

# ELICITING DIVERSE THINKING SCHEMATA FOR LARGE REASONING MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large reasoning models (LRMs) have attracted increasing attention for their ability to solve complex mathematical problems by generating extended reasoning chains. In this work, we highlight a critical yet underexplored aspect of their reasoning process—thinking schemata, which we define as the distinct transitions between reasoning steps and the variety of solution paths the model produces. We observe a correlation between the diversity of thinking schemata and model performance, which motivates us to enhance diversity as a means to further improve reasoning potential and generalization ability. To this end, we propose **Diverse Schemata Policy Optimization (DiScO)**, a method to elicit diverse thinking schemata by first endowing the model with the capabilities to be aware of the thinking schemata in its reasoning chain and then encouraging their diversity through reinforcement learning. Experiments on multiple mathematical reasoning benchmarks demonstrate that DiScO consistently outperforms standard group relative policy optimization, with particularly pronounced gains on challenging datasets such as AIME, [where our 7B and 32B DiScO models surpass the closed-source frontier LRMs by 2.9%-8.5%](#). Overall, our work suggests the important role that diversity of the reasoning procedure plays and points to scaling along the diversity dimension as a promising research direction.

## 1 INTRODUCTION

Large reasoning models (LRMs) have emerged as a promising paradigm for complex problem solving, where performance is achieved through generating multi-step chains of thought that decompose tasks into intermediate reasoning steps. Recent advances in scaling and reinforcement learning have further strengthened their ability to tackle challenging reasoning tasks, making them a central focus of research in artificial intelligence. Nevertheless, current approaches often impose overly rigid reasoning structures—such as strictly chain-like sequences (Wei et al., 2022; Zhang et al.) or explicitly organized trees and graphs (Yao et al., 2023; Long, 2023; Besta et al., 2024; Yao et al., 2024). While intuitive, these structures risk introducing strong priors that may not align with how human reasoning naturally unfolds, which is typically flexible, associative, and non-linear—more like a cloud of evolving perspectives than a fixed path.

To better capture this cognitive flexibility, we introduce the concept of *thinking schemata*<sup>1</sup>. Inspired by schema theory in cognitive neuroscience (Axelrod, 1973; Arbib, 1992; Fischbein & Grossman, 1997), we define thinking schemata as the distinct transitions (e.g., “alternatively”, “on the other hand”, etc.) between intermediate reasoning steps (i.e., *Reasoning Transitions*) and the variety of solution trajectories and outcomes that emerge from them (i.e., *Answer Candidates*), shown in Figure 1. Unlike previous views that regard reasoning as a linear deduction, our schemata definition centers on the dynamics of the reasoning process itself, abstracting away from rigid structures to instead measure the number of potential intermediate answers and the transitions that connect them. This abstraction encourages models to not only “think deep” but also “think diverse”, promoting richer, more human-like reasoning trajectories.

<sup>1</sup>Compared to the previous parlance of schemata (Agarwal et al., 2025; Wen et al., 2025b; Chen et al., 2025; Shen et al., 2025) which mostly focus on the data with graph structure, the thinking schemata we defined in this work consider more regarding the thinking structure.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

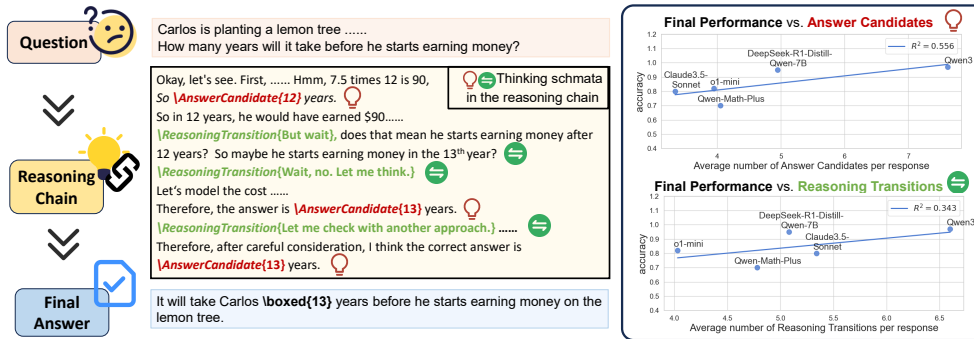


Figure 1: An illustration of the thinking schemata, including answer candidates and the reasoning transitions, that we defined and their relation with the final performance.

To concretize the above intuition, we further observe that the diversity of reasoning trajectories plays a crucial role in model performance, shown in Figure 1 and Section 2.2, and then introduce DiScO (**D**iverse **S**chemata **P**olicy **O**ptimization), a framework designed to enhance thinking schemata during both training and inference. In training, we equip the model with a basic self-awareness of its reasoning process by having it annotate points of thought transition and enumerate possible intermediate answers. We then apply reinforcement learning to encourage the generation of diverse transitions and distinct solution paths, thereby strengthening its reasoning capacity. At inference time, we further promote diversity through two complementary strategies: (i) truncating the initial 20% of the reasoning chain to allow the model to “restart” from a fresh perspective, and (ii) filtering out duplicated reasoning by detecting the first point of repetition and removing the redundant span. Together, these mechanisms push the model to “think deep and think diverse,” yielding more flexible and human-like reasoning trajectories.

We evaluate DiScO on a range of challenging mathematical reasoning benchmarks against both open-source and frontier reasoning models. Across settings, DiScO consistently outperforms comparable baselines at the 7B- and 32B-parameter scales, demonstrating substantial improvements in pass@1 accuracy. Notably, we also observe a clear relationship between the diversity of thinking schemata and model performance: models trained with DiScO generate more varied intermediate reasoning steps, which strongly aligns with gains in both accuracy and generalization. These results validate that promoting diverse thinking schemata is an effective and scalable way to strengthen reasoning ability, enabling DiScO not only to set new state-of-the-art results among open-source models but also to narrow the gap with leading proprietary systems. Furthermore, our findings suggest that scaling along the diversity dimension represents a promising and necessary direction for future research in advancing large reasoning models.

## 2 THINKING SCHEMATA

### 2.1 DEFINITIONS

*Schemata*, originally introduced in cognitive science (Casson, 1983; Rumelhart, 1984; Bartlett, 1995), refer to structured frameworks or mental templates that guide perception, understanding, and reasoning (Fischbein, 1999; Fischbein & Grossman, 1997; Cheng & Holyoak, 1985). Inspired by this theory, we define **Thinking Schemata** in the context of LRMs as the latent reasoning structure that governs how the model transitions between intermediate thoughts and arrives at potential solutions during multi-step problem solving.

Thinking schemata provide a higher-level abstraction of the reasoning dynamics within a model. Rather than focusing solely on final answers or local token-level behavior, we view thinking schemata as encompassing the broader cognitive pattern of the model’s inference process. This includes how it shifts perspectives, explores solution space, and generates candidate reasoning paths.

*Diversity of Thinking Schemata:* In this work, we particularly focus on two key properties that characterize the diversity and richness of thinking schemata:

Table 1: Diversity results of sampled reasoning chains on the AMC 2023 benchmark.

Model	Accuracy	Answer Candidates-avg	Reasoning Transition-avg
Qwen 3 <sup>2</sup>	97.5	7.62	6.60
DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI, 2025)	95.0	4.95	5.08
o1-mini (Jaech et al., 2024)	82.5	3.95	4.03
Claude3.5-Sonnet (Anthropic, 2025)	80.0	3.34	5.34
Qwen-Math-Plus <sup>2</sup>	70.0	4.05	4.78

**Reasoning Transition.** We define *Reasoning Transition* as the transition between semantically distinct reasoning steps or perspectives within a reasoning chain. These transitions reflect changes in reasoning mode—such as shifting from numerical computation to algebraic manipulation, or from a geometric intuition to a formal derivation. A model exhibiting more diverse Reasoning Transitions is considered to have greater cognitive flexibility and is better positioned to explore alternative reasoning routes.

**Answer Candidate.** The *Answer Candidate* refers to the set of distinct final or intermediate solutions the model proposes during its reasoning process. A model with diverse thinking schemata is expected to raise multiple plausible answers or sub-conclusions when reasoning, particularly in tasks with ambiguity, multi-path solutions, or when under uncertainty. The presence of multiple answer hypotheses is indicative of a broader inferential manifold being explored.

By capturing and promoting variation in both *Reasoning Transition* and *Answer Candidate*, our approach aims to elicit a richer distribution of thinking schemata, thus enhancing the reasoning capability and robustness of LRMs. An illustrative example is provided in Appendix E.2

## 2.2 DIVERSE THINKING SCHEMATA HELP LARGE REASONING MODELS

While LRMs have made substantial progress on complex tasks, we hypothesize that their performance can be further improved by increasing the diversity of their thinking schemata. Concretely, we expect that (1) a larger set of distinct answer candidates and (2) richer sequences of semantic transitions afford more opportunities for discovering correct solution trajectories, whereas (3) repetitive reasoning segments reflect a bias toward recycled templates and thus harm final performance.

**Experimental setup.** To evaluate these hypotheses, we sample multiple reasoning chains per input from a range of models (o1-mini (Jaech et al., 2024), DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI, 2025), Claude3.5-sonnet (Anthropic, 2025), Qwen-math-plus<sup>2</sup> and Qwen3<sup>2</sup>) on the AMC 2023 benchmark (Art of Problem Solving, 2023). Each sampled chain is annotated by Qwen3 using the labeling prompt described in Appendix C.2, producing per-instance metrics including the number of Answer Candidates (“answerCandidates-avg”) and the number of semantic transitions (“reasoningTransition-avg”). Descriptive statistics for these metrics are reported in Table 1.

**Results.** We quantify the relationship between model performance (measured as accuracy) and diversity metrics by fitting simple linear regression models across multiple model–dataset points. Across models, we observe a consistent pattern: higher Answer Candidates-avg and more reasoning transitions are associated with improved performance. The coefficients of determination  $R^2$  are shown in Figure 1.

**Discussion.** The empirical relationship between diversity metrics and accuracy supports the view that there is a potential that LRMs could benefit from exploring varied inferential trajectories rather than relying on a narrow set of cognitive templates. These findings highlight a gap in current approaches: although models can produce multiple reasoning chains, they rarely cultivate structured diversity in how solutions are pursued. To address this, we introduce a method that explicitly promotes diverse schemata in reasoning. By equipping models with mechanisms to recognize and regulate their own thinking schemata, and by reinforcing the generation of distinct transitions and solution paths, our approach aims to transform the observed correlations between diversity and performance into concrete improvements.

<sup>2</sup><https://bailian.console.aliyun.com/>

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

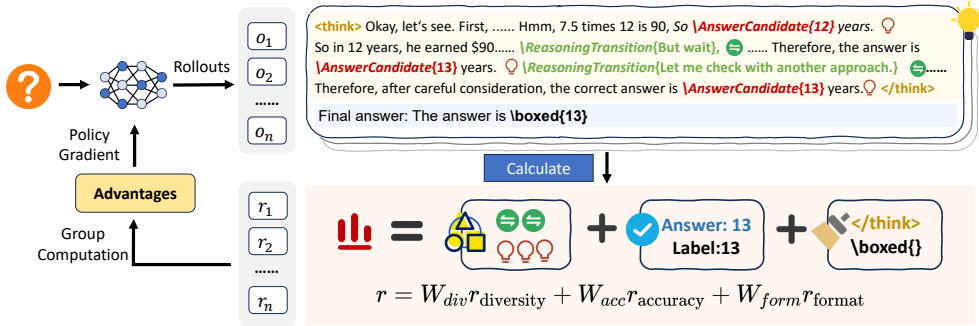


Figure 2: An illustration of the DiScO framework. During training, the model generates reasoning chains and receives rewards for the diversity of answer candidates and reasoning transitions, structured formatting, and accurate final outputs.

### 3 METHOD

We aim to enhance the reasoning ability of LRMs by explicitly modeling and encouraging diversity in their thinking schemata. Our approach consists of two key components. First, we equip the model with an annotation ability, enabling it to mark transitions in its reasoning process and identify possible intermediate answers, thereby fostering self-awareness of its own thought dynamics. Second, we introduce a reinforcement learning framework—Diverse Schemata Policy Optimization—that builds upon this annotation ability to encourage the generation of distinct reasoning transitions and diverse solution trajectories. Together, these components endow the model with the capacity to “think deep and think diverse”, leading to more flexible and robust reasoning.

#### 3.1 ANNOTATION ABILITY EQUIPMENT

We distill from Qwen-max (Qwen et al., 2024) to equip our models with the basic ability of self-awareness about the thinking schemata in its generated reasoning chain. We sampled 799 “question-reasoning-answer” triplets from the OpenR1-Math-220k (Hugging Face, 2025) dataset and have the Qwen-max to annotate the thinking schemata in the reasoning chains. The prompt we used for this task is shown in Appendix C.2 and the hyperparameters are shown in Section 4.1.

#### 3.2 DIVERSE SCHEMATA POLICY OPTIMIZATION (DISCO)

To encourage large reasoning models to “think deep and think diverse,” we adopt Group Relative Policy Optimization (GRPO) as the reinforcement learning backbone. GRPO (Shao et al., 2024; Guo et al., 2025) is a variant of PPO (Schulman et al., 2017), which avoids training an additional value function by leveraging groupwise relative comparisons of sampled outputs. Specifically, given a question  $q$ , GRPO samples a group of outputs  $\{o_1, \dots, o_G\}$  from the old policy  $\pi_{\theta_{old}}$  and optimizes the policy  $\pi_{\theta}$  with the following objective:

$$r_{i,t} = \frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t} | q, o_{i,<t})} \quad (1)$$

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{[q, \{o_i\}]} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min [r_{i,t} \hat{A}_{i,t}, \text{clip}(r_{i,t}, 1 - \epsilon, 1 + \epsilon) \hat{A}_{i,t}] - \beta \mathbb{D}_{KL}[\pi_{\theta} || \pi_{ref}] \right\} \right] \quad (2)$$

where  $\hat{A}_{i,t}$  is the advantage estimated by comparing rewards across outputs in the same group,  $\epsilon$  and  $\beta$  are hyper-parameters,  $\pi_{ref}$  is a frozen reference model.

**Reward Design.** Unlike standard GRPO, which typically relies on accuracy- and format-based rewards, we introduce a composite reward that explicitly encourages diverse schemata in reasoning:

$$R = W_{div}R_{diversity} + W_{acc}R_{accuracy} + W_{form}R_{format}, \quad (3)$$

where  $W_*$  denotes the corresponding weight assigned to each reward component.

**Diversity Reward.** To capture richness in reasoning trajectories, we measure four complementary aspects:

$$R_{\text{diversity}} = S_{\text{ansCnt}} + S_{\text{thoughtCnt}} + S_{\text{ansDiv}} + S_{\text{ansAcc}}, \quad (4)$$

$$S_{\text{ansCnt}} = \min(N_{\text{ans}}, \text{Max}_{\text{ansCnt}}) \cdot w_{\text{cnt}}, \quad (5)$$

$$S_{\text{thoughtCnt}} = \min(N_{\text{thought}}, \text{Max}_{\text{thoughtCnt}}) \cdot w_{\text{cnt}}, \quad (6)$$

$$S_{\text{ansDiv}} = \min(N_{\text{ans}}^{\text{uniq}}, \text{Max}_{\text{ansDiv}}) \cdot w_{\text{div}}, \quad (7)$$

$$S_{\text{ansAcc}} = \min(N_{\text{ans}}^{\text{true}}, \text{Max}_{\text{ansAcc}}) \cdot w_{\text{true}}. \quad (8)$$

Here  $N_{\text{ans}}$  counts answer candidates,  $N_{\text{thought}}$  counts reasoning transitions,  $N_{\text{ans}}^{\text{uniq}}$  measures unique solution candidates, and  $N_{\text{ans}}^{\text{true}}$  counts correct intermediate candidates.  $\text{Max}_*$  and  $\omega_*$  are hyperparameters, where  $\text{Max}_*$  denotes the maximum truncation limit for the corresponding  $N_*$ , and  $\omega_*$  represents its associated weight.

The above terms directly reflect the two core aspects of thinking schemata:  $N_{\text{thought}}$  captures the richness of **Reasoning Transitions**, while  $N_{\text{ans}}$  and its variants ( $N_{\text{ans}}^{\text{uniq}}$ ,  $N_{\text{ans}}^{\text{true}}$ ) operationalize the diversity and reliability of **Answer Candidates**. By rewarding both dimensions jointly, the model is encouraged to explore alternative inferential routes while maintaining plausible and accurate solution hypotheses.

**Accuracy Reward.** The model is rewarded for producing a correct final answer:

$$R_{\text{accuracy}} = \mathbb{I}(\text{ans}) = \begin{cases} 1 & \text{if the final answer is correct,} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

**Format Reward.** To ensure well-structured reasoning outputs (using the `</think>` tag to mark the end of thought and the `\boxed{\}` to represent the final answer), we provide a lightweight reward for following the annotation format:

$$R_{\text{format}}(o) = 0.5 \cdot \mathbb{I}["</think>" \in o] + 0.5 \cdot \mathbb{I}["\boxed" \in o]. \quad (10)$$

By combining these components, DiScO aligns reinforcement learning not only with accuracy but also with the promotion of diverse and structured thinking schemata, bridging the empirical link between diversity and performance observed in our analysis.

### 3.3 INFERENCE-TIME INTERVENTIONS FOR ENHANCING REASONING DIVERSITY

To encourage deeper and more diverse reasoning during inference on complex tasks, we introduce two simple yet effective strategies that dynamically optimize the input reasoning chain to the model at inference time.

**Initial Truncation:** We allow the model to “forget” the initial 20% of its prior reasoning chain by truncating the beginning of the reasoning sequence and continuing generation. This strategy simulates a cognitive reset, giving the model more test-time “thinking space” to reframe the problem and potentially explore new reasoning trajectories.

**Truncation with Repetition Elimination:** In cases where the model exceeds the max generation length and produces highly repetitive reasoning—defined as the repetition of the same 15-word phrase multiple times—we apply a filtering step that removes duplicated segments. Specifically, we identify the first point of repetition and truncate the input to remove the duplicated span. This method prevents the model from being trapped in local loops and encourages it to explore more diverse and meaningful reasoning paths.

Together, these techniques aim to expand the inference-time search space, mitigate redundancy, and promote the generation of varied thinking schemata.

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETUP

**Datasets** We use DeepSeek-R1-Distill-Qwen-7B/32B as our base models and build our implementation on the VeRL framework (Sheng et al., 2024). During the supervised fine-tuning (SFT) stage, we leveraged a subset of the OpenR1-Math-220k dataset. The first 962 samples were annotated using the Qwen-Max<sup>3</sup> and Qwen3<sup>3</sup> APIs, with the annotation prompt provided in the Appendix C.2. After filtering, we retained 840 samples for training and 45 samples for validation. For the GRPO stage, we randomly sampled 8k instances from the DeepScaler (Shi et al., 2025) dataset within an accuracy range of 20%–80%, reserving 0.5% as the validation set. We use 8 rollouts per prompt, and a max length of 32768 tokens. Other training hyperparameters can be found in Appendix B.

To evaluate the model’s reasoning capabilities, we select benchmarks of varying difficulty levels. For mathematical reasoning, we adopt GSM8K (Cobbe et al., 2021), MATH500 (Hendrycks et al., 2024), and the competition-level AIME 2024 (Art of Problem Solving, 2025a), AIME 2025 (Art of Problem Solving, 2025b), and AMC 2023 (Art of Problem Solving, 2023). For STEM reasoning, we use the challenging GPQA-Diamond (Rein et al., 2024) benchmark.

**Baselines** We leverage a diverse range of Frontier LLMs, Open-Sourced Reasoning LLMs as our baseline models. Frontier LLMs include leading proprietary models such as GPT-4o (Hurst et al., 2024), Claude 3.5-Sonnet, and the GPT-o1 series, which are recognized for their advanced reasoning capabilities. Open-Sourced Reasoning LLMs feature prominent models like DeepSeek-Coder-V2-Instruct (Zhu et al., 2024), Mathstral-7B (Mistral AI, 2024), NuminaMath-72B (Li et al., 2024), LLaMA3.1 series (Grattafiori et al., 2024), and Qwen2.5-Math-72B (Yang et al., 2024), which are widely benchmarked in mathematical tasks. The 7B and 32B parameter-scaled cohorts include base and instruction-tuned variants such as Qwen2.5-Math-7B/32B, ReasonFlux-7B/32B (Yang et al., 2025), and QwQ-32B (Qwen, 2024), which are evaluated to explore the trade-offs between model size and performance.

In the ablation study, we fine-tune DeepSeek-R1-Distill-Qwen-7B with standard SFT to obtain Qwen2.5-SFT-7B, which serves as the baseline for comparison with Qwen2.5-Anno-7B (with annotation capability). We further perform GRPO training without the diversity reward on Qwen2.5-Anno-7B and Qwen2.5-Anno-32B, obtaining Qwen2.5-Anno-GRPO-7B and Qwen2.5-Anno-GRPO-32B, which are compared with our DiScO models.

### 4.2 RESULTS

**Main results** Table 2 reports Pass@1 accuracy across multiple mathematical reasoning benchmarks. Among frontier LLMs, GPT-o1-mini and GPT-o1-preview perform strongly but show limited accuracy on competition datasets such as AIME. For open-source models, math-specialized systems like Qwen2.5-Math-72B and DeepSeek-V3 outperform earlier instruction-tuned baselines, while ReasonFlux-32B sets a strong baseline with 91.2% on MATH-500 and 56.7% on AIME 2024.

Our proposed DiScO models achieve the best results at both scales. DiScO-7B obtains 95.4% on MATH-500, 86.7% on AIME 2025, and an average of 78.1%, clearly surpassing other 7B models. DiScO-32B further improves to 93.8% on MATH-500, 86.7% on AIME 2025, 66.7% on AIME 2024, and 83.2% average, outperforming both frontier and open-source baselines.

These results demonstrate that DiScO substantially advances mathematical reasoning performance, particularly on challenging competition-level tasks. Conversely, on datasets where baselines already perform well, the gains are more modest. Taken together, this pattern suggests that promoting thinking schemata diversity is especially valuable when models face problems that require extensive exploration of the reasoning space.

**Ablation results for diverse reward design** Table 3 presents Pass@1 accuracy across six mathematical reasoning benchmarks, comparing the full DiScO model with ablated variants and the Qwen2.5-Anno-GRPO-7B baseline.

<sup>3</sup><https://bailian.console.aliyun.com/>

Table 2: Pass@1 accuracy comparison on various reasoning benchmarks.

Model	MATH-500	AIME 2024	AIME 2025	AMC 2023	GPQA-Diamond	GSM8K	Average
<b>Frontier LLMs</b>							
GPT-4o	76.6	16.7	26.7	47.5	53.6	89.5	51.7
Claude3.5-Sonnet	78.3	16.0	43.3	87.5	40.4	96.4	60.3
GPT-o1-preview	85.5	44.6	46.7	90.0	<b>73.3</b>	94.9	72.5
GPT-o1-mini	90.0	56.7	50.8	95.0	60.0	95.8	74.7
<b>Open-Sourced Reasoning LLMs</b>							
DeepSeek-Coder-V2-Instruct	75.3	13.3	-	57.5	44.3	94.9	57.1
Mathstral-7B-v0.1	57.8	0.0	0.0	37.5	9.1	77.1	30.3
NuminaMath-72B-CoT	64.0	3.3	-	70.0	35.3	91.4	52.8
LLaMA3.1-8B-Instruct	51.4	6.7	0.0	25.0	30.4	82.4	32.7
LLaMA3.1-70B-Instruct	65.4	16.7	3.3	50.0	48.0	91.7	45.9
LLaMA3.1-405B-Instruct	73.8	23.3	10.0	50.0	49.0	<b>96.8</b>	50.5
Qwen2.5-Math-72B-Instruct	85.6	30.0	26.7	70.0	42.9	95.5	58.5
DeepSeek-V3	90.2	47.7	39.2	80.0	59.1	94.2	68.4
ReasonFlux-32B	91.2	56.7	37.2	85.0	61.2	79.3	68.4
rStar-Math	89.4	50.0	-	87.5	-	95.0	80.5
<b>DiScO-32B</b>	<b>93.8</b>	<b>86.7</b>	<b>66.7</b>	<b>97.5</b>	58.1	96.4	<b>83.2</b>
<b>7B-Level Base Model</b>							
Qwen2.5-Math-7B	58.8	16.7	3.3	22.5	28.3	<b>95.0</b>	37.4
SuperCorrect-7B	70.2	26.7	13.3	37.5	20.7	84.7	42.2
DeepSeek-R1-Distill-Qwen-7B	64.0	36.7	13.3	95.0	10.6	70.2	48.3
Qwen2.5-Math-7B-Instruct	82.6	13.3	16.7	28.3	<b>62.5</b>	<b>95.0</b>	49.7
ReasonFlux-7B	88.6	36.7	36.7	80.0	35.9	83.9	60.3
<b>DiScO-7B</b>	<b>95.4</b>	<b>83.3</b>	<b>53.3</b>	<b>95.0</b>	46.0	92.5	<b>77.6</b>
<b>32B-Level Base Model</b>							
Qwen2.5-32B-Instruct	79.4	16.5	13.3	64.0	49.5	94.4	52.9
Sky-T1-32B-preview	89.5	43.3	36.7	-	56.8	94.8	64.2
ReasonFlux-32B	91.2	56.7	37.2	85.0	61.2	79.3	68.4
QwQ-32B-preview	90.6	50.0	46.7	75.0	<b>65.2</b>	91.2	69.8
DeepSeek-R1-Distill-Qwen-32B	93.4	80.0	60.0	97.5	63.6	94.0	81.4
<b>DiScO-32B</b>	<b>93.8</b>	<b>86.7</b>	<b>66.7</b>	<b>97.5</b>	58.1	<b>96.4</b>	<b>83.2</b>

Table 3: Ablation results of diverse reward designs, showing Pass@1 accuracy across benchmarks.

Model	MATH-500	AIME 2024	AIME 2025	AMC 2023	GPQA-Diamond	GSM8K	Average
DiScO-7B	<b>94.8</b>	<b>83.3</b>	<b>50.0</b>	<b>95.0</b>	43.4	92.3	<b>77.0</b>
– AnswerCandidate	93.0	76.7	<b>50.0</b>	92.5	49.5	93.3	75.8 (-1.2)
– ReasoningTransition	93.6	80.0	46.7	92.5	<b>51.0</b>	93.3	76.2 (-0.8)
Qwen2.5-Anno-GRPO-7B	92.6	56.7	<b>50.0</b>	90.0	49.5	<b>93.6</b>	72.1 (-4.9)

We observe that both dimensions of schemata diversity—*Answer Candidate* and *Reasoning Transition*—contribute substantially to performance. Removing either reward leads to consistent drops in average accuracy: –AnswerCandidate decreases from 77.0% to 75.8%, while –ReasoningTransition reduces it to 76.2%. This indicates that the two objectives capture complementary aspects of reasoning diversity: Answer Candidate rewards encourage exploration of multiple solution hypotheses, whereas Reasoning Transition rewards foster richer transitions between intermediate reasoning steps.

**Ablation results for inference** Figure 3 summarizes the average Pass@1 accuracy across six mathematical reasoning benchmarks. Detailed results can be found in Appendix D. At the 7B scale, truncation consistently improves baseline models: for example, DeepSeek-R1-Distill-Qwen-7B rises from 48.3% to 56.2% with repetition elimination. Among all 7B variants, DiScO-7B achieves the highest accuracy (77.6%) when combined with initial truncation.

At the 32B scale, truncation further amplifies performance. DeepSeek-R1-Distill-Qwen-32B improves steadily from 81.4% to 81.6%, while Qwen2.5-Anno-32B and Qwen2.5-Anno-GRPO-32B also show moderate gains. The best overall results come from DiScO-32B, which reaches 83.2% with repetition elimination, establishing the strongest performance across all models.

Overall, the results show that truncation strategies yield stable gains across scales and benchmarks, with the largest benefits observed on challenging datasets such as GPQA-Diamond and AIME 2025. These findings further confirm that our lightweight truncation methods effectively reduce

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

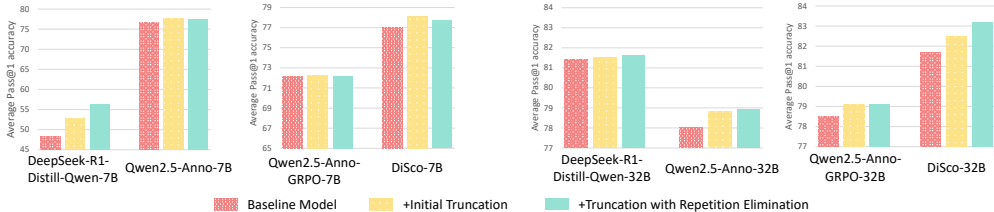


Figure 3: Ablation results for three inference strategies, where each bar denotes the average Pass@1 accuracy across six mathematical reasoning benchmarks (MATH-500, AIME 2024, AIME 2025, AMC 2023, GPQA-Diamond, and GSM8K).

redundancy in reasoning and improve overall accuracy, particularly when combined with diversity-oriented training on larger scales.

### 4.3 ANALYSIS

**Thinking Schemata Analysis** We analyze the diversity of reasoning chains generated by baseline models (Qwen2.5-Anno-GRPO-7B) and our diversity-enhanced variants (DiScO-7B) across multiple datasets (Table 4). Compared to the baseline, DiScO-7B consistently increases both the average number of reasoning transitions (#RT-avg) and answer candidates (#AC-avg), indicating richer exploration of the solution space. However, it is crucial to note that these metrics are not monotonically beneficial, but quality matters more than quantity. For instance, DeepSeek-R1-Distill-Qwen-7B generates high #RT-avg and #AC-avg values on certain datasets, such as AIME 2025, but these inflated numbers stem from repetitive content generation rather than meaningful exploration, as detailed in Appendix E.1. The true-answer ratio (TAR) which measures the proportion of correct answers among all generated candidates, serves as a key indicator of robustness, i.e. model’s ability to consistently identify correct solutions when exploring diverse reasoning paths. Notably, DiScO-7B improves TAR compared to baselines in 5 out of 6 datasets, demonstrating that greater schemata

Table 4: Diversity measurement of sampled reasoning chains. “GSM8k-100”, “GPQA-100”, and “MATH-100” denote subsets of 100 samples randomly drawn from GSM8k, GPQA-Diamond, and MATH 500, respectively. “Distill-7B” and “GRPO-7B” refer to DeepSeek-R1-Distill-Qwen-7B and DeepSeek-R1-GRPO-Qwen-7B, respectively. Here, #RT-avg = average number of Reasoning Transitions, #AC-avg = average number of Answer Candidates, RC = Repetitive Content, UA = Unique Answer, and TAR = True-Answer Ratio.

Dataset	Model	#RT-avg ↑	#AC-avg ↑	RC ↓	UA ↑	TAR ↑
AIME 2024	Distill-7B	5.73	4.07	<b>0.80</b>	2.33	0.48
	GRPO-7B	<b>12.86</b>	<b>15.21</b>	2.76	<b>7.07</b>	0.23
	DiScO-7B	12.24	8.34	1.31	4.45	<b>0.51</b>
AIME 2025	Distill-7B	<b>57.76</b>	<b>22.00</b>	2.21	<b>6.28</b>	0.08
	GRPO-7B	9.96	9.11	<b>1.36</b>	4.89	0.26
	DiScO-7B	9.76	8.62	1.52	5.00	<b>0.27</b>
AMC 2023	Distill-7B	<b>8.72</b>	7.80	1.57	2.42	<b>0.65</b>
	GRPO-7B	5.95	6.90	<b>1.45</b>	3.55	0.45
	DiScO-7B	8.64	<b>11.26</b>	1.97	<b>4.36</b>	0.50
GSM8k-100	Distill-7B	4.79	6.44	<b>1.44</b>	2.40	0.41
	GRPO-7B	6.13	8.52	1.87	<b>3.67</b>	0.45
	DiScO-7B	<b>6.32</b>	<b>10.26</b>	1.93	3.35	<b>0.50</b>
GPQA-100	Distill-7B	2.68	<b>18.21</b>	1.96	<b>4.22</b>	0.02
	GRPO-7B	10.99	6.01	<b>1.10</b>	3.82	0.06
	DiScO-7B	<b>14.59</b>	6.92	1.45	4.02	<b>0.08</b>
MATH-100	Distill-7B	<b>14.86</b>	<b>10.56</b>	1.69	2.73	0.43
	GRPO-7B	7.98	6.64	1.38	<b>2.79</b>	<b>0.47</b>
	DiScO-7B	11.64	7.05	<b>1.33</b>	2.55	0.46

diversity not only enriches reasoning dynamics but also enhances the robustness. Our analysis reveals that meaningful diversity, characterized by reduced repetitive content (RC) and increased unique answers (UA), is more important than raw counts of transitions or candidates. DiScO successfully reduces homogenization in reasoning where it maintains lower repetitive content while generating more unique intermediate perspectives. This balance between exploration and focus enables the model to avoid both the trap of repetitive reasoning loops and the inefficiency of unfocused exploration. Overall, these findings confirm that encouraging controlled variation in both transitions and candidate answers expands the model’s effective reasoning space, leading to more robust and flexible problem-solving behavior, particularly on challenging mathematical reasoning tasks.

## 5 DISCUSSION

**Diversity-based scaling** An important implication of our work is that scaling reasoning diversity itself could be a promising direction. Encouraging diverse thinking schemata through reinforcement

learning requires several times more computation at the same parameter scale, but the improvements on challenging benchmarks justify this cost. Similar to recent findings that LRMs increasingly rely on compute-intensive post-training, such as RL and reward modeling (Ji et al., 2025), we argue that diversity-based scaling is a worthwhile yet demanding direction, calling for both greater resources and more efficient optimization strategies, making it a crucial and challenging path for the next generation of reasoning models.

## 6 RELATED WORK

### 6.1 LARGE REASONING MODELS (LRMs)

Recent advances in Large Reasoning Models, such as OpenAI o1 and DeepSeek-R1, highlight the central role of reinforcement learning (RL) in enabling reasoning capabilities that pre-training alone cannot provide, with especially notable gains in mathematical reasoning (He et al., 2025) and code generation (Zhuo et al., 2024). DeepSeek-R1, for instance, demonstrates that large-scale RL with structured accuracy or test-based rewards, implemented via Group Relative Policy Optimization (GRPO), can induce sophisticated reasoning behaviors even before downstream alignment. As RL becomes a standard mechanism for improving LRMs, recent work has increasingly emphasized not only correctness but also diversity in reasoning, inspired by Yao et al. (2025). Consequently, a growing line of research investigates RL methods that explicitly encourage diverse multi-step reasoning. Recently, some work focuses on improving GRPO via more effective reward design. Zhang & Zuo (2025) introduces enhancements to GRPO for mathematical reasoning, incorporating a length-dependent accuracy reward, explicit penalties for incorrect answers, and a difficulty-aware advantage reweighting strategy, collectively improving learning efficiency. Meanwhile, Hu et al. (2020) study how to adaptively leverage a shaping reward by formulating its use as a bi-level optimization problem. Given GRPO’s widespread use and demonstrated effectiveness across LRM training pipelines, we also build on this framework in our method.

### 6.2 REASONING PATTERNS IN LARGE REASONING MODELS

Researchers have investigated reasoning patterns in LRMs and their effects on solving math problems. Minegishi et al. (2025) introduces “Topology of Reasoning,” where hidden-state clusters at each reasoning step form structures whose cycle frequency, graph diameter, and small-world characteristics correlate with model capacity and task difficulty, offering interpretable graph-theoretic insights into why reasoning-optimized LLMs perform better. Tian et al. (2025) presents a multi-round test-time thinking strategy where a model, given only its previous final answer (not its chain of thought), re-answers the same question across rounds—yielding consistent accuracy improvements without extra training. Marjanović et al. (2025) reveals that LRMs construct structured, multi-stage chains—starting with problem definition, followed by iterative “blooming” breakdowns and “reconstruction” reflections, before a final decision—showing that overly long reasoning degrades performance. An et al. (2025) argues LRMs often bloat reasoning steps—“overthinking”—and proposes dynamically pruning inefficient sub-patterns in multi-stage chains, yielding more concise, resource-efficient, and accurate outcomes. (Yang et al., 2025) enhances LRMs by guiding them through efficient multi-step reasoning with scalable, hierarchical thought templates, outperforming flat chain-of-thought methods. Wen et al. (2025a) shows that while smaller LLMs benefit from structured “thinking” patterns like decomposition, self-ask, self-debate, and self-critic, larger models perform best with simpler, unstructured monologue-style reasoning. Lee et al. (2025) proposes automatically extracting, clustering, and interpreting diverse chain-of-thought strategies from outputs to predict and steer LRMs toward more effective patterns.

Our definition of schemata differs from, and is more general than, the aforementioned thinking patterns, as it does not rely on any manual attribution or classification of reasoning chains. Instead, we focuses on the answer candidates and the transitions between reasoning trajectories, aiming to improve the performance of LRMs by encouraging the diverse generation of these thinking schemata.

### 6.3 DIVERSITY IN LARGE REASONING MODELS

To improve the training performance of LRMs, recent studies have proposed methods to explore diversity in LRMs during training. Yao et al. (2025) investigate the importance of promoting diversity during RL training and introduce Potential@k, a metric quantifying an LLM’s reasoning potential after RL training. Their work demonstrates a strong correlation between solution diversity and performance. To leverage this, the method integrates a token-level diversity objective into R1-zero training, enhancing exploration while maintaining stability. Wang et al. (2025) focus on local branching by identifying high-entropy tokens as key divergence points and optimizing RL updates around them. Token-level methods, however, operate on localized entropy signals and therefore encourage diversity only at the micro-level. In contrast, an emerging line of trajectory-level approaches seeks to promote diversity at the scale of entire reasoning paths. One prominent direction draws on GFlowNets, which sample structured objects, such as reasoning traces, with probability proportional to reward, thereby enabling diverse and globally consistent exploration. Hu et al. (2023) use Sub-trajectory Balance to amortize posterior sampling in LLMs. Yu et al. and Nair et al. (2025) formulate multi-step reasoning as flows on Directed Acyclic Graphs(DAGs) to encourage varied reasoning paths under minimal supervision; and Younsi et al. (2025) trains a PRM from MCTS data and fine-tunes a step-level GFlowNet with Subtrajectory Balance and PRM-based multiplicative rewards to generate accurate, diverse trajectories. Despite these successes, trajectory-level GFlowNet methods typically rely on structural assumptions such as a DAG-based formulation of trajectories and the intrinsic flow-matching constraints required by the framework. These inductive biases guide exploration but can constrain the reasoning space, limiting the model’s ability to capture free-form, spontaneous reasoning patterns that do not align with the assumed structure.

While token-level methods improve local branching behavior, their reliance on localized entropy signals restricts diversity to micro-level perturbations rather than capturing meaningful variations in reasoning strategy. GFlowNet-based approaches operate at the trajectory level but primarily optimize the probability of generating high-reward trajectories, achieving “reward-proportional” diversity rather than truly semantic diversity. In contrast, our method directly evaluates and shapes complete rollout trajectories, focusing on how they shift perspectives, explore alternative solution avenues, and form globally distinct reasoning paths. This trajectory-level view captures richer forms of exploratory reasoning and leads to genuine semantic diversity that better reflects the breadth of human-like problem-solving.

## 7 CONCLUSION

In this work, we introduce the notion of thinking schemata as a new perspective for understanding and improving the reasoning process of large reasoning models. We show that existing models tend to follow narrow and repetitive schemata, limiting their ability to generalize and explore alternative reasoning paths. To address this, we proposed DiScO, a reinforcement learning framework that equips models with self-awareness of their reasoning dynamics and explicitly promotes diversity in transitions and solution trajectories. Extensive experiments across multiple mathematical reasoning benchmarks demonstrated that DiScO substantially improves both accuracy and flexibility, narrowing the gap with frontier models. We believe that encouraging diverse thinking schemata offers a promising direction for developing more robust and human-like reasoning systems.

540 ETHICS STATEMENT

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we have made several efforts that are documented throughout this paper. Our experiments utilize the open-source models described in Section 4.1. The prompts used in our experiments are detailed in Appendix C. The complete code for our implementation, including inference processes, is provided in the supplementary materials. All datasets used in our experiments are described comprehensively in Section 4.1, and the supplementary code includes all data processing steps and any preprocessing applied. We encourage other researchers to consult these references for replicating our findings.

REFERENCES

- Bhavik Agarwal, Ishan Joshi, and Viktoria Rojkova. Think inside the json: Reinforcement strategy for strict llm schema adherence. arXiv preprint arXiv:2502.14905, 2025.
- Sohyun An, Ruochen Wang, Tianyi Zhou, and Cho-Jui Hsieh. Don’t think longer, think wisely: Optimizing thinking dynamics for large reasoning models. arXiv preprint arXiv:2505.21765, 2025.
- Anthropic. Introducing claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, August 2025.
- Michael A Arbib. Schema theory. The encyclopedia of artificial intelligence, 2:1427–1443, 1992.
- Art of Problem Solving. 2023 amc 10a problems and solutions. [https://wiki.artofproblemsolving.com/wiki/index.php/2023\\_AMC\\_10A](https://wiki.artofproblemsolving.com/wiki/index.php/2023_AMC_10A), 2023.
- Art of Problem Solving. Aime problems and solutions. [https://artofproblemsolving.com/wiki/index.php/AIME\\_Problems\\_and\\_Solutions](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions), 2025a.
- Art of Problem Solving. Aime problems and solutions. [https://artofproblemsolving.com/wiki/index.php/2025\\_AIME\\_I](https://artofproblemsolving.com/wiki/index.php/2025_AIME_I), 2025b.
- Robert Axelrod. Schema theory: An information processing model of perception and cognition. American political science review, 67(4):1248–1266, 1973.
- Frederic Charles Bartlett. Remembering: A study in experimental and social psychology. Cambridge university press, 1995.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI conference on artificial intelligence, volume 38, pp. 17682–17690, 2024.
- Ronald W Casson. Schemata in cognitive anthropology. Annual review of anthropology, pp. 429–462, 1983.
- Yiye Chen, Harpreet Sawhney, Nicholas Gydé, Yanan Jian, Jack Saunders, Patricio Vela, and Ben Lundell. A schema-guided reason-while-retrieve framework for reasoning on scene graphs with large-language-models (llms). arXiv preprint arXiv:2502.03450, 2025.

- 594 Patricia W Cheng and Keith J Holyoak. Pragmatic reasoning schemas. *Cognitive psychology*, 17  
595 (4):391–416, 1985.
- 596
- 597 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
598 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
599 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 600 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,  
601 2025. URL <https://arxiv.org/abs/2501.12948>.
- 602 Efraim Fischbein. Intuitions and schemata in mathematical reasoning. *Educational studies in*  
603 *mathematics*, 38(1):11–50, 1999.
- 604
- 605 Efraim Fischbein and Aline Grossman. Schemata and intuitions in combinatorial reasoning.  
606 *Educational studies in Mathematics*, 34(1):27–47, 1997.
- 607 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
608 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,  
609 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Ko-  
610 renev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava  
611 Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux,  
612 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret,  
613 Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius,  
614 Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary,  
615 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab  
616 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco  
617 Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind That-  
618 tai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Kore-  
619 vaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra,  
620 Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-  
621 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,  
622 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jong-  
623 soo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala,  
624 Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid  
625 El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren  
626 Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin,  
627 Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi,  
628 Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew  
629 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Ku-  
630 mar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoy-  
631 chev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan  
632 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,  
633 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ra-  
634 mon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Ro-  
635 hit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan  
636 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell,  
637 Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng  
638 Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer  
639 Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman,  
640 Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mi-  
641 haylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor  
642 Kerkez, Vincent Gougeon, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei  
643 Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang  
644 Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-  
645 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning  
646 Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh,  
647 Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria,  
Ahava Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein,  
Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, An-  
drew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, An-  
nie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,

- 648 Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leon-  
649 hardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu  
650 Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Mon-  
651 talvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao  
652 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia  
653 Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide  
654 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le,  
655 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily  
656 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smoth-  
657 ers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni,  
658 Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia  
659 Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan,  
660 Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harri-  
661 son Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaç,  
662 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James  
663 Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-  
664 nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang,  
665 Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Jun-  
666 jie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy  
667 Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang,  
668 Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell,  
669 Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa,  
670 Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias  
671 Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L.  
672 Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike  
673 Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari,  
674 Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan  
675 Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong,  
676 Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent,  
677 Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar,  
678 Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Ro-  
679 driguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,  
680 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin  
681 Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon,  
682 Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ra-  
683 maswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,  
684 Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal,  
685 Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satter-  
686 field, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj  
687 Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo  
688 Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook  
689 Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Ku-  
690 mar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov,  
691 Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiao-  
692 jian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,  
693 Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao,  
694 Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhao-  
695 duo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL  
696 <https://arxiv.org/abs/2407.21783>.  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000
- 694 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
695 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
696 via reinforcement learning. [arXiv preprint arXiv:2501.12948](https://arxiv.org/abs/2501.12948), 2025.  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian  
700 Yu, Zhenwen Liang, Wenxuan Wang, et al. Deepmath-103k: A large-scale, challenging, de-  
701 contaminated, and verifiable mathematical dataset for advancing reasoning. [arXiv preprint  
arXiv:2504.11456](https://arxiv.org/abs/2504.11456), 2025.

- 702 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
703 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021.  
704 URL <https://arxiv.org/abs/2103.03874>, 2, 2024.  
705
- 706 Edward J Hu, Moksh Jain, Eric Elmoznino, Younesse Kaddar, Guillaume Lajoie, Yoshua Bengio,  
707 and Nikolay Malkin. Amortizing intractable inference in large language models. [arXiv preprint  
708 arXiv:2310.04363](https://arxiv.org/abs/2310.04363), 2023.
- 709 Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and  
710 Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping. [Advances  
711 in Neural Information Processing Systems](https://arxiv.org/abs/2005.15931), 33:15931–15941, 2020.  
712
- 713 Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL [https:  
714 //github.com/huggingface/open-r1](https://github.com/huggingface/open-r1).
- 715 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
716 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. [arXiv preprint  
717 arXiv:2410.21276](https://arxiv.org/abs/2410.21276), 2024.  
718
- 719 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec  
720 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. [arXiv  
721 preprint arXiv:2412.16720](https://arxiv.org/abs/2412.16720), 2024.  
722
- 723 Jiaming Ji, Kaile Wang, Tianyi Alex Qiu, Boyuan Chen, Jiayi Zhou, Changye Li, Hantao Lou,  
724 Josef Dai, Yunhuai Liu, and Yaodong Yang. Language models resist alignment: Evidence from  
725 data compression. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher  
726 Pilehvar (eds.), [Proceedings of the 63rd Annual Meeting of the Association for Computational  
727 Linguistics \(Volume 1: Long Papers\)](https://arxiv.org/abs/2507.23411), pp. 23411–23432, Vienna, Austria, July 2025. Association  
728 for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.1141.  
729 URL <https://aclanthology.org/2025.acl-long.1141/>.
- 730 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
731 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
732 serving with pagedattention. In [Proceedings of the ACM SIGOPS 29th Symposium on Operating  
733 Systems Principles](https://arxiv.org/abs/2305.13011), 2023.
- 734 Seongyun Lee, Seungone Kim, Minju Seo, Yongrae Jo, Dongyoung Go, Hyeonbin Hwang, Jinho  
735 Park, Xiang Yue, Sean Welleck, Graham Neubig, et al. The cot encyclopedia: Analyzing, pre-  
736 dicting, and controlling how a reasoning model will think. [arXiv preprint arXiv:2505.10185](https://arxiv.org/abs/2505.10185),  
737 2025.  
738
- 739 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
740 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in  
741 ai4maths with 860k pairs of competition math problems and solutions. [Hugging Face repository](https://arxiv.org/abs/2409.13099),  
742 13(9):9, 2024.
- 743 Jieyi Long. Large language model guided tree-of-thought. [arXiv preprint arXiv:2305.08291](https://arxiv.org/abs/2305.08291), 2023.  
744
- 745 Sara Vera Marjanović, Arkil Patel, Vaibhav Adlakha, Milad Aghajohari, Parishad BehnamGhader,  
746 Mehar Bhatia, Aditi Khandelwal, Austin Kraft, Benno Krojer, Xing Han Lù, et al. Deepseek-r1  
747 thoughtology: Let’s think about llm reasoning. [arXiv preprint arXiv:2504.07128](https://arxiv.org/abs/2504.07128), 2025.  
748
- 749 Gouki Minegishi, Hiroki Furuta, Takeshi Kojima, Yusuke Iwasawa, and Yutaka Matsuo. Topology  
750 of reasoning: Understanding large reasoning models through reasoning graph properties. [arXiv  
751 preprint arXiv:2506.05744](https://arxiv.org/abs/2506.05744), 2025.
- 752 Mistral AI. MathΣtral. <https://mistral.ai/news/mathstral>, July 2024. URL [https:  
753 //mistral.ai/news/mathstral](https://mistral.ai/news/mathstral).
- 754  
755 Lakshmi Nair, Ian Trase, and Mark Kim. Flow-of-options: Diversified and improved llm reasoning  
by thinking through options. [arXiv preprint arXiv:2502.12929](https://arxiv.org/abs/2502.12929), 2025.

- 756 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
757 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
758 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin  
759 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,  
760 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,  
761 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. [arXiv  
preprint arXiv: 2412.15115](#), 2024.
- 763 Team Qwen. Qwq: Reflect deeply on the boundaries of the unknown, November 2024. URL  
764 <https://qwenlm.github.io/blog/qwq-32b-preview/>.
- 765
- 766 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-  
767 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-  
768 mark. In [First Conference on Language Modeling](#), 2024.
- 769 David E Rumelhart. Schemata and the cognitive system. 1984.
- 770
- 771 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
772 optimization algorithms. [arXiv preprint arXiv:1707.06347](#), 2017.
- 773
- 774 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,  
775 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemati-  
776 cal reasoning in open language models. [arXiv preprint arXiv:2402.03300](#), 2024.
- 777
- 778 Jiaxin Shen, Jinan Xu, Huiqi Hu, Luyi Lin, Fei Zheng, Guoyang Ma, Fandong Meng, Jie Zhou,  
779 and Wenjuan Han. A law reasoning benchmark for llm with tree-organized structures including  
factum probandum, evidence and experiences. [arXiv preprint arXiv:2503.00841](#), 2025.
- 780
- 781 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,  
782 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. [arXiv preprint  
arXiv: 2409.19256](#), 2024.
- 783
- 784 Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. Efficient reinforcement fine-  
785 tuning via adaptive curriculum learning, 2025. URL [https://arxiv.org/abs/2504.  
05520](https://arxiv.org/abs/2504.05520).
- 786
- 787 Xiaoyu Tian, Sitong Zhao, Haotian Wang, Shuaiting Chen, Yunjie Ji, Yiping Peng, Han Zhao, and  
788 Xiangang Li. Think twice: Enhancing llm reasoning by scaling multi-round test-time thinking.  
789 [arXiv preprint arXiv:2503.19855](#), 2025.
- 790
- 791 Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,  
792 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive  
793 effective reinforcement learning for llm reasoning. [arXiv preprint arXiv:2506.01939](#), 2025.
- 794
- 795 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
796 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. [Advances in  
neural information processing systems](#), 35:24824–24837, 2022.
- 797
- 798 Pengcheng Wen, Jiaming Ji, Chi-Min Chan, Juntao Dai, Donghai Hong, Yaodong Yang, Sirui Han,  
799 and Yike Guo. Thinkpatterns-21k: A systematic study on the impact of thinking patterns in llms.  
[arXiv preprint arXiv:2503.12918](#), 2025a.
- 800
- 801 Wuzhenghong Wen, Su Pan, et al. Schema-rl: A reasoning training approach for schema linking in  
802 text-to-sql task. [arXiv preprint arXiv:2506.11986](#), 2025b.
- 803
- 804 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-  
805 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2.5-math technical report: Toward mathematical  
expert model via self-improvement. [arXiv preprint arXiv:2409.12122](#), 2024.
- 806
- 807 Ling Yang, Zhaochen Yu, Bin Cui, and Mengdi Wang. Reasonflux: Hierarchical llm reasoning via  
808 scaling thought templates. [arXiv preprint arXiv:2502.06772](#), 2025.
- 809
- Jian Yao, Ran Cheng, Xingyu Wu, Jibin Wu, and Kay Chen Tan. Diversity-aware policy optimization  
for large language model reasoning. [arXiv preprint arXiv:2505.23433](#), 2025.

810 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
811 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances  
812 in neural information processing systems, 36:11809–11822, 2023.

813  
814 Yao Yao, Zuchao Li, and Hai Zhao. Got: Effective graph-of-thought reasoning in language models.  
815 In Findings of the Association for Computational Linguistics: NAACL 2024, pp. 2901–2921,  
816 2024.

817 Adam Younsi, Ahmed Attia, Abdalgader Abubaker, Mohamed El Amine Seddik, Hakim Hacid,  
818 and Salem Lahlou. Accurate and diverse llm mathematical reasoning via automated prm-guided  
819 gflownets. arXiv preprint arXiv:2504.19981, 2025.

820 Fangxu Yu, Lai Jiang, Haoqiang Kang, Shibo Hao, and Lianhui Qin. Flow of reasoning: Training  
821 llms for divergent reasoning with minimal examples. In Forty-second International Conference  
822 on Machine Learning.

823  
824 Jixiao Zhang and Chunsheng Zuo. Grpo-lead: A difficulty-aware reinforcement learning approach  
825 for concise mathematical reasoning in language models. arXiv preprint arXiv:2504.09696, 2025.

826  
827 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in  
828 large language models. In The Eleventh International Conference on Learning Representations.

829 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and  
830 Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In  
831 Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics  
832 (Volume 3: System Demonstrations), Bangkok, Thailand, 2024. Association for Computational  
833 Linguistics. URL <http://arxiv.org/abs/2403.13372>.

834 Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li,  
835 Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models  
836 in code intelligence. arXiv preprint arXiv:2406.11931, 2024.

837  
838 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam  
839 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code  
840 generation with diverse function calls and complex instructions. arXiv preprint arXiv:2406.15877,  
841 2024.

842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

## A THE USE OF LARGE LANGUAGE MODELS

We disclose that Large Language Models (LLMs) were utilized solely for our paper writing, including grammar correction and wording refinement. LLMs were not used to generate research ideas or perform analyses.

## B HYPERPARAMETERS

During the SFT stage, we adopt the LLaMA-Factory framework (Zheng et al., 2024) with a context length of 32K, learning rate of 1e-5, and train each model for 3 epochs. For reinforcement learning, we employ the VeRL framework (Sheng et al., 2024) with a 16K context length, 8 rollouts, and batch sizes of 256 for the 7B model and 512 for the 32B model. The 7B model is trained for 4 epochs (112 steps) on 8 nodes with 8×H20 GPUs, requiring approximately 55 hours, while the 32B model is trained for 4 epochs (72 steps) on 16 nodes with 8×H20 GPUs, taking about 36 hours. We use checkpoints at step 50 (7B) and step 40 (32B), where both models achieve peak accuracy. Hyperparameter details are summarized in Table 5. During inference, we set temperature to 0.0 and top.k to 1.0, applying greedy decoding with vLLM (Kwon et al., 2023). The prompts used for training and inference are demonstrated in Appendix C.3 and Appendix C.1, respectively.

Hyperparameter	Value
$\beta$	0.001
$\epsilon$	0.2
$W_{acc}$	2.0
$W_{div}$	1.0
$W_{form}$	1.0
$\omega_{wnt}$	0.1
$\omega_{div}$	0.2
$\omega_{true}$	0.3
$Max_{ansCnt}$	15
$Max_{thoughtCnt}$	20
$Max_{ansDiv}$	15
$Max_{ansAcc}$	15

Table 5: Hyperparameters used in GRPO training.

## C PROMPT DESIGN

### C.1 INFERENCE PROMPT DESIGN

The adopted prompt for inference is shown below. To ensure that the model engages in thorough reasoning, we follow the recommendation to enforce the model to initiate its response with “<think>\n” at the beginning of every output<sup>4</sup>.

#### Inference prompt design

```
Try to solve the following question step by step. If the final answer is obtained, use \\boxed{} to represent it.
### Question: {question}
<|Assistant|><think>\n
```

### C.2 LABELING PROMPT DESIGN

As mentioned in Section 2.2, we sample multiple reasoning chains from a range of models on the AMC 2023 benchmark, with each chain annotated by Qwen3. The adopted labeling prompt is shown

<sup>4</sup><https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B>

below. In the code implementation, we use the two tags “possibleAnswer” and “thoughtchange” to represent “Answer Candidate” and “Reasoning Transition” respectively.

### Labeling prompt design

The following is the reasoning chain that is used to answer a difficult math problem. Please process the reasoning chain according to the following rules:

1. Label all segments that are potentially final results in the reasoning chain with `\\possibleAnswer{}` format. DO NOT label all the possible intermediate results, ONLY label the ones that could be the final answers, no matter it’s correct or wrong. Label as many as you could.
2. An example of the `\\possibleAnswer{}` annotation: “Wait, 5 times 360 is 1800, and 1800 divided by 36. Let’s do that division:  $1800 \div 36$ . Hmm, 36 times 50 is 1800, right? Because  $36 \times 50$  is 1800. So,  $1800 \div 36 = \\possibleAnswer{50}$ . Therefore, the degrees for cherry pie would be `\\possibleAnswer{50}` degrees.”
3. Label all segments that indicate a shift in reasoning within the text reasoning chain using the `\\thoughtchange{}` format. Label as many as you could.
4. An example of the `\\thoughtchange{}` annotation: “`\\thoughtchange{Wait, maybe}` I messed up the dailyprogress.`\\thoughtchange{Wait, hold on}`. If the original totaltime is T days, then when they switch to the newequipment after 1/3 of the tunnel is done, which took T/3 days, and then the remaining 2/3 is doneat a slower daily rate”
5. DO NOT change other parts and keep them exactly the same as the the original solution.

```
### original solution:
{solution}
### Result:
```

### C.3 TRAINING PROMPT DESIGN

The adopted prompt for SFT and GRPO training is shown below. In the code implementation, we use the two tags “possibleAnswer” and “thoughtchange” to represent “Answer Candidate” and “Reasoning Transition” respectively.

### Training prompt design

Try to solve the following question step by step. Please show your reasoning chain according to the following rules:

1. First thinks about the reasoning chain in the mind and then provides the user with the answer. The reasoning chain is enclosed within `<think>` `</think>` tags, i.e., `<think>` reasoning chain here `</think>`. If the final answer is obtained, use `\\boxed{}` to represent it.
2. Label all segments that are potentially final results in the reasoning chain with `\\possibleAnswer` format. DO NOT label all the possible intermediate results, ONLY label the ones that could be the final answers, no matter it’s correct or wrong. Label as many as you could.
3. An example of the `\\possibleAnswer` annotation: “Wait, 5 times 360 is 1800, and 1800 divided by 36. Let’s do that division:  $1800 \div 36$ . Hmm, 36 times 50 is 1800, right? Because  $36 \times 50$  is 1800. So,  $1800 \div 36 = \\possibleAnswer50$ . Therefore, the degrees for cherry pie would be `\\possibleAnswer50` degrees.”
4. Label all segments that indicate a shift in reasoning within the text reasoning chain using the `\\thoughtchange` format. Label as many as you could.
5. An example of the `\\thoughtchange` annotation: “`\\thoughtchangeWait, maybe` I messed up the dailyprogress. `\\thoughtchangeWait, hold on`. If the original totaltime is T days, then when they switch to the newequipment after 1/3 of the tunnel is done, which took T/3 days, and then the remaining 2/3 is doneat a slower daily rate.”

```
### Question:{question}
<|Assistant|><think>
```

## D ABLATION RESULTS FOR INFERENCE

As mentioned in Section 4.2, Table 6 demonstrates detailed ablation results for two inference strategies. Overall, the results show that truncation strategies yield stable gains across scales and benchmarks. In particular, DiScO-7B improves from 50.0% to 53.3% on AIME 2025 with truncation, while DiScO-32B achieves the most pronounced gain on GPQA-Diamond, rising from 53.5% to 58.1% with truncation combined with repetition elimination. These findings further confirm that our lightweight methods effectively reduce redundancy in reasoning and improve overall accuracy.

Table 6: Pass@1 accuracy comparison on various mathematical reasoning benchmarks.

Model	MATH-500	AIME 2024	AIME 2025	AMC 2023	GPQA-Diamond	GSM8K	Average
Qwen-plus-latest	85.0	46.7	26.7	97.5	8.1	94.3	59.7
<b>7B-Level Base Model</b>							
DeepSeek-R1-Distill-Qwen-7B	64.0	36.7	13.3	95.0	10.6	70.2	48.3
+Initial Truncation	78.8	36.7	20.0	95.0	13.6	72.5	52.8
+Truncation with Repetition Elimination	81.8	43.3	26.7	95.0	17.2	72.9	56.2
Qwen2.5-SFT-7B	93.0	63.3	40.0	92.5	48.0	91.0	71.3
+Initial Truncation	93.0	63.3	40.0	92.5	48.0	92.9	71.6
+Truncation with Repetition Elimination	93.2	66.7	40.0	92.5	50.0	92.7	72.5
Qwen2.5-Anno-7B	93.6	<b>83.3</b>	46.7	95.0	48.0	93.5	76.7
+Initial Truncation	93.6	<b>83.3</b>	46.7	<b>97.5</b>	<b>51.0</b>	<b>93.7</b>	<b>77.6</b>
+Truncation with Repetition Elimination	93.6	<b>83.3</b>	46.7	<b>97.5</b>	49.5	<b>93.7</b>	77.4
Qwen2.5-Anno-GRPO-7B	92.6	56.7	50.0	90.0	49.5	93.6	72.1
+Initial Truncation	92.8	56.7	50.0	90.0	50.0	93.6	72.2
+Truncation with Repetition Elimination	92.8	56.7	50.0	90.0	50.0	93.6	72.2
DiScO-7B	94.8	<b>83.3</b>	50.0	95.0	43.4	92.3	76.5
+Initial Truncation	95.4	<b>83.3</b>	<b>53.3</b>	95.0	46.0	92.5	<b>77.6</b>
+Truncation with Repetition Elimination	<b>95.6</b>	<b>83.3</b>	50.0	95.0	46.0	92.7	77.1
<b>32B-Level Base Model</b>							
DeepSeek-R1-Distill-Qwen-32B	93.4	80.0	60.0	<b>97.5</b>	63.6	94.0	81.4
+Initial Truncation	93.6	80.0	60.0	<b>97.5</b>	<b>64.1</b>	94.0	81.5
+Truncation with Repetition Elimination	<b>93.8</b>	80.0	60.0	<b>97.5</b>	<b>64.1</b>	94.0	81.6
Qwen2.5-Anno-32B	92.2	73.3	63.3	95.0	49.0	95.3	78.0
+Initial Truncation	92.2	73.3	63.3	95.0	53.5	95.3	78.8
+Truncation with Repetition Elimination	92.2	73.3	63.3	95.0	54.0	95.5	78.9
Qwen2.5-Anno-GRPO-32B	92.8	76.7	46.7	<b>97.5</b>	61.1	96.0	78.5
+Initial Truncation	93.6	76.7	46.7	<b>97.5</b>	<b>64.1</b>	96.0	79.1
+Truncation with Repetition Elimination	<b>93.8</b>	76.7	46.7	<b>97.5</b>	<b>64.1</b>	96.0	79.1
DiScO-32b	91.4	<b>86.7</b>	<b>66.7</b>	<b>97.5</b>	53.5	94.5	81.7
+Initial Truncation	92.0	<b>86.7</b>	<b>66.7</b>	<b>97.5</b>	57.1	94.8	82.5
+Truncation with Repetition Elimination	<b>93.8</b>	<b>86.7</b>	<b>66.7</b>	<b>97.5</b>	58.1	<b>96.4</b>	<b>83.2</b>

## E CASE STUDY

### E.1 EXAMPLE OF REPETITIVE GENERATION

We observe that DeepSeek-R1-Distill-Qwen-7B tends to recycle its prior reasoning and generate highly repetitive content, especially on challenging problems.

#### Example: Response generated by DeepSeek-R1-Distill-Qwen-7B

Okay, so I'm trying to figure out this organic chemistry problem. The starting material is (2,2-dimethylbut-3-en-1-yl)oxybenzene, ..... So the structure is: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH<sub>2</sub>-CH<sub>2</sub>. But with a double bond between carbons 3 and 4, so it's actually: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH=CH<sub>2</sub>\n Wait, no, that's not correct. The butenyl group is but-3-en-1-yl, so the structure is: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH<sub>2</sub>-CH<sub>2</sub>. But with a double bond between carbons 3 and 4, so it's: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH=CH<sub>2</sub>.  
Wait, no, that's not correct. The butenyl group is but-3-en-1-yl, so the structure is: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH<sub>2</sub>-CH. But with a double bond between carbons 3 and 4, so it's: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH=CH<sub>2</sub>.  
Wait, no, that's not correct. The butenyl group is but-3-en-1-yl, so the structure is: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH<sub>2</sub>-CH<sub>2</sub>. But with a double bond between carbons 3 and 4, so it's: CH<sub>2</sub>-CH(CH<sub>3</sub>)-CH=CH<sub>2</sub>.  
.....

