2.4 illustrates that the length of Chain-of-Thought (CoT) reasoning significantly influences the reasoning performance of distilled smaller models during the Direct Preference Optimization (DPO) phase. In contrast, our preliminary experiments did not indicate such an effect during SFT. Therefore, we propose using the longest CoT data in the SFT phase and the shortest CoT data in the DPO phase.

Specifically, we extract all valid reasoning paths that lead to correct answers from the CoT trees, as multiple correct paths often exist. From these, we select paths categorized by their relative length(short, medium, and long) based on token count, correct paths serve as positive examples, whereas negative examples are generated by identifying incorrect paths that share the shortest common prefix with their corresponding positive paths.

We recognize that "short" and "long" are inherently relative and that CoT lengths vary significantly with problem complexity. For instance, a simple arithmetic problem like "1+1=?" naturally involves a shorter CoT compared to more complex integrals such as " $\int_0^1 \frac{\ln(x+1)}{x^2+1} dx = ?$ ". To address these variations systematically, we employ Monte

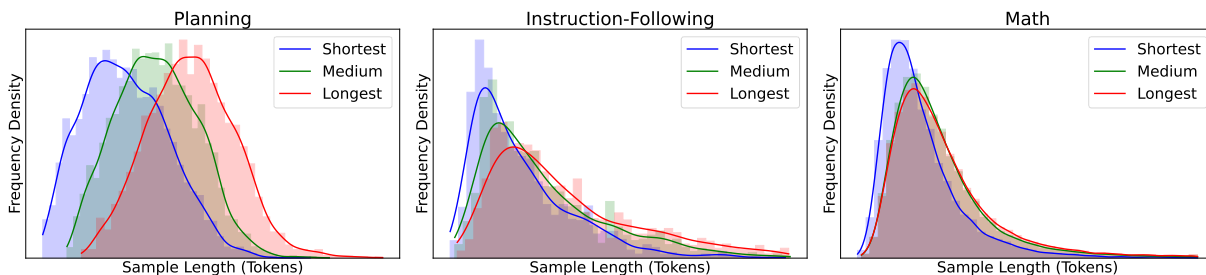


Figure 3: Distribution of data lengths by token count: A comparative analysis of sampling strategies and their correspondence to short, medium, and long data lengths illustrated with histograms and KDE curves

|         |         |   |   |      |    |       |      |      |      |       |      |        |      |      |      |
|---------|---------|---|---|------|----|-------|------|------|------|-------|------|--------|------|------|------|
| Chosen: | Alright | , | I | need | to | break | down | this | task | ..... | The  | answer | is:  | 23   | cm   |
| Logp:   | 0       | 0 | 0 | 0    | 0  | 0     | 0    | 0    | 0    | ..... | 3.39 | 10.51  | 5.11 | 9.12 | 4.65 |
| Reject: | Alright | , | I | need | to | break | down | this | task | ..... | So   | answer | is:  | 28   | cm   |
| Logp:   | 0       | 0 | 0 | 0    | 0  | 0     | 0    | 0    | 0    | ..... | 2.41 | 6.21   | 2.15 | 4.32 | 5.72 |

} DPO Loss

Figure 4: The illustration of masking-based DPO, by setting the log probabilities of the common prefix tokens in preference pairs to zero.

Carlo Tree Search (MCTS) to sample multiple reasoning paths per query and select representative paths by their relative lengths, rather than imposing rigid token-count thresholds.

### 3.2 Fine-grained DPO

Recent studies highlight that DPO is sensitive to the length of responses, which can lead to biased reward assessments (Lu et al., 2024; Liu et al., 2024). Longer chosen responses increase the model’s tendency to generate longer outputs, while longer rejected responses push the model to move away from such outputs, potentially without reducing their length. These issues are particularly pronounced in long CoT reasoning tasks, where length disparities may undermine DPO’s effectiveness in fine-tuning reasoning models.

**Conservative DPO** Conservative DPO (cDPO) (Mitchell, 2023) adapts the standard DPO framework to handle noisy preference labels, typically encountered when labels may be flipped with a small probability  $\epsilon$ . The key innovation of cDPO is to modify the target distribution to account for potential label noise, setting the preference probability to  $p(y_w \succ y_l) = 1 - \epsilon$ . This adjustment reduces the impact of noisy labels by softening the gradient updates, making the model less sensitive to incorrect preferences. Formally, the cDPO loss is defined as:

$$\mathcal{L}_{\text{DPO}}^\epsilon(\theta, y_w, y_l) = -(1 - \epsilon) \log \hat{p}_\theta(y_w \succ y_l) - \epsilon \log(1 - \hat{p}_\theta(y_w \succ y_l)),$$

where  $\hat{p}_\theta(y_w \succ y_l)$  is the predicted preference probability. The gradient of this loss combines weighted contributions from both the correct and incorrect paths, facilitating more stable training under label noise. By upweighting correct preferences and downweighting incorrect ones, cDPO prevents overfitting to noisy data, resulting in more reliable optimization and improved model robustness.

**Masking-based DPO** To mitigate the adverse effects of shared prefixes inherent in tree-based data, we modify the DPO loss computation by masking out the shared prefix tokens. Specifically, prior to loss calculation, we identify the number of tokens constituting the common prefix between the positive and negative samples and adjust the loss mask by setting the corresponding entries for these shared tokens to zero-analogous to the treatment of padding tokens in standard loss formulations, as shown in Figure 4. This ensures that the shared prefix tokens do not contribute to the gradient computation, allowing the model to focus on the differentiating segments of the outputs and better distinguish between valid and invalid reasoning paths. This strategy provides a fine-grained adjustment to the DPO objective, enhancing optimization in settings with substantial prefix overlap.

### 3.3 Joint Post-training Objective

In model training, the typical approach follows a sequential training paradigm, first conducting SFT followed by RLHF or DPO. However, this sequential training process is suboptimal due to the inherent trade-off between SFT and RLHF/DPO, where the model tends to forget the content learned in the first stage as it progresses to the second. Even regularization methods like KL divergence cannot fully mitigate the forgetting caused by the distribution shift from the SFT dataset to the preference-based dataset, as highlighted by (Fernando et al., 2025). A similar phenomenon is observed in our work,

| Model                        | Math        |             |             | Planning    | Instruction-Following |             |             |
|------------------------------|-------------|-------------|-------------|-------------|-----------------------|-------------|-------------|
|                              | GSM8K       | MATH        | AIME        | Blocksworld | Zh                    | En          | Other       |
| <b>Llama-3.1-8B-Instruct</b> | 85.5        | 47.0        | 11.7        | 10.0        | 61.5                  | 76.2        | 67.1        |
| + Sky-T1                     | 84.8        | 44.0        | 6.7         | 2.0         | 25.4                  | 31.6        | 29.7        |
| + Our Data                   | <b>87.4</b> | <b>51.4</b> | <b>15.0</b> | <b>12.4</b> | <b>69.2</b>           | <b>76.6</b> | <b>79.1</b> |
| <b>Llama-3.2-1B-Instruct</b> | 43.4        | 33.0        | 0.0         | 0.2         | 31.9                  | 47.8        | 33.1        |
| + Sky-T1                     | 37.2        | 19.0        | 0.0         | 0.0         | 8.1                   | 10.0        | 7.5         |
| + Our Data                   | <b>50.0</b> | <b>38.6</b> | <b>0.0</b>  | <b>5.6</b>  | <b>41.1</b>           | <b>50.6</b> | <b>47.1</b> |
| <b>Qwen2.5-7B-Instruct</b>   | 90.4        | 62.0        | 15.0        | 10.6        | 69.6                  | 72.8        | 74.4        |
| + Sky-T1                     | 89.6        | 61.6        | 9.4         | 0.4         | 26.2                  | 24.5        | 30.6        |
| + Our Data                   | <b>90.7</b> | <b>64.0</b> | <b>15.0</b> | <b>12.0</b> | <b>73.1</b>           | <b>73.4</b> | <b>78.8</b> |
| <b>Qwen2.5-1.5B-Instruct</b> | 67.5        | 38.4        | 0.0         | 1.0         | 49.6                  | 43.8        | 41.0        |
| + Sky-T1                     | 66.7        | 34.4        | 0.0         | 0.0         | 9.2                   | 11.8        | 14.8        |
| + Our Data                   | <b>74.6</b> | <b>46.8</b> | <b>0.0</b>  | <b>5.4</b>  | <b>51.5</b>           | <b>51.6</b> | <b>54.8</b> |

Table 4: Comparison of SFT results across various models on multiple different tasks, including math, planning and instruction-following Benchmarks.

where such forgetting contributes to the emergence of formalistic long-time thinking in distilled models. To address this, we introduce SFT loss during the DPO training phase to alleviate the performance degradation resulting from the switch in training methodologies. The final loss function is thus modified as:  $\mathcal{L} = \mathcal{L}_{\text{DPO}} + \alpha \mathcal{L}_{\text{SFT}}$ , where the hyperparameter  $\alpha$  enables a better trade-off between SFT and preference learning, helping to maintain consistency in the model’s reasoning patterns throughout the training process. This adjustment ensures more robust and stable performance across stages.

## 4 Experiments

### 4.1 Experimental Setup

**Models** We start with the baseline model, "Our LRM (SFT)," which is Llama-3.1-8B fine-tuned on our CoT data. Direct Preference Optimization (DPO) is applied next, followed by Data Length Balance. Conservative DPO (cDPO) is then added, and a Joint Loss function combining DPO and Supervised Fine-Tuning (SFT) loss is incorporated. Finally, masking-based DPO is applied. Each of these methods is sequentially added to the baseline, as shown in Table 5.

**Benchmark** We evaluate our approach on five benchmarks, each capturing different reasoning challenges. AIME focuses on higher-level math with 60 questions from 2023 and 2024, while GSM8K (Cobbe et al., 2021a) features elementary-to-intermediate arithmetic tasks. MATH500 (Lightman et al., 2023) presents a wide range of advanced

mathematical problems, testing deeper analytical thinking. For sequential decision making, we adopt the classical Blocksworld (Valmeekam et al., 2022) planning domain from the International Planning Competitions (IPC). Lastly, Multi-IF (He et al., 2024) assesses multi-turn instruction following in eight languages, encompassing 4,501 multilingual, three-turn conversations.

### 4.2 Experimental Results

**Constructed Data Validation** We apply our constructed CoT data to smaller-scale models from the Llama and Qwen families, comparing the results against the Sky-T1 dataset, which employs a QwQ-based distillation pipeline effective on larger models (32B). While Sky-T1 demonstrates competitive performance on large models, it faces challenges when scaled down to 8B models due to the inherent limitations in context processing and reasoning capabilities. In contrast, our CoT data, specifically designed to address these limitations, leads to substantial improvements in smaller models, particularly in tasks involving arithmetic reasoning such as GSM8K and MATH, as well as more complex open-ended tasks like AIME and Blocksworld, as shown in Table 4. This validates the effectiveness of our constructed data in advancing the performance of small models across a wide range of reasoning tasks.

**Main Results** As shown in Figure 5, we progressively adding various techniques described in Section 3 to address the challenges identified during DPO. Initially, we observe that DPO causes a

| Model              | Math                     |                                 |                                  | Planning                        | Instruction-Following    |                          |                                 |
|--------------------|--------------------------|---------------------------------|----------------------------------|---------------------------------|--------------------------|--------------------------|---------------------------------|
|                    | GSM8K                    | MATH                            | AIME                             | Blocksworld                     | Zh                       | En                       | Other                           |
| <i>Baseline</i>    |                          |                                 |                                  |                                 |                          |                          |                                 |
| Our LRM (SFT)      | 87.4<br>( <u>0.23%</u> ) | <b>51.4</b><br>( <b>5.40%</b> ) | <b>15.0</b><br>( <b>30.00%</b> ) | <u>12.4</u><br>( <b>1.80%</b> ) | 69.2<br>( <b>0.77%</b> ) | 76.6<br>( <b>1.69%</b> ) | <b>79.1</b><br>( <b>1.08%</b> ) |
| + DPO              | 86.2<br>(6.37%)          | 41.8<br>(31.80%)                | 8.3<br>(55.00%)                  | 2.0<br>(93.60%)                 | 5.7<br>(91.54%)          | 6.3<br>(90.93%)          | 6.7<br>(92.22%)                 |
| <i>Our Methods</i> |                          |                                 |                                  |                                 |                          |                          |                                 |
| + Data Balance     | 86.8<br>(5.08%)          | 28.0<br>(46.40%)                | 6.6<br>(65.00%)                  | 6.8<br>(44.60%)                 | 43.4<br>(30.77%)         | 44.7<br>(44.73%)         | 42.4<br>(45.28%)                |
| + cDPO             | <b>87.5</b><br>(3.71%)   | 48.6<br>(15.00%)                | <b>15.0</b><br>(45.00%)          | 4.4<br>(47.40%)                 | 61.9<br>(11.15%)         | 66.4<br>(15.61%)         | 67.7<br>(15.40%)                |
| + Joint Loss       | 86.8<br>(0.38%)          | 48.6<br>(8.60%)                 | <u>10.0</u><br>(31.67%)          | 8.6<br>(9.00%)                  | <b>72.3</b><br>(1.15%)   | <b>78.9</b><br>(1.90%)   | <u>78.1</u><br>(2.22%)          |
| + Masking          | 87.2<br>( <b>0.15%</b> ) | <u>51.0</u><br>(5.80%)          | 8.0<br>(38.33%)                  | <b>12.6</b><br>(10.20%)         | <u>72.0</u><br>(1.15%)   | <u>77.2</u><br>(1.90%)   | <b>79.1</b><br>(1.36%)          |

Table 5: Performance comparison among different methods. The best performance is boldfaced, while the second best is underlined. The numbers in parentheses indicates the ratio of instances where no answer is obtained in the specified format.

significant increase in output length, which results in a marked drop in model performance due to the high proportion of samples without answers. To mitigate this issue, we explore several strategies, including the data balance, applying cDPO to reduce the impact of noisy labels, SFT Loss for multi-objective training to prevent catastrophic forgetting, and masking shared prefixes during loss calculation to reduce overemphasis on redundant tokens. Our results show that these adjustments achieves a notable improvement in reasoning tasks, particularly in planning and instruction-following, while maintaining competitive performance on mathematical benchmarks. The performance improvements can be primarily attributed to the reduction of formalistic long-time thinking, which is a major source of inefficiency in reasoning. By addressing this issue, our model exhibits a stronger ability to generate meaningful and correct answers instead of producing excessive, irrelevant reasoning steps. This leads to a substantial enhancement in overall model effectiveness, with improvements of coherent reasoning and accurate outputs. The integration of these methods ensures that our model achieves robust and efficient performance across a wide range of reasoning tasks, Contributing to the application of DPO technology in LRMs.

### 4.3 Effects of Joint DPO Loss

In this study, we explore the effects of combining DPO loss with SFT loss within a joint post-training

| Datasets | CDPO                    | +0.5                    | +1.0                   | +1.5                   | +2.0                   |
|----------|-------------------------|-------------------------|------------------------|------------------------|------------------------|
| GSM8K    | <b>87.5</b><br>(3.71%)  | 86.5<br>(0.53%)         | 86.8<br>(0.38%)        | 85.5<br>(0.08%)        | 85.6<br>(0.08%)        |
| MATH     | <u>48.6</u><br>(15.00%) | <b>50.0</b><br>(10.20%) | <u>48.6</u><br>(8.60%) | 48.4<br>(0.08%)        | 48.0<br>(9.00%)        |
| AIME     | <b>15.0</b><br>(45.00%) | 11.6<br>(35.00%)        | 10.0<br>(31.67%)       | 6.6<br>(36.67%)        | 11.6<br>(31.67%)       |
| Plan.    | 4.4<br>(47.40%)         | 7.8<br>(12.80%)         | <b>8.6</b><br>(9.00%)  | 7.6<br>(6.40%)         | 8.4<br>(7.40%)         |
| IF (Zh)  | 61.9<br>(11.15%)        | 68.8<br>(2.31%)         | <b>72.3</b><br>(1.15%) | 68.4<br>(3.46%)        | <u>70.7</u><br>(0.77%) |
| IF (En)  | 66.4<br>(15.61%)        | 76.1<br>(4.22%)         | <b>78.9</b><br>(1.90%) | 77.2<br>(2.74%)        | <u>78.2</u><br>(2.32%) |
| IF (Ot.) | 67.7<br>(15.40%)        | 78.0<br>(1.99%)         | 78.1<br>(2.22%)        | <b>79.2</b><br>(1.82%) | <u>78.9</u><br>(1.93%) |

Table 6: Effects of joint loss (combining DPO loss and SFT loss with a weight factor  $\alpha$  in  $\mathcal{L} = \mathcal{L}_{\text{DPO}} + \alpha \mathcal{L}_{\text{SFT}}$ ) on model performance across different datasets with varying hyperparameter settings.

objective, as outlined in Section 3.4. Our goal is to stabilize the training process and mitigate issues like over-optimization observed in pure DPO. To this end, we experiment with varying values of the hyperparameter alpha, which controls the weight balance between the DPO and SFT losses. As shown in Table 6, our results indicate that alpha=1 provides the best trade-off, as smaller values still lead to some degree of catastrophic forgetting, while larger values reduce the effectiveness of preference alignment, thereby diminishing the efficiency of the valuable preference dataset. Consequently, we adopt the combined  $\mathcal{L} = \mathcal{L}_{\text{DPO}} + \mathcal{L}_{\text{SFT}}$  as the configuration for subsequent experiments.



#### 4.4 MCTS Inference Exploration

As an additional experiment, we explored the impact of applying MCTS at the inference stage. As shown in Table 7, we use Test@N to denote the percentage of problems solved correctly at least once when allowing the model to make N separate guesses for each problem (Cobbe et al., 2021b). We evaluated solve rates at Test@1, Test@8, and Test@32 on the MATH dataset. Specifically, at Test@8 and Test@32, the MCTS-based approach outperforms the model without MCTS inference, demonstrating its ability to expand the solution space and leverage test-time scaling effectively.

### 5 Related Work

**Reasoning Models** Recent advancements in reasoning models, like OpenAI o1 (OpenAI, 2024), DeepSeek-R1 (DeepSeek-AI et al., 2025), and Qwen QwQ (Qwen Team, 2024), make significant strides in complex reasoning tasks through CoT generation and increased test-time compute. These models scale reasoning depth by expanding their thinking processes, generating step-by-step solutions to problems, which significantly boosts performance in domains such as mathematics and coding (Zhong et al., 2024a; DeepSeek-AI et al., 2025). However, the challenge remains that these advancements largely depend on large model sizes, and smaller models struggle to replicate this reasoning behavior (Fu et al., 2023b).

**Knowledge Distillation** Researchers focus on transferring the reasoning capabilities of LRMs into smaller models, through distillation techniques. Distillation methods, such as fine-tuning smaller models on CoT data generated by larger models (Huang et al., 2024a), show that small models can benefit from reasoning data generated by large models. However, recent findings reveal that small models often fail to capture intricate reasoning patterns due to their limited capacity, leading to suboptimal performance when directly distilled from large models (Li et al., 2025a; Wang et al., 2025a).

**Monte Carlo Tree Search** MCTS has been proposed as a solution to improve reasoning by exploring multiple reasoning paths during inference (Zhao et al., 2024b; Tian et al., 2024). Moreover, MCTS provides a powerful mechanism for generating high-quality, diverse reasoning data that can subsequently be harnessed to fine-tune reasoning models. For instance, (Tian et al., 2024)

| Model                 | Test@1 | Test@8 | Test@32     |
|-----------------------|--------|--------|-------------|
| Llama-3.1-8B-Instruct | 47.0   | 67.6   | <b>75.8</b> |
| Our Best Model        | 51.0   | 70.2   | <b>79.2</b> |
| + MCTS Decode         | 51.0   | 70.8   | <b>82.8</b> |

Table 7: Performance on MATH Dataset: Test@1, Test@8, and Test@32 Results. Test@N denotes the percentage of problems solved correctly at least once when the model is allowed to make N separate guesses for each problem.

employs MCTS to synthesize candidate reasoning paths that capture varied solution strategies. Similarly, the RStar introduced in (Qi et al., 2024) utilizes MCTS to construct structured data, ensuring that generated reasoning chains are both comprehensive and coherent. In parallel, Math-Shepherd (Wang et al., 2024a) and OmegaPRM (Luo et al., 2024) employ MCTS to collect high-quality reasoning data, which is subsequently used to train PRM. Extending these ideas, (Zhang et al., 2024) combines the strengths of MCTS and PRM to guide policy updates.

### 6 Conclusion

We explore strategies to transfer long CoT reasoning to smaller models, addressing learning difficulty and bias inheritance in distillation. We propose a MCTS framework that independently generates flexible, tree-based reasoning data, reducing reliance on large teacher models. Enhanced by CoT-aware post-training, our approach effectively mitigates overly formalistic long-time thinking. Experiments across diverse tasks (math, planning and instruction-following) demonstrate robust performance improvements. These findings highlight the importance of well-designed data and post-training strategies in improving the efficiency and reliability of smaller-scale reasoning models. In the future, we will continue to explore how reasoning techniques can improve a broader range of non-mathematical tasks, including machine translation (Liu et al., 2025), multilingual and multimodal reasoning (Zeng et al., 2025; Li et al., 2025b).

### References

Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfschagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. 2012. A survey of monte carlo tree search methods. *IEEE Transac-*

- tions on Computational Intelligence and AI in games, 4(1):1–43.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, et al. 2024. Do not think that much for  $2+3=?$  on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187*.
- 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. 2021a. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- 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. 2021b. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- DeepSeek-AI, Daya Guo, and et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Heshan Fernando, Han Shen, Parikshit Ram, Yi Zhou, Horst Samulowitz, Nathalie Baracaldo, and Tianyi Chen. 2025. Mitigating forgetting in llm supervised fine-tuning and preference learning. *Preprint*, arXiv:2410.15483.
- Heshan Devaka Fernando, Han Shen, Parikshit Ram, Yi Zhou, Horst Samulowitz, Nathalie Baracaldo, and Tianyi Chen. 2024. Mitigating forgetting in LLM supervised fine-tuning and preference learning. *CoRR*, abs/2410.15483.
- Y. Fu, L. Zhang, and M. Liu. 2023a. Specializing large language models for reasoning tasks. *International Journal of AI*, 56:134–145.
- Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023b. Specializing smaller language models towards multi-step reasoning. *Preprint*, arXiv:2301.12726.
- Z. Guo, T. Shi, and K. Zhang. 2025. Deepseek-r1: Scaling test-time compute for enhanced reasoning. *Machine Learning Journal*, 72:25–45.
- Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, et al. 2024. Multi-if: Benchmarking llms on multi-turn and multilingual instructions following. *arXiv preprint arXiv:2410.15553*.
- Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. 2024a. O1 replication journey – part 2: Surpassing o1-preview through simple distillation, big progress or bitter lesson? *Preprint*, arXiv:2411.16489.
- Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei Qin, Weizhe Yuan, and Pengfei Liu. 2024b. O1 replication journey—part 2: Surpassing o1-preview through simple distillation, big progress or bitter lesson? *arXiv preprint arXiv:2411.16489*.
- Yuetai Li, Xiang Yue, Zhangchen Xu, Fengqing Jiang, Luyao Niu, Bill Yuchen Lin, Bhaskar Ramasubramanian, and Radha Poovendran. 2025a. Small models struggle to learn from strong reasoners. *Preprint*, arXiv:2502.12143.
- Yunxin Li, Zhenyu Liu, Zitao Li, Xuanyu Zhang, Zhenran Xu, Xinyu Chen, Haoyuan Shi, Shenyan Jiang, Xintong Wang, Jifang Wang, Shouzheng Huang, Xinpeng Zhao, Borui Jiang, Lanqing Hong, Longyue Wang, Zhuotao Tian, Baoxing Huai, Wenhan Luo, Weihua Luo, Zheng Zhang, Baotian Hu, and Min Zhang. 2025b. Perception, reason, think, and plan: A survey on large multimodal reasoning models. *arXiv preprint arXiv:2505.04921*.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*.
- Sinuo Liu, Chenyang Lyu, Minghao Wu, Longyue Wang, Weihua Luo, Kaifu Zhang, and Zifu Shang. 2025. New trends for modern machine translation with large reasoning models. *arXiv preprint arXiv:2503.10351*.
- Wei Liu, Yang Bai, Chengcheng Han, Rongxiang Weng, Jun Xu, Xuezhi Cao, Jingang Wang, and Xunliang Cai. 2024. Length desensitization in direct preference optimization. *Preprint*, arXiv:2409.06411.
- Junru Lu, Jiazheng Li, Siyu An, Meng Zhao, Yulan He, Di Yin, and Xing Sun. 2024. Eliminating biased length reliance of direct preference optimization via down-sampled kl divergence. *Preprint*, arXiv:2406.10957.
- Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. 2024. Improve mathematical reasoning in language models by automated process supervision. *Preprint*, arXiv:2406.06592.
- Eric Mitchell. 2023. A note on dpo with noisy preferences & relationship to ipo.
- OpenAI. 2024. Learning to reason with llms. <https://openai.com/index/learning-to-reason-with-llms/>. [Accessed 19-09-2024].
- Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lina Zhang, Fan Yang, and Mao Yang. 2024. Mutual reasoning makes smaller llms stronger problem-solvers. *Preprint*, arXiv:2408.06195.

- Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector Liu, Yuanzhi Li, et al. 2024. O1 replication journey: A strategic progress report—part 1. *arXiv preprint arXiv:2410.18982*.
- Qwen Team. 2024. [Qwq: Reflect deeply on the boundaries of the unknown](#).
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems*, volume 36, pages 53728–53741. Curran Associates, Inc.
- Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. 2024. [Toward self-improvement of llms via imagination, searching, and criticizing](#). *Preprint*, arXiv:2404.12253.
- Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2022. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. In *Neural Information Processing Systems*.
- Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. 2024a. [Math-shepherd: Verify and reinforce llms step-by-step without human annotations](#). *Preprint*, arXiv:2312.08935.
- Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025a. [Thoughts are all over the place: On the underthinking of o1-like llms](#). *Preprint*, arXiv:2501.18585.
- Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, et al. 2025b. Thoughts are all over the place: On the underthinking of o1-like llms. *arXiv preprint arXiv:2501.18585*.
- Zhichao Wang, Bin Bi, Zixu Zhu, Xiang-Bo Mao, Jun Wang, and Shiyu Wang. 2024b. [UFT: unifying fine-tuning of SFT and RLHF/DPO/UNA through a generalized implicit reward function](#). *CoRR*, abs/2410.21438.
- Bo Zeng, Chenyang Lyu, Sinuo Liu, Mingyan Zeng, Minghao Wu, Xuanfan Ni, Tianqi Shi, Yu Zhao, Yefeng Liu, Chenyu Zhu, Ruizhe Li, Jiahui Geng, Qing Li, Yu Tong, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2025. Marco-Bench-MIF: On multi-lingual instruction-following capability of large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*.
- Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. [Rest-mcts\\*: Llm self-training via process reward guided tree search](#). *Preprint*, arXiv:2406.03816.
- Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi Shi, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2024a. [Marco-o1: Towards open reasoning models for open-ended solutions](#). *arXiv preprint arXiv:2411.14405*.
- Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi Shi, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2024b. [Marco-o1: Towards open reasoning models for open-ended solutions](#). *Preprint*, arXiv:2411.14405.
- Tianyang Zhong, Zhengliang Liu, and et al. 2024a. [Evaluation of openai o1: Opportunities and challenges of agi](#). *Preprint*, arXiv:2409.18486.
- Y. Zhong, R. Li, and X. Wang. 2024b. Evaluation of large reasoning models. *Journal of AI Research*, 64:1–15.