

TAILORED PRIMITIVE INITIALIZATION IS THE SECRET KEY TO REINFORCEMENT LEARNING

Anonymous authors
 Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) has emerged as a powerful paradigm for enhancing the reasoning capabilities of large language models (LLMs). While RL has demonstrated substantial performance gains, it still faces key challenges, including low sampling efficiency and a strong dependence on model initialization: some models achieve rapid improvements with minimal RL steps, while others require significant training data to make progress. In this work, we investigate these challenges through the lens of reasoning token coverage and argue that initializing LLMs with diverse, high-quality reasoning primitives is essential for achieving stable and sample-efficient RL training. We propose Tailor, a finetuning pipeline that automatically discovers and curates novel reasoning primitives, thereby expanding the coverage of reasoning-state distributions before RL. Extensive experiments on mathematical and logical reasoning benchmarks demonstrate that Tailor generates more diverse and higher-quality warm-start data, resulting in higher downstream RL performance.

1 INTRODUCTION

Large Language Models (LLMs) have demonstrated remarkable reasoning capabilities across a broad range of application domains, including mathematical problem solving (Yu et al., 2023a; Yue et al., 2025b; Zeng et al., 2025b; Shen et al., 2025b), code generation (Xia et al., 2024; Yang et al., 2024b), and complex decision-making in agentic tasks such as API calling (Liu et al., 2024b; Prabhakar et al., 2025), autonomous driving (Li et al., 2024; Wei et al., 2024), and robotics (Yu et al., 2023b; Team et al., 2025a). Reinforcement learning (RL) has emerged as a promising paradigm for enhancing reasoning abilities by exploiting feedback from environments and leveraging verifiable reward signals (Ouyang et al., 2022; Wen et al., 2025). Such training approaches have been shown to improve model performance through the generation of long chains of thought (CoTs) (Wei et al., 2022) and test-time scaling (Yu et al., 2025c), thereby producing outputs of higher quality.

Although RL has demonstrated strong capabilities, it still faces significant challenges. First, RL suffers from low sampling efficiency (Haarnoja et al., 2018; Du et al., 2019; Shi & Chi, 2024), a limitation that is further exacerbated in the context of LLMs due to their large parameter space and the high computational cost associated with policy rollouts and updates (Sun et al., 2025; Zheng et al., 2025). Second, empirical studies (Gandhi et al., 2025) show that different LLMs respond inconsistently to RL: While some exhibit substantial performance gains, others show minimal or no improvement. This suggests that successful RL training on LLMs is highly sensitive to model initialization (Gandhi et al., 2025; Yue et al., 2025a; Cen et al., 2025; Chen et al., 2025c). To address these challenges, one line of research focuses on improving RL algorithms through dynamic sampling (Yu et al., 2025c) and data selection techniques (Zheng et al., 2025; Sun et al., 2025). Another direction adopts a data-centric perspective, emphasizing enhancements to data quality (Yu

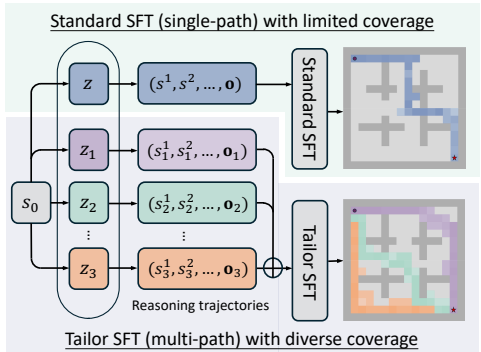


Figure 1: Coverage comparison. The goal in the maze refers to the correct answer, and the trajectories refer to the thinking tokens.

et al., 2025a). Prior data-centric efforts have identified emergent behavioral motifs, such as self-verification and reflection, that correlate with successful RL training (Gandhi et al., 2025). However, these patterns have typically been identified within narrowly scoped domains or tasks.

Our key insight is that successful RL training in LLMs requires initializing models with *reasoning primitives* that offer high thinking trajectory coverage. Specifically, reasoning primitive refers to reasoning patterns (e.g., bottom-up reasoning, top-down reasoning, etc.), which fundamentally govern the reasoning token distribution. Most existing warm-start pipelines, where supervised fine-tuning (SFT) precedes the RL stage, include only a limited set of primitive patterns, such as rule-based traces or distilled target behaviors from the SFT dataset. However, these often exhibit low coverage of the reasoning token distribution, as illustrated in Figure 1, leading to inefficient exploration and slower RL improvement. To overcome this limitation, we propose the Tailor (Task-Adaptive, Learning-Primitive-Oriented Reinforcement) finetuning pipeline, which automatically discovers diverse reasoning primitives. Our contributions are summarized as follows:

1. Analysis of reasoning primitive coverage. We investigate the role of reasoning primitive diversity in warm-start RL for LLMs, shifting the focus from hand-crafting demonstrations with specific reasoning styles.

2. The Tailor finetuning pipeline. We introduce the Tailor pipeline, which automatically curates a reasoning-diverse SFT dataset to better initialize models for subsequent RL training.

3. Comprehensive experimental evaluation. We evaluate Tailor on math and logical reasoning tasks across multiple LLM models. Our experiments and ablation studies show that Tailor improves the quality and diversity of reasoning demonstrations, leading to a stronger downstream RL.

2 RELATED WORK

LLM Reasoning. Reasoning models were first conceptualized in the OpenAI series models (Jaech et al., 2024), referring to LLMs that leverage test-time scaling by generating long chains of thought (CoTs) before final answers to improve output quality (Guo et al., 2025; Team et al., 2025b). The success of OpenAI-o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025) demonstrated that large-scale RL can incentivize reasoning capabilities in LLM training. Considerable efforts have been devoted to developing RL algorithms, such as VinePPO (Feng et al., 2023), Reinforce++ (Hu, 2025), and DAPO (Yu et al., 2025c), to advance the frontier of reasoning models. RL-based training has also been adopted in a variety of domain-specific LLMs in addition to general-purpose models (Liu et al., 2025; Shen et al., 2025a; Chu et al., 2025; Nguyen et al., 2025). Furthermore, LLM agents apply RL to enhance reasoning in textual tool use and multi-step planning (Song et al., 2025; Chen et al., 2025b; Jin et al., 2025; Qian et al., 2025), mobile and web environments (Chen et al., 2025a; Qi et al., 2024), and code generation (Wei et al., 2025; Zeng et al., 2025a; Pan et al., 2024).

Data-Centric Methods for RL Training. Data-centric approaches focus on improving the quality of training data rather than modifying the training algorithms (Zhou et al., 2023; Li et al., 2025; Guha et al., 2025; Yu et al., 2025b; Liang et al., 2025). Several works (Hong et al., 2023; Yao et al., 2024; Lee et al., 2024b) aim to shift the behavior policy distribution to facilitate more effective RL training (Yu et al., 2025a). Within the warm-start RL pipeline, where SFT precedes RL tuning, recent studies have shown that SFT induces coarse-grained changes in the LLM’s thinking pattern distribution (Fu et al., 2025), and that RL post-training tends to amplify patterns learned during SFT (Zhao et al., 2025). These findings highlight the critical role of SFT demonstrations (Yan et al., 2025; Ma et al., 2025), which serve as a “format teacher” (He et al., 2025) to guide policy rollouts and exploration during RL. Our method falls within the data-centric category, focusing specifically on curating the SFT dataset to better prepare LLMs for downstream RL training.

Reasoning Primitives. Primitives are often used to represent to capture trajectory features (Goyal et al., 2020; Peng et al., 2019). In the context of LLM reasoning, the internal thinking patterns manifested in model completions have been referred to as reasoning primitives (Li et al., 2025), behaviors (Zhao et al., 2025; Cen et al., 2025), or reasoning strategies (Qu et al., 2025). Prior works have identified specific primitives, such as reflection and backtracking, as correlates of effective test-time scaling and RL performance improvements (Yeo et al., 2025; Shen et al., 2025b; Kim et al., 2025). Additionally, by comparing models that show large versus marginal gains during RL, Gandhi et al. (2025) identified four key primitives: verification, backtracking, backward chaining, and subgoaling, which are critical to RL success. In contrast to these studies, which focus on an-

108 analyzing specific reasoning patterns, we explore how to automatically discover novel primitives and
 109 investigate how the diversity and quality of reasoning primitives affect the effectiveness of RL.
 110

111 3 PRELIMINARY

112 **Markov Decision Process (MDP).** We model the reasoning and action process of a large language
 113 model (LLM) under reinforcement learning as a Markov Decision Process (MDP) (Puterman, 2014),
 114 defined by the tuple $M = (\mathcal{S}, \mathcal{A}, P, r, s_0)$. Here, \mathcal{S} denotes the state space, where each state $s \in \mathcal{S}$
 115 encodes the current reasoning context of the LLM, including the history of reasoning tokens. The
 116 initial state $s_0 \in \mathcal{S}_0$ corresponds to a query from the query set \mathcal{S}_0 . The action space \mathcal{A} consists
 117 of all possible reasoning steps, where an action $a \in \mathcal{A}$ represents either an intermediate reasoning
 118 token or the final answer. The transition function is deterministic and defined as: $P(s_{t+1} | s_t, a_t) =$
 119 $\mathbb{I}[s_{t+1} = [s_t, a_t]]$, where each new state is formed by appending the chosen action to the current
 120 state. The reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ specifies the immediate reward received upon taking
 121 action a in state s . A policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ maps each state to a probability distribution over actions.
 122 A trajectory τ is defined as a sequence of states $\tau = (s_1, \dots, s_t, \dots, s_{T-1}, \mathbf{o})$, where s_1 through
 123 s_{T-1} represent intermediate reasoning steps (thinking tokens) and \mathbf{o} is the final answer. T denotes
 124 the number of steps.

125 **Reinforcement Learning Fine-Tuning.** The goal of Reinforcement Learning Fine-Tuning in LLMs
 126 is to optimize the expected reward over a set of queries \mathcal{S}_0 :

$$127 \max_{\theta} \mathcal{J} = \mathbb{E}_{s_0 \sim \mathcal{S}_0, s \sim \pi_\theta} \left[\sum_{t=1}^T r(s_0, \tau^{(<t)}) \right], \quad \text{where } r(s_0, \tau) = \mathbf{1}(\mathbf{o} = \mathbf{o}_{\text{gold}}), \quad (1)$$

130 where \mathbf{o}_{gold} is the ground-truth answer to the query. In this work, we primarily adopt a rule-based
 131 outcome reward that assigns a binary signal to the final answer tokens in the generated output. We
 132 study the objective function of *KL-regularized clipped policy optimization*, which incorporates regu-
 133 larization toward a reference policy, applies clipping to the policy ratio, and optimizes the following
 134 surrogate objective:

$$135 \mathcal{J}_{\text{RL}}(\mathcal{S}_0; \theta) = \mathbb{E}_{s_0 \sim \mathcal{S}_0, \tau \sim \pi_\theta(\cdot | s_0)} \left[\frac{1}{N} \sum_{i=1}^N \left[\min \left\{ \frac{\pi_\theta(\cdot | \tau^{(<t)})}{\pi_{\text{old}}(\cdot | \tau^{(<t)})} A_i, \text{clip} \left(\frac{\pi_\theta(\cdot | \tau^{(<t)})}{\pi_{\text{old}}(\cdot | \tau^{(<t)})}, 1 - \epsilon_1, 1 + \epsilon_2 \right) A_i \right\} - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \right] \right], \quad (2)$$

140 where ϵ_1, ϵ_2 are the clipping ratios, β is the KL regularization coefficient, and A_i is the advantage
 141 term determined by the specific RL algorithm. This objective is used in many RL frameworks such
 142 as GRPO (Shao et al., 2024) and DAPO (Yu et al., 2025c).

143 **Warm-Start RL.** This work focuses on the warm-start RL pipeline, which first applies Supervised
 144 Fine-Tuning (SFT) followed by RL training (Shao et al., 2024). The SFT stage enables LLMs
 145 to follow formatting instructions, become familiar with the dataset, and acquire initial reasoning
 146 capabilities within the target domain using a demonstration dataset \mathcal{D}_{SFT} . During SFT, the policy π_θ
 147 is trained by minimizing the negative log-likelihood loss:

$$148 \min_{\theta} \mathcal{L}_{\text{SFT}}(\mathcal{D}; \theta) = -\mathbb{E}_{(s_0, \tau, a) \sim \mathcal{D}_{\text{SFT}}} \left[\sum_{t=1}^T \log \pi_\theta(a_t | s_0, \tau^{(<t)}) \right] \quad (3)$$

149 For simplicity, we use the term SFT model to refer to the model after SFT.

150 4 METHOD

151 4.1 REASONING PRIMITIVE

152 In warm-start RL, LLMs initially adopt the thinking token distribution from the SFT dataset and
 153 progressively refine their reasoning patterns in RL through interaction with verifiers. Empirical
 154 studies (Gandhi et al., 2025; Cen et al., 2025) show that the distribution of thinking tokens in the
 155 SFT dataset strongly influences downstream reasoning performance: when warm-starting from two
 156 SFT datasets with different reasoning chain patterns, the resulting RL performance can diverge
 157 dramatically, even though both patterns are valid and interpretable to humans. To investigate this
 158

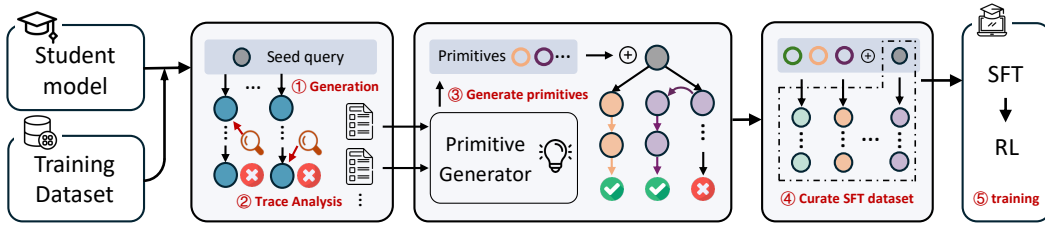


Figure 2: Overview of the Tailor finetuning pipeline.

phenomenon and identify better CoT patterns, we introduce the notion of *primitives*, which control the trajectory-wise distribution ρ of thinking tokens to model textual reasoning patterns. Formally, we augment the MDP tuple $M' = (\mathcal{S}, \mathcal{A}, P, r, s_0, z)$ with a reasoning primitive z (Qu et al., 2025):

$$\rho(\pi_\theta | z) = \mathbb{E}_{s_0 \sim \mathcal{S}_0} \prod_t \pi_\theta(a_t | s_t, z) \mathbb{I}[s_{t+1} = [s_t, a_t]], \quad (4)$$

The distribution of thinking tokens ρ is influenced by the primitives z . In the context of LLMs, primitives z can correspond to prompt instructions that are combined with the query s_0 to guide subsequent reasoning. We provide several textual examples of reasoning primitives in Figure 3. For instance, when constructing reasoning chains for a math problem, different primitives such as *Top-Down* and *Bottom-Up* reasoning, or strategies for self-verification and error recovery, can be applied. Primitives also encompass broader behaviors such as *reflection* and *backtracking*, which enable the model to detect and recover from failures during the reasoning process. This concept naturally extends beyond mathematics to other reasoning domains, such as software engineering tasks (Zhang et al., 2025), where primitives emerge from variations in agentic and prompt designs.

Primitive initialization plays a key role in warm-start RL: due to the regularization term and policy clip mechanism in the RL objective (2), LLMs are hard to automatically discover novel primitives themselves during the RL process (Zhao et al., 2025), especially when we are training specialized LLMs with limited capability and pre-training data distribution. Prior works have identified specific patterns, including *backtracking*, and *reflection* (Gandhi et al., 2025), that have a correlation with performance improvement in the subsequent RL stage. We interpret the success of specific primitives as the increase in the thinking token distribution coverage, as *reflection* and *backtracking*, which means revisit of previous context, implicitly increasing the probability to explore other thinking states.

4.2 TAILOR PIPELINE

Building on the above analysis, we argue that initializing LLMs with broad thinking token coverage leads to stronger exploration capabilities and higher RL performance. Our objective is therefore to discover a more diverse set of reasoning primitives, thereby expanding the coverage of thinking tokens and improving both the potential and sampling efficiency of the subsequent RL process. Beyond diversity, the quality of primitives is equally critical to the efficiency of RL. As noted by (Qu et al., 2025), high-quality primitives enable LLMs to achieve superior performance on specific problem sets. We assess the quality of primitives from two perspectives: (1) they should be consistent with and aligned to the target domain tasks, and (2) they should be easy for LLMs to learn, meaning that the SFT dataset distribution remains close to the pre-training distribution.

Math question: The number of each Manager's Backpack equals the difference between ... How many Manager Backpack does each Graphic Design Studio have?

Primitive 1: Top-Down	Primitive 2: Bottom-Up
1) Dependency chain discover	1) Variable Definition
2) Problem breakdown	2) Build Equations
⋮	⋮
7) Information gathering	3) Solve intermediate results
8) Final answer derivation	⋮
	7) Final answer derivation

Figure 3: Examples of reasoning primitives.

Having established the importance of diverse and high-quality reasoning primitives, we now introduce our method for initializing LLMs prior to RL. To meet the objectives of diversity and quality, we propose the **Task-AdaptIve, Learning-Primitive-Oriented Reinforcement (Tailor)** fine-tuning pipeline, illustrated in Figure 2. The key components of our approach are described as follows:

Demonstration Trace Analysis. The first step is to prompt the instructed student LLMs to generate completions on a small training data subset and label the answers using verifiers. These demonstrations reflect thinking tokens close to the student models’ pre-training distribution and exhibit a wide range of failures and errors. We then employ an LLM-based teacher model¹ to analyze these traces and propose corrections and alternative reasoning paths toward the correct answers when failures are detected. This step is designed to uncover diverse repair-oriented reasoning patterns that remain close to the student models’ pre-training distributions, making them easier to adopt and follow.

Reasoning Primitive Synthesis. Building on the summarized traces from the previous step, we synthesize reasoning primitives by prompting the teacher model to generate corresponding primitives z that guide LLMs to produce failure-recovery reasoning patterns during generation. Given that the observed failures span a wide range of types, we are able to curate a broad collection of reasoning instructions to support the subsequent SFT dataset curation.

Tailor SFT Dataset Curation. After obtaining the set of reasoning primitives \mathcal{Z} , we prompt the LLM teacher model to generate reasoning traces that explicitly follow each primitive.

After curating the Tailor SFT dataset, we fine-tune the student models on it and subsequently apply RL. A simplified overview of the Tailor process is provided in Algorithm 1, and additional details are included in the Appendix A.

Intuition. In the Tailor pipeline, student models often exhibit a variety of reasoning failures in their demonstration traces. These failures include skipping intermediate steps, making unsupported assumptions, or mishandling mathematical relationships. By analyzing these traces, the teacher model constructs a targeted set of repair primitives \mathcal{Z} that directly address the observed error patterns. This set is not only tailored to the student model’s generation behaviors and the target domain but is also enriched with diverse instructions that correspond to distinct types of reasoning errors. As a result, the dataset generated using \mathcal{Z} includes a broader range of reasoning primitives with greater coverage. This increased coverage improves alignment with the task distribution and provides a more effective foundation for subsequent reinforcement learning, leading to higher sampling efficiency and performance gains in warm-start RL.

5 EXPERIMENTS

We now evaluate the effectiveness of Tailor in improving RL performance. We begin by presenting the main results, followed by a comparison of primitive diversity, and conclude with ablation studies.

5.1 EXPERIMENT SETUP

Dataset and Benchmarks. We conduct experiments on two reasoning benchmarks: (1) iGSM (Ye et al., 2024), a grade-school math benchmark that tests mathematical and commonsense reasoning skills, containing difficulty of medium and hard tasks; and (2) KK (Knights & Naves) (Xie et al., 2024), a logical reasoning benchmark based on dynamically generated knights and knaves puzzles. We choose these two benchmarks to assess problem-solving and reasoning capabilities in LLMs while minimizing the influence of factual knowledge retrieval and reducing the risk of data contamination from the pre-training stage (Wu et al., 2025), as the synthetic nature of iGSM and KK helps

Algorithm 1 Tailor Finetuning Pipeline

Input: SFT seed dataset \mathcal{D}_s , subsets of RL dataset \mathcal{D}_{RL} , student model π_θ , teacher model π_t .
Output: Tailor RL model $\pi_{\theta'}$.

- 1: # Demonstration trace analysis.
- 2: Generate completions: $\tau \sim \pi_\theta(\cdot | s_0)$, $s_0 \sim \mathcal{D}_{RL}$
- 3: Summarize Failures $\mathcal{F} = \{f_k\}$, $f_k \leftarrow \pi_t(\tau_k)$.
- 4: # Reasoning primitive synthesis.
- 5: Generate primitives $\mathcal{Z} = \{z_m\}$: $z_m \leftarrow \pi_t(f_m)$
- 6: # Tailor SFT dataset curation.
- 7: Initialize SFT $\mathcal{D}_{SFT} \leftarrow \emptyset$
- 8: **for** $(s_0, a) \in \mathcal{D}_s$ **do**
- 9: Sample a primitive $z \sim \mathcal{Z}$;
- 10: Curate trace: $\tau \leftarrow \pi_t(\cdot | s_0, z)$;
- 11: Update $\mathcal{D}_{SFT} \leftarrow \mathcal{D}_{SFT} \cup (s_0, \tau)$
- 12: **end for**
- 13: # SFT & RL.
- 14: Perform SFT (3) on \mathcal{D}_{SFT} and RL (2) on \mathcal{D}_{RL} .
- 15: **Return:** RL policy $\pi_{\theta'}$

¹We deploy DeepSeek-V3-0324 as the teacher model in our experiments.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

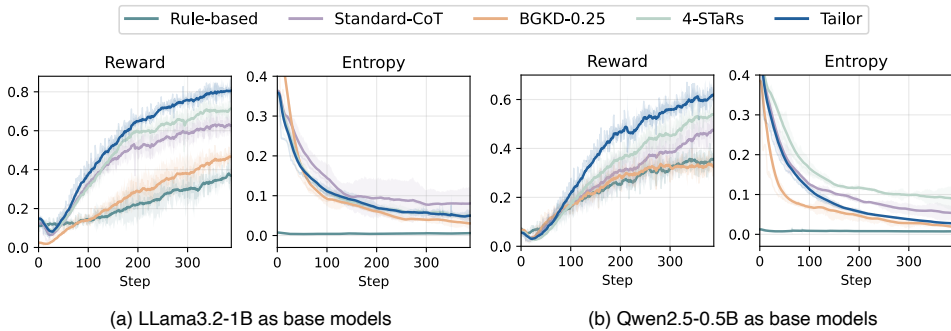


Figure 4: Training Curves of the KK tasks. We average curves with 3 random seeds.

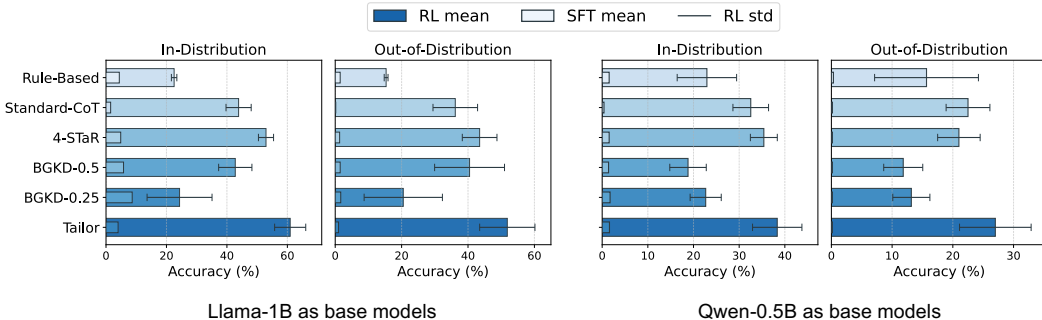


Figure 5: The evaluation results (%) on the KK reasoning tasks. We train SFT models for 4 epochs and finetune them with RL for 5 epochs. The mean value and standard deviation are calculated over 3 random seeds.

mitigate such issues (Ye et al., 2024; 2025). For both datasets, we evaluate methods in in-distribution (In-Dist) data and out-of-distribution (OOD) data. See Appendix A for more details.

Baselines. We compare Tailor with three types of baselines: (1) Rule-based Demonstration. Both the iGSM and KK benchmarks provide ground-truth functions to generate rule-based reasoning traces. We construct CoTs using these rule-based traces, perform SFT on the data, and then apply RL. We refer to this baseline as **Rule-based**. (2) Human-Crafted Primitives. Gandhi et al. (2025) identifies four critical STaR behaviors for RL and injects these reasoning patterns by prompting a teacher model with specialized instructions. We use their prompts to collect demonstrations and perform distillation, referring to this baseline as **4-STaR**. Additionally, we include a **Standard-CoT** baseline in this category, where teacher demonstrations are generated using CoT prompting (Wei et al., 2022). (3) Re-distillation. Chen et al. (2025c) propose a re-distillation strategy in which a trained model is used to regenerate the SFT dataset, producing data that better aligns with the model’s pre-training distribution. This idea is consistent with Generalized Knowledge Distillation (GKD) (Agarwal et al., 2024), which aims to reduce the distributional gap between the finetuning dataset and the student model’s pre-training distribution. We apply re-distillation to the 4-STaR SFT model to curate a new dataset, which we then combine with the original 4-STaR dataset. After SFT on this mixed data, we apply RL. We refer to this baseline as Batch-GKD (**BGKD- λ**), where $\lambda \in [0, 1]$ denotes the proportion of re-distilled data in the SFT dataset.

Other Experiment Settings. We evaluate our method using the Llama (Grattafiori et al., 2024) and Qwen (Yang et al., 2024a) model families. For Llama, we use Llama3.2-1B and Llama3.2-3B as base models; for Qwen, we use Qwen2.5-0.5B and Qwen2.5-3B. The demonstration dataset size for the SFT stage is set to 8,000, and the RL dataset contains 10,000 examples. Unless otherwise specified, we use DeepSeek-V3 (Liu et al., 2024a) as the teacher model with a decoding temperature of 0.5. The evaluation temperature for both SFT and RL models is set to 1.0. For RL, we adopt the DAPO algorithm (Yu et al., 2025c) with a rule-based outcome reward based on final answer

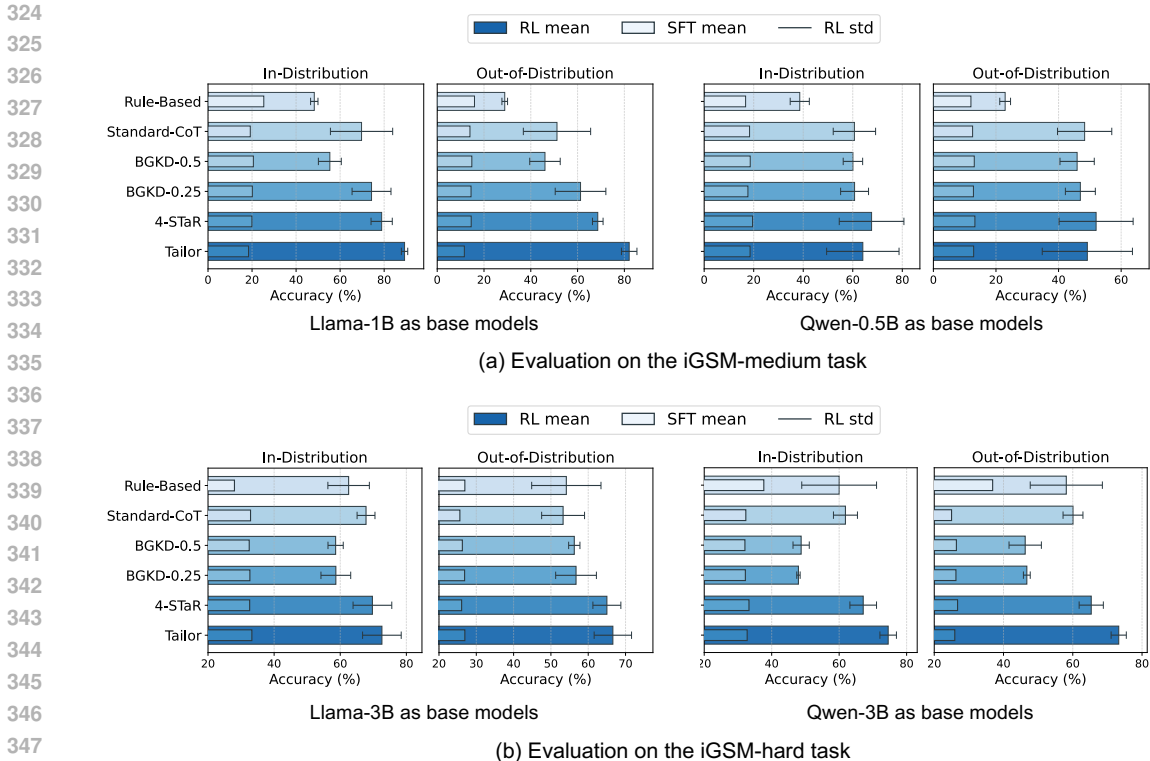


Figure 6: The evaluation results (%) on the iGSM reasoning tasks. We train SFT models for 4 epochs and finetune them with RL for 5 epochs for the iGSM-medium tasks, and finetune with RL for 100 steps for the iGSM-hard tasks. The mean value and standard deviation are calculated over 3 random seeds.

accuracy. All experiments are conducted using the Verl framework (Sheng et al., 2025). Additional details are provided in Appendix A.

5.2 MAIN RESULTS

We present the RL performance on the KK and iGSM tasks in Figure 5 and Figure 6, respectively. For the KK task, we simulate a practical setting in which the teacher model performs well, while the student model fails to achieve satisfactory performance. The teacher model achieves approximately 95% accuracy on the SFT dataset. We do not perform rejection sampling in this setting, as incorrect answers are rare and still provide useful learning signals (Xie et al., 2024). In contrast, for the iGSM tasks, we simulate a scenario where the teacher model performs poorly, achieving only around 25% accuracy on the SFT dataset. In this case, we apply rejection sampling during SFT data collection to filter out incorrect answers, while keeping the overall dataset size unchanged. The SFT dataset size is fixed at 8,000 for both KK and iGSM tasks.

KK Experiments. From Figure 5, we observe that both Llama and Qwen models struggle to achieve competitive performance after RL with rule-based CoTs. This suggests that accuracy and reasoning consistency alone are insufficient indicators of whether an SFT dataset provides good demonstrations for efficient RL. For example, rule-based annotations yield perfectly accurate and coherent reasoning traces, yet still fail to support effective RL training. These findings highlight the importance of curating SFT datasets that better prepare LLMs for RL.

Re-distillation methods (BGKD- λ) show modest performance improvements during the SFT stage. By using self-generated data, these methods reduce the distributional gap between post-training and pre-training data, resulting in more learnable patterns and improved SFT efficiency. However, BGKD still underperforms after the RL stage, as relying solely on re-distilled data diminishes the exploration benefits provided by the teacher model. In contrast, baselines such as Standard-CoT and

4-STaR achieve substantial RL performance gains by incorporating specific reasoning behaviors from teacher demonstrations. Finally, our method Tailor achieves the best overall RL performance and the largest improvements on both in-distribution and OOD evaluation sets, benefiting from the more diverse reasoning primitives it provides. As shown in Figure 4, Tailor also exhibits faster learning and higher sampling efficiency, attributed to its higher initial thinking token coverage. The entropy of the baseline methods, Standard-CoT and 4-STaRs, also remains high, but their training rewards are lower. This suggests that they may be exploring only superficial variations, such as changing individual words, rather than engaging in meaningful strategy-level exploration. As a result, the benefit for RL is limited.

iGSM Experiments. In the iGSM experiments, we observe that although rule-based methods enable good SFT performance, achieving top accuracy among some methods, they fail to provide a suitable starting point for efficient RL. This observation is consistent with the findings in the KK experiments and further reinforces that initial SFT model accuracy is not the sole factor influencing RL effectiveness; the underlying thinking token distribution also plays a crucial role. BGKD- λ yields only modest improvements in RL performance. In contrast, the 4-STaR baseline enhances RL outcomes by incorporating behaviors such as *subgoaling* and *reflection*, which have been shown to correlate with successful RL. Our method, Tailor, achieves top post-RL performance by leveraging the diverse reasoning primitives introduced during the SFT stage. In experiments using Qwen-0.5B as the base model, 4-STaR performs slightly better, likely due to the limited capacity of smaller models to learn and generalize from the diverse reasoning primitives introduced by Tailor during SFT. Nevertheless, Tailor consistently demonstrates substantial RL performance gains over all other baselines across other iGSM testing datasets.

Diversity Analysis. To evaluate the diversity of reasoning primitives in the curated SFT datasets, we use NV-Embed-v2 (Lee et al., 2024a) to compute embeddings of the thinking tokens and calculate the average pairwise cosine similarity for responses to the same seed query. Results on the KK dataset are shown in Figure 7. We observe that 4-STaR improves diversity over Standard-CoT, as reflected by a lower median similarity and a longer tail extending toward lower values. This improvement can be attributed to the inclusion of exploration behaviors such as *reflection* and *backtracking*, which increase variation in the generated reasoning traces. Moreover, our proposed method Tailor further reduces the similarity scores compared to 4-STaR, with an even lower median and a heavier tail in the low-similarity region. This indicates that Tailor significantly enhances high-level reasoning diversity in the SFT dataset, better preparing LLMs for effective exploration during RL.

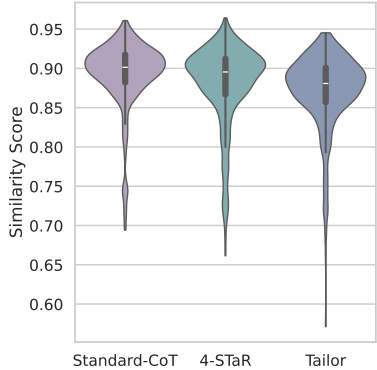


Figure 7: Primitive similarity: lower score means higher diversity.

Takeaways

- (1) Reasoning correctness in \mathcal{D}_{SFT} alone does not guarantee RL success: although rule-based CoT achieves the highest accuracy in \mathcal{D}_{SFT} , it does not result in an efficient RL process.
- (2) Thinking-token coverage is crucial for RL initialization: Tailor increases reasoning-trajectory diversity in \mathcal{D}_{SFT} , enhancing exploration and improving RL efficiency.

5.3 ABLATION STUDY

Ablation on Reasoning Primitives. In our main experiments, we use 25 reasoning primitives for SFT data collection. To investigate the effect of diversity and coverage in reasoning primitives, we vary the size of the primitive set from 4 to 25 while keeping the overall SFT dataset size fixed. We perform ablations on the iGSM-medium task using Llama-1B as the base model. The results are shown in Figure 8 (a). We observe that when the primitive set is small and less diverse, the resulting RL performance is lower. This is likely because the limited primitive set covers only a small subset of effective strategies, making it difficult to generalize across a wide range of questions and hard for exploration. As the number of reasoning primitives increases to 25, post-RL performance improves, highlighting the importance of primitive coverage when preparing LLMs for RL.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

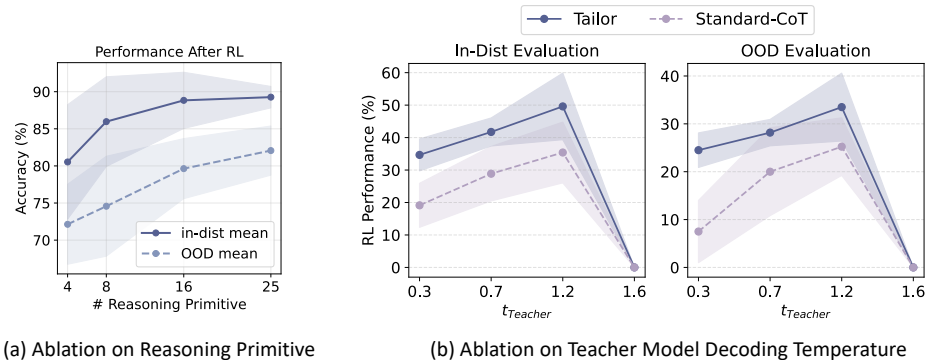


Figure 8: Ablation Study. (a) Effect of the number of reasoning primitives in the SFT dataset. Experiments are conducted on the iGSM-medium dataset using Llama3.2-1B as the base model. (b) Effect of the teacher model’s decoding temperature. Experiments are conducted on the KK dataset using Qwen2.5-0.5B as the base model.

Ablation on Teacher Model Temperature.

Adjusting the teacher model’s sampling temperature $t_{Teacher}$ is a common technique for controlling the diversity and coverage of demonstration datasets (Mukherjee et al., 2023; Hong et al., 2024). We vary the decoding temperature of the teacher model and conduct ablations on the KK task using Llama-1B as the base model. The results are shown in Figure 8 (b). We observe that as the decoding temperature increases, RL performance generally improves, suggesting that higher diversity in teacher outputs benefits downstream training. However, our method Tailor consistently outperforms the baselines across all temperature settings, demonstrating a better balance between quality and diversity. Notably, when the temperature is set too high (e.g., $t = 1.6$), the quality of demonstrations deteriorates significantly, and the student model fails to learn reasonable traces in SFT, resulting in near-zero accuracy after RL.

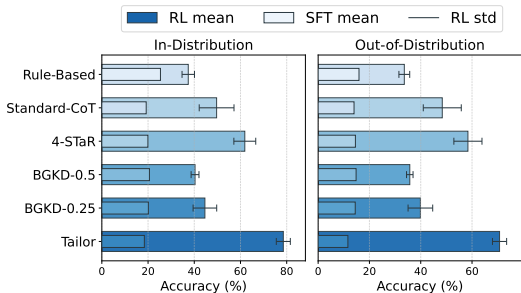


Figure 9: GRPO experiments on the iGSM-medium task with LLama-1B as base models.

Ablation on Alternative RL Algorithms. In addition to DAPO, we also combine our Tailor pipeline with GRPO (Shao et al., 2024) on the iGSM-medium task using Llama-1B as the base model. The results are shown in Figure 9, where we observe that the conclusions remain consistent with those from the DAPO: Tailor achieves the highest RL improvement and final performance. These results demonstrate the compatibility of our approach with a broader range of RL algorithms.

6 CONCLUSION

In this paper, we interpret the variation in RL performance across different initial models and the issue of low sampling efficiency through the lens of reasoning primitive quality and coverage. To address these challenges, we propose Tailor, a pipeline that discovers diverse and high-quality reasoning primitives to construct demonstration data for warming up LLMs in RL. By increasing the coverage of the thinking token distribution, Tailor facilitates faster exploration and unlocks greater performance potential during RL tuning. Extensive experiments across multiple benchmarks demonstrate that our method effectively curates a more diverse training corpus and significantly improves downstream RL performance. One limitation of Tailor is current evaluation focuses on logical and mathematical reasoning tasks. As future work, we plan to extend the pipeline to other domains such as code generation and agent-based decision-making. A potential negative impact is that misuse of our method could lead to the generation of unsafe or toxic primitives. Nevertheless, we believe Tailor sheds light on data-centric RL, emphasizing that a successful RL process depends on well-initialized models with diverse and high-quality thinking trajectory distribution.

486 REPRODUCIBILITY STATEMENT
487

488 Descriptions of our method are provided in Section 4, with implementation details and prompts
489 included in Appendices A and B. Experiments shown in Section 5 were run with multiple random
490 seeds, and we report mean values and standard deviations across runs. We use publicly available
491 datasets and benchmarks, with descriptions also provided in the appendix.
492

493 REFERENCES
494

- 495 Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu
496 Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-
497 generated mistakes. In *The twelfth international conference on learning representations*, 2024.
498 6
- 499 Zhepeng Cen, Yihang Yao, William Han, Zuxin Liu, and Ding Zhao. Behavior injection: Preparing
500 language models for reinforcement learning. *arXiv preprint arXiv:2505.18917*, 2025. 1, 2, 3, 16
501
- 502 Kevin Chen, Marco Cusumano-Towner, Brody Huval, Aleksei Petrenko, Jackson Hamburger,
503 Vladlen Koltun, and Philipp Krähenbühl. Reinforcement learning for long-horizon interactive
504 llm agents. *arXiv preprint arXiv:2502.01600*, 2025a. 2
- 505 Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Fan Yang, Zenan Zhou,
506 Weipeng Chen, Haofen Wang, Jeff Z Pan, et al. Learning to reason with search for llms via
507 reinforcement learning. *arXiv preprint arXiv:2503.19470*, 2025b. 2
508
- 509 Yutong Chen, Jiandong Gao, and Ji Wu. Towards revealing the effectiveness of small-scale fine-
510 tuning in rl-style reinforcement learning. *arXiv preprint arXiv:2505.17988*, 2025c. 1, 6
- 511 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V
512 Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation
513 model post-training. *arXiv preprint arXiv:2501.17161*, 2025. 2
514
- 515 Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient
516 for sample efficient reinforcement learning? *arXiv preprint arXiv:1910.03016*, 2019. 1
- 517 Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and
518 Jun Wang. Alphazero-like tree-search can guide large language model decoding and training.
519 *arXiv preprint arXiv:2309.17179*, 2023. 2
520
- 521 Yuqian Fu, Tinghong Chen, Jiajun Chai, Xihuai Wang, Songjun Tu, Guojun Yin, Wei Lin, Qichao
522 Zhang, Yuanheng Zhu, and Dongbin Zhao. Srft: A single-stage method with supervised and
523 reinforcement fine-tuning for reasoning. *arXiv preprint arXiv:2506.19767*, 2025. 2
- 524 Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D Goodman. Cogni-
525 tive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv*
526 *preprint arXiv:2503.01307*, 2025. 1, 2, 3, 4, 6
- 527 Anirudh Goyal, Shagun Sodhani, Jonathan Binas, Xue Bin Peng, Sergey Levine, and Yoshua Ben-
528 gio. Reinforcement learning with competitive ensembles of information-constrained primitives.
529 In *International Conference on Learning Representations*, 2020. 2
530
- 531 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhishek Pandey, Ankur Kadian, Ahmad
532 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
533 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Ko-
534 renev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava
535 Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, ..., et al. The llama 3 herd of
536 models. *arXiv preprint arXiv:2407.21783*, 2024. URL [https://arxiv.org/abs/2407.](https://arxiv.org/abs/2407.21783)
537 [21783](https://arxiv.org/abs/2407.21783). “The Llama 3 Herd of Models”. 6
- 538 Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna
539 Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, et al. Openthoughts: Data recipes for reason-
ing models. *arXiv preprint arXiv:2506.04178*, 2025. 2

- 540 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
541 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
542 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 2
- 543
- 544 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
545 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-
546 cations. *arXiv preprint arXiv:1812.05905*, 2018. 1
- 547
- 548 Lixuan He, Jie Feng, and Yong Li. Amft: Aligning llm reasoners by meta-learning the optimal
549 imitation-exploration balance. *arXiv preprint arXiv:2508.06944*, 2025. 2
- 550
- 551 Zhang-Wei Hong, Pulkit Agrawal, Rémi Tachet des Combes, and Romain Laroche. Harness-
552 ing mixed offline reinforcement learning datasets via trajectory weighting. *arXiv preprint*
553 *arXiv:2306.13085*, 2023. 2
- 554
- 555 Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Wang, Yung-Sung Chuang, Aldo Pareja, James Glass,
556 Akash Srivastava, and Pulkit Agrawal. Curiosity-driven red-teaming for large language models.
557 *arXiv preprint arXiv:2402.19464*, 2024. 9
- 558
- 559 Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv*
560 *preprint arXiv:2501.03262*, 2025. 2
- 561
- 562 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
563 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv*
564 *preprint arXiv:2412.16720*, 2024. 2
- 565
- 566 Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and
567 Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement
568 learning. *arXiv preprint arXiv:2503.09516*, 2025. 2
- 569
- 570 Joongwon Kim, Anirudh Goyal, Liang Tan, Hannaneh Hajishirzi, Srinivasan Iyer, and Tianlu Wang.
571 Astro: Teaching language models to reason by reflecting and backtracking in-context. *arXiv*
572 *preprint arXiv:2507.00417*, 2025. 2
- 573
- 574 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
575 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
576 serving with pagedattention. In *Proceedings of the 29th symposium on operating systems princi-*
577 *ples*, pp. 611–626, 2023. 18
- 578
- 579 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catan-
580 zaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embed-
581 ding models. *arXiv preprint arXiv:2405.17428*, 2024a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2405.17428)
582 [2405.17428](https://arxiv.org/abs/2405.17428). Includes NV-Embed-v2; ranks No. 1 on the Massive Text Embedding Bench-
583 mark (MTEB) as of August 30, 2024. 8
- 584
- 585 Jaewoo Lee, Sujin Yun, Taeyoung Yun, and Jinkyoo Park. Gta: Generative trajectory augmentation
586 with guidance for offline reinforcement learning. *Advances in Neural Information Processing*
587 *Systems*, 37:56766–56801, 2024b. 2
- 588
- 589 Boyi Li, Yue Wang, Jiageng Mao, Boris Ivanovic, Sushant Veer, Karen Leung, and Marco
590 Pavone. Driving everywhere with large language model policy adaptation. In *Proceedings of the*
591 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14948–14957, 2024. 1
- 592
- 593 Yang Li, Youssef Emad, Karthik Padthe, Jack Lanchantin, Weizhe Yuan, Thao Nguyen, Jason We-
ston, Shang-Wen Li, Dong Wang, Iliia Kulikov, et al. Naturalthoughts: Selecting and distilling
reasoning traces for general reasoning tasks. *arXiv preprint arXiv:2507.01921*, 2025. 2
- Yiqing Liang, Jieli Qiu, Wenhao Ding, Zuxin Liu, James Tompkin, Mengdi Xu, Mengzhou Xia,
Zhengzhong Tu, Laixi Shi, and Jiacheng Zhu. Modomodo: Multi-domain data mixtures for
multimodal llm reinforcement learning. *arXiv preprint arXiv:2505.24871*, 2025. 2

- 594 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
595 Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen,
596 Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen,
597 Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding,
598 Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Ji-
599 ash Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao
600 Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong
601 Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang,
602 Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang,
603 Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen,
604 R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi
605 Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye,
606 Shirong Ma Ye, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang,
607 Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu,
608 Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang,
609 Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha
610 Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xingchao Xie, Xingkai
611 Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng
612 Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanping Huang, Yao
613 Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yuchen Yu, Yi Zheng, Yichao Zhang, Yi-
614 fan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma,
615 Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng
616 Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxuan You,
617 Yuyang Liu, Zhongyu Zhou, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie,
618 Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report. December 2024a. URL
<https://arxiv.org/abs/2412.19437>. 6
- 619 Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong.
620 Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models.
621 *arXiv preprint arXiv:2505.24864*, 2025. 2
- 622 Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Juntao Tan, Weiran Yao, Zhiwei Liu,
623 Yihao Feng, Rithesh RN, et al. Apigen: Automated pipeline for generating verifiable and diverse
624 function-calling datasets. *Advances in Neural Information Processing Systems*, 37:54463–54482,
625 2024b. 1
- 627 Lu Ma, Hao Liang, Meiyi Qiang, Lexiang Tang, Xiaochen Ma, Zhen Hao Wong, Junbo Niu,
628 Chengyu Shen, Runming He, Bin Cui, et al. Learning what reinforcement learning can't: In-
629 terleaved online fine-tuning for hardest questions. *arXiv preprint arXiv:2506.07527*, 2025. 2
- 630 Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and
631 Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv*
632 *preprint arXiv:2306.02707*, 2023. 9
- 634 Xuan-Phi Nguyen, Shrey Pandit, Revanth Gangi Reddy, Austin Xu, Silvio Savarese, Caiming Xiong,
635 and Shafiq Joty. Sfr-deepresearch: Towards effective reinforcement learning for autonomously
636 reasoning single agents. *arXiv preprint arXiv:2509.06283*, 2025. 2
- 637
638 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
639 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
640 low instructions with human feedback. *Advances in neural information processing systems*, 35:
641 27730–27744, 2022. 1
- 642 Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe
643 Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv preprint*
644 *arXiv:2412.21139*, 2024. 2
- 645
646 Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mpc: Learning
647 composable hierarchical control with multiplicative compositional policies. *Advances in neural*
information processing systems, 32, 2019. 2

- 648 Akshara Prabhakar, Zuxin Liu, Ming Zhu, Jianguo Zhang, Tulika Awalgaonkar, Shiyu Wang, Zhiwei
649 Liu, Haolin Chen, Thai Hoang, Juan Carlos Niebles, et al. Apigen-mt: Agentic pipeline for multi-
650 turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601*,
651 2025. 1
- 652 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
653 Wiley & Sons, 2nd edition, 2014. ISBN 9781118625873. 3
- 654 Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Wenyi Zhao, Yu Yang, Xinyue Yang,
655 Jiadai Sun, Shuntian Yao, et al. Webrl: Training llm web agents via self-evolving online curricu-
656 lum reinforcement learning. *arXiv preprint arXiv:2411.02337*, 2024. 2
- 657 Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan
658 Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*,
659 2025. 2
- 660 Yuxiao Qu, Anikait Singh, Yoonho Lee, Amrith Setlur, Ruslan Salakhutdinov, Chelsea Finn, and
661 Aviral Kumar. Learning to discover abstractions for llm reasoning. In *Proceedings of the
662 AI4Math Workshop, ICML 2025*, July 2025. URL [https://openreview.net/forum?
663 id=rPzbDxcErp](https://openreview.net/forum?id=rPzbDxcErp). Poster. 2, 4
- 664 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
665 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical
666 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 3, 9
- 667 Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun
668 Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large
669 vision-language model. *arXiv preprint arXiv:2504.07615*, 2025a. 2
- 670 Maohao Shen, Guangtao Zeng, Zhenting Qi, Zhang-Wei Hong, Zhenfang Chen, Wei Lu, Gre-
671 gory Wornell, Subhro Das, David Cox, and Chuang Gan. Satori: Reinforcement learning
672 with chain-of-action-thought enhances llm reasoning via autoregressive search. *arXiv preprint
673 arXiv:2502.02508*, 2025b. 1, 2
- 674 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
675 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings
676 of the Twentieth European Conference on Computer Systems*, pp. 1279–1297, 2025. 7, 18
- 677 Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with
678 near-optimal sample complexity. *Journal of Machine Learning Research*, 25(200):1–91, 2024. 1
- 679 Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang,
680 and Ji-Rong Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement
681 learning. *arXiv preprint arXiv:2503.05592*, 2025. 2
- 682 Yifan Sun, Jingyan Shen, Yibin Wang, Tianyu Chen, Zhendong Wang, Mingyuan Zhou, and Huan
683 Zhang. Improving data efficiency for llm reinforcement fine-tuning through difficulty-targeted
684 online data selection and rollout replay. *arXiv preprint arXiv:2506.05316*, 2025. 1
- 685 Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montser-
686 rat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza,
687 Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint
688 arXiv:2503.20020*, 2025a. 1
- 689 Kimi Team, Angang Du, Bofei Gao, Bofei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun
690 Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with
691 llms. *arXiv preprint arXiv:2501.12599*, 2025b. 2
- 692 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
693 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
694 neural information processing systems*, 35:24824–24837, 2022. 1, 6

- 702 Yuxi Wei, Zi Wang, Yifan Lu, Chenxin Xu, Changxing Liu, Hao Zhao, Siheng Chen, and Yanfeng
703 Wang. Editable scene simulation for autonomous driving via collaborative llm-agents. In *Pro-*
704 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15077–
705 15087, 2024. 1
- 706 Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried,
707 Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. Swe-rl: Advancing llm reasoning via rein-
708 forcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*, 2025. 2
- 709 Xumeng Wen, Zihan Liu, Shun Zheng, Zhijian Xu, Shengyu Ye, Zhirong Wu, Xiao Liang, Yang
710 Wang, Junjie Li, Ziming Miao, et al. Reinforcement learning with verifiable rewards implicitly
711 incentivizes correct reasoning in base llms. *arXiv preprint arXiv:2506.14245*, 2025. 1
- 712 Mingqi Wu, Zhihao Zhang, Qiaole Dong, Zhiheng Xi, Jun Zhao, Senjie Jin, Xiaoran Fan, Yuhao
713 Zhou, Huijie Lv, Ming Zhang, et al. Reasoning or memorization? unreliable results of reinforc-
714 e-ment learning due to data contamination. *arXiv preprint arXiv:2507.10532*, 2025. 5
- 715 Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying
716 llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024. 1
- 717 Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih
718 Ghazi, and Ravi Kumar. On memorization of large language models in logical reasoning. *arXiv*
719 *preprint arXiv:2410.23123*, 2024. 5, 7, 16
- 720 Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang.
721 Learning to reason under off-policy guidance. *arXiv preprint arXiv:2504.14945*, 2025. 2
- 722 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
723 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,
724 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin
725 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,
726 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,
727 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report.
728 arXiv preprint arXiv:2412.15115, 2024a. URL <https://arxiv.org/abs/2412.15115>.
729 “Qwen2.5 Technical Report”. 6
- 730 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,
731 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering.
732 *Advances in Neural Information Processing Systems*, 37:50528–50652, 2024b. 1
- 733 Yihang Yao, Zhepeng Cen, Wenhao Ding, Haohong Lin, Shiqi Liu, Tingnan Zhang, Wenhao Yu,
734 and Ding Zhao. Oasis: Conditional distribution shaping for offline safe reinforcement learning.
735 *Advances in Neural Information Processing Systems*, 37:78451–78478, 2024. 2
- 736 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1,
737 grade-school math and the hidden reasoning process. In *The Thirteenth International Conference*
738 *on Learning Representations*, 2024. 5, 6, 16
- 739 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of Language Models: Part 2.2,
740 How to Learn From Mistakes on Grade-School Math Problems. In *Proceedings of the 13th Inter-*
741 *national Conference on Learning Representations*, ICLR ’25, April 2025. Full version available
742 at <https://ssrn.com/abstract=5250631>. 6
- 743 Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-
744 of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025. 2
- 745 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhen-
746 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions
747 for large language models. *arXiv preprint arXiv:2309.12284*, 2023a. 1
- 748 Ping Yu, Jack Lanchantin, Tianlu Wang, Weizhe Yuan, Olga Golovneva, Ilya Kulikov, Sainbar
749 Sukhbaatar, Jason Weston, and Jing Xu. Cot-self-instruct: Building high-quality synthetic
750 prompts for reasoning and non-reasoning tasks. *arXiv preprint arXiv:2507.23751*, 2025a. 1, 2

- 756 Ping Yu, Weizhe Yuan, Olga Golovneva, Tianhao Wu, Sainbayar Sukhbaatar, Jason Weston, and
757 Jing Xu. Rip: Better models by survival of the fittest prompts. *arXiv preprint arXiv:2501.18578*,
758 2025b. [2](#)
- 759
- 760 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian
761 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system
762 at scale. *arXiv preprint arXiv:2503.14476*, 2025c. [1](#), [2](#), [3](#), [6](#), [18](#)
- 763
- 764 Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Are-
765 nas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language to
766 rewards for robotic skill synthesis. *arXiv preprint arXiv:2306.08647*, 2023b. [1](#)
- 767
- 768 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does re-
769 inforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv
preprint arXiv:2504.13837*, 2025a. [1](#)
- 770
- 771 Yu Yue, Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiase Chen, Chengyi
772 Wang, TianTian Fan, Zhengyin Du, et al. Vapo: Efficient and reliable reinforcement learning for
773 advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025b. [1](#)
- 774
- 775 Guangtao Zeng, Maohao Shen, Delin Chen, Zhenting Qi, Subhro Das, Dan Gutfreund, David Cox,
776 Gregory Wornell, Wei Lu, Zhang-Wei Hong, et al. Satori-swe: Evolutionary test-time scaling for
777 sample-efficient software engineering. *arXiv preprint arXiv:2505.23604*, 2025a. [2](#)
- 778
- 779 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplert-
zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv
preprint arXiv:2503.18892*, 2025b. [1](#)
- 780
- 781 Kexun Zhang, Weiran Yao, Zuxin Liu, Yihao Feng, Zhiwei Liu, Rithesh RN, Tian Lan, Lei Li,
782 Renze Lou, Jiacheng Xu, et al. Diversity empowers intelligence: Integrating expertise of software
783 engineering agents. In *The Thirteenth International Conference on Learning Representations*,
784 2025. [4](#)
- 785
- 786 Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach.
787 Echo chamber: RL post-training amplifies behaviors learned in pretraining. *arXiv preprint
arXiv:2504.07912*, 2025. [2](#), [4](#)
- 788
- 789 Haizhong Zheng, Yang Zhou, Brian R Bartoldson, Bhavya Kailkhura, Fan Lai, Jiawei Zhao, and
790 Beidi Chen. Act only when it pays: Efficient reinforcement learning for llm reasoning via selective
rollouts. *arXiv preprint arXiv:2506.02177*, 2025. [1](#)
- 791
- 792 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
793 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information
794 Processing Systems*, 36:55006–55021, 2023. [2](#)
- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809

810 A MORE DETAILS OF EXPERIMENTS

811 A.1 EVALUATION BENCHMARK

812 In this section, we introduce the KK and iGSM benchmarks used for evaluation.

813 **KK** (Xie et al., 2024) is a logical reasoning benchmark composed of dynamically generated Knights
814 and Knaves (KK) puzzles. In each puzzle, every character is either a Knight (who always tells the
815 truth) or a Knave (who always lies). The puzzle presents a set of natural language statements made
816 by the characters, referring to themselves or others. The task is to infer the identity (Knight or
817 Knave) of each character based on these statements. Solving the puzzle requires ensuring global
818 logical consistency: the assigned identities must be compatible with the logical implications of all
819 statements. The input is the full puzzle text (i.e., all character statements), and the expected output
820 is a complete assignment of identities for all characters. A prediction is considered correct only if
821 the identities of all characters are correctly labeled. We provide an example of a query below.

822 **Query:**

823 A very special island is inhabited only by knights and knaves. Knights always tell the truth, and
824 knaves always lie. You meet 7 inhabitants: Abigail, Aria, Jacob, Liam, James, Amelia, and David.
825 In Abigail’s words: “David is a knight if and only if Liam is a knight”. Aria noted, “If Abigail
826 is a knight then Liam is a knight”. “Abigail is a knave” - Jacob. As Liam put it, “If Amelia is a
827 knight then Abigail is a knight”. “Jacob is a knight or Amelia is a knave,” James claimed. “David
828 is not a knight” - Amelia. David said that Jacob is a knave and David is a knight. So who is a
829 knight and who is a knave?
830

831 **Rule-based CoT:**

832 <think>

833 Assume Abigail is a knight. No contradiction is found in their claim that David is a knight if and
834 only if Liam is a knight. Assume David is a knight. No contradiction is found in their claim that
835 Jacob is a knave and David is a knight. Jacob cannot be a knight, because this would contradict
836 the claim of their own that Abigail is a knave. Assume Jacob is a knave. No contradiction is
837 found in their false claim that Abigail is a knave. Assume Liam is a knight. No contradiction is
838 found in their claim that If Amelia is a knight then Abigail is a knight. Amelia cannot be a knight,
839 because this would contradict the claim of their own that David is not a knight. Assume Amelia is
840 a knave. No contradiction is found in their false claim that David is not a knight. Assume Aria
841 is a knight. No contradiction is found in their claim that If Abigail is a knight then Liam is a
842 knight. Assume James is a knight. No contradiction is found in their claim that Jacob is a knight
843 or Amelia is a knave. < /think>

844 <answer>

845 Thus, the final answer is boxed{Abigail is a knight. David is a knight. Jacob is a knave. Liam is
846 a knight. Amelia is a knave. Aria is a knight. James is a knight.}
847 < /answer>

848 **iGSM** (Ye et al., 2024) is an infinite grade-school math problem benchmark designed to evaluate
849 LLMs’ mathematical and commonsense reasoning abilities. We apply slight modifications to iGSM
850 to improve fine-tuning stability, including the removal of modulo operations (Cen et al., 2025). This
851 is motivated by prior findings that base models struggle to compute modulo with high accuracy,
852 and such capability cannot be substantially improved through training on a small SFT dataset (Ye
853 et al., 2024). The iGSM benchmark is constructed under the assumption that ground-truth reasoning
854 chains follow a Directed Acyclic Graph (DAG) structure. It is important to note that our proposed
855 Tailor does not leverage any prior knowledge of this DAG representation. Implementation details,
856 including definitions of nodes and edges, their generation process, and the natural language trans-
857 lation procedure, are thoroughly introduced in the original paper (Ye et al., 2024). We provide a
858 sample query and rule-based CoT reasoning as an example below.

859 **Query:**

860 The number of each Grab & Go Store’s Canned Peaches equals the sum of each QuickMart’s
861 Corned Beef, each Colonial Quarter’s Product and each Colonial Quarter’s QuickMart.
862 The number of each QuickMart’s Canned Peaches equals 0.
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

The number of each QuickMart's Corned Beef equals 4 more than each QuickMart's Canned Peaches.
The number of each Colonial Quarter's QuickMart equals 6 more than the difference of each QuickMart's Corned Beef and each QuickMart's Product.
The number of each Vintage Architecture District's Grab & Go Store equals the sum of each QuickMart's Canned Peaches and each Colonial Quarter's Product.
The number of each Grab & Go Store's Corned Beef equals 0 times each QuickMart's Corned Beef.
The number of each Vintage Architecture District's QuickMart equals the sum of each Colonial Quarter's QuickMart and each Colonial Quarter's Product.
How many Product does each Vintage Architecture District have?

Rule-based CoT:

<think>

Let's compute the answer step by step. - According to the information given, the number of each QuickMart's Canned Peaches is 0. Let's call it D. So $D = 0$.

- Now, we can find the number of each QuickMart's Corned Beef. Let's call it b. Then $b = 4 + D = 4 + 0 = 4$.

- We can then calculate the number of each QuickMart's Product. Let's call it Y. Then $Y = b + D = 4 + 0 = 4$.

- We can then calculate the number of each Colonial Quarter's QuickMart. Let's call it d. Then $d = 6 + (b - Y) = 6 + (4 - 4) = 6$.

- Then, let's denote the number of each Colonial Quarter's Product as V. Then $V = d * Y = 6 * 4 = 24$.

- Then, let's denote the number of each Vintage Architecture District's Grab & Go Store as v. Then $v = D + V = 0 + 24 = 24$.

- Now, we can find the number of each Grab & Go Store's Canned Peaches. Let it be c. Then $c = b + V + d = 4 + 24 + 6 = 34$.

- We can then calculate the number of each Grab & Go Store's Corned Beef. Let it be l. Then $l = 0 * b = 0 * 4 = 0$.

- Now, we can find the number of each Grab & Go Store's Product. Let's denote it as r. Then $r = c + l = 34 + 0 = 34$.

- We can then calculate the number of each Vintage Architecture District's QuickMart. Let's call it i. Then $i = d + V = 6 + 24 = 30$.

- We can then calculate the number of each Vintage Architecture District's Product. Let's call it x. Then $x = v * r + i * Y = 24 * 34 + 30 * 4 = 936$.

Thus, the answer is 936.

< /think>

<answer>

The final answer is 936.

< /answer>

Difficulty Control. Both KK and iGSM offer controllable difficulty. In iGSM, difficulty is determined by the *number of operations* required to reach the correct answer in the ground-truth rule-based reasoning trace. In KK, difficulty is controlled by the *total number of characters* (knights and knaves) involved in the puzzle. For KK, the SFT dataset includes puzzles with 4–8 characters, and the RL dataset spans 7–11. The in-distribution and out-of-distribution (OOD) evaluations use puzzles with 7–11 and 12–13 characters, respectively. For the iGSM-medium task, the SFT dataset contains problems with 15–20 operations, and the RL dataset contains 15–20. The in-distribution and OOD evaluation sets use 15–20 and 21–25 operations, respectively. For the iGSM-hard task, the SFT dataset again includes 15–20 operations, while the RL dataset contains 25–30. The in-distribution and OOD evaluation sets use 25–30 and 31–35 operations, respectively.

A.2 TRAINING DETAILS

SFT Training. We perform supervised fine-tuning (SFT) on datasets curated using both Tailor and baseline methods. The key hyperparameters for SFT in the iGSM and KK tasks are summarized in Table 1. For dataset construction, we use 1,000 seed queries from the iGSM benchmark and 500 from the KK dataset.

Table 1: Configurations of SFT training

Configurations	Value
training epochs	4
global batch size	32
learning rate	5×10^{-6}
learning rate scheduler	constant
padding	not removed

RL training. We adopt DAPO (Yu et al., 2025c) for the main experiments using the Verl training framework (Sheng et al., 2025). The key hyperparameters for DAPO are listed in Table 2. For the GRPO experiments shown in Figure 9, we also set the overlong buffer length to 3,072 tokens and the overlong penalty factor to 1.0 to mitigate CUDA out-of-memory issues during RL training.

Table 2: Key hyperparameters for RL (DAPO) training

Configurations	Value
training epochs	5
batch size	$128 \times 16 = 2048$
sequence parallel size	2 (actor), 1 (ref)
gradient clipping	1.0
entropy coefficient	0.0
learning rate	1×10^{-6}
weight decay	0.1
warmup steps	10
optimizer	Adam
responses per prompt	16
max response length	4000 tokens
temperature	1.0
top-p	1.0
top-k	-1
KL loss coefficient	0.0
clip ratio (low)	0.2
clip ratio (high)	0.28
remove padding	True
overlong penalty factor	1.0
overlong buffer length	3072 tokens
rollout backend	vllm (Kwon et al., 2023)

RL reward verifiers. As described in Section 3, we use an outcome-based reward in the RL process. In both the KK and iGSM tasks, the final answer is extracted via string matching. For the KK experiments, the SFT demonstration template is as follows:

```
<think>
... Thinking process ...
</think>
<answer>
Thus, the final answer is boxed{ {name} is a {role}, ...}.
</answer>
```


972 Here, name refers to the character’s full name, and role refers to either Knave or Knight. To verify
 973 the correctness of the final answer, we examine each name–role pair using string matching. A
 974 completion is rewarded with $r = 1$ only if all roles are correctly assigned. Otherwise, the trace is
 975 labeled with $r = 0$.

976 For the iGSM task, we also use string matching in the reward function. The SFT demonstration
 977 template is as follows:
 978

```

979 <think>
980 ... Thinking process ...
981 </think>
982 <answer>
983 The final answer is {answer}.
984 </answer>
985

```

986 where the answer is always an integer. We assign $r = 1$ if the answer matches the ground-truth
 987 value, and 0 otherwise.
 988

989 **SFT curation pipeline details.** In the first step of our Tailor pipeline, where student models gener-
 990 ate demonstrations and the teacher model analyzes the traces, we provide the teacher with three ran-
 991 domly selected incorrect traces and one randomly selected correct trace. The teacher is instructed to
 992 analyze the failure cases and identify the reasoning behind the correct trace. We prompt the teacher
 993 model to generate Chain-of-Thought (CoT) reasoning enclosed between special tokens, followed by
 994 the generation of reasoning primitives (i.e., prompts used to curate the SFT dataset in the next step).
 995 The prompt used in this process is provided in Appendix B.1. Additionally, we include examples of
 996 primitives in Appendix B.2.
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

B PROMPTS

B.1 PROMPTS FOR TRACE ANALYSIS AND PRIMITIVE SYNTHESIS.

You are a large language model. Follow the instructions below carefully. Your goal is to generate instruction prompts that guide student models to produce high-quality reasoning.

(1) Below is a correct demonstration generated by a student model for the given query. Analyze the reasoning process and identify which behaviors or strategies are particularly effective and beneficial for the model’s reasoning.

Query (Correct Case): {correct query}
Correct CoT: {correct response}

(2) Below are three failure cases, each consisting of the original query, the incorrect response from the student model, and the expected ground-truth answer. Analyze why each response is incorrect by referencing specific steps or reasoning patterns that led to the failure. Based on this analysis, provide instructions on how to identify and revise the failed demonstration to arrive at the correct answer.

Failure Case 1:
Query: {fail1 query}
Response: {fail1 response}
Ground Truth: {fail1 gt}

Failure Case 2:
Query: {fail2 query}
Response: {fail2 response}
Ground Truth: {fail2 gt}

Failure Case 3:
Query: {fail3 query}
Response: {fail3 response}
Ground Truth: {fail3 gt}

(3) Based on the comparison between the correct and failed demonstrations, explain what kind of reasoning behavior is essential for robust and correct student model performance. You should explicitly go through each failure case one by one using the ground-truth answer, identify the correct reasoning pattern, uncover any implicit assumptions present in the correct reasoning path that lead to the correct final answers, and analyze how to correct the incorrect reasoning chains.

Then, new instruction prompts are designed to: (a) preserve the reasoning pattern exhibited in the correct demonstration; (b) incorporate the reasoning strategies and implicit assumptions necessary to reach the correct answers in the failure cases; and (c) guide the language model to self-identify and self-correct when similar failure patterns arise, steering it toward the correct reasoning path.

Please perform the entire thinking process described above within the `<prompt_think>...</prompt_think>` tag.

Please output the modified instruction prompt clearly inside `<generated_prompt>...</generated_prompt>` tags.

Please do not include any demonstration example inside `<generated_prompt>...</generated_prompt>`.

You should be aware that every query has a valid answer. So the generated prompt must encourage the language model to conclude a valid answer.

Do not include unrelated words such as `<|im_start| >` inside `<generated_prompt>...</generated_prompt>`, just instructions for solving the problem.

B.2 EXAMPLES OF SYNTHETIC REASONING PRIMITIVE

Example primitives for the iGSM Dataset. The generated primitive examples for solving this math dataset are shown in the following boxes with simplification for display. Example 1 aims to prevent errors where the model jumps to answers without properly resolving dependencies or dropping intermediate relationships. It addresses these issues by requiring the model to explicitly construct the dependency graph, clearly identify the target value to solve, and conduct a more detailed exami-

1080 nation, including recomputing key steps even when no errors are detected. Example 2 focuses on
 1081 preventing misinterpretation of the given conditions or overlooking parts of the problem description.
 1082 For instance, it helps the model avoid skipping constraints or variables that may appear redundant.
 1083 The most distinctive parts in Examples 1 and 2 are highlighted in red and blue, respectively. All
 1084 generated primitives are combined with format instructions to ensure the teacher model generations
 1085 comply with the required output format.

1086 **iGSM Example prompt 1:**

1087 Problem-Solving Instructions for Robust Reasoning

1088 1. Define Variables and Equations:

- 1089 - List every entity and attribute mentioned in the query.
- 1090 - Assign a unique variable to each quantity.
- 1091 - Translate all given statements into equations using these variables.

1092 2. Prioritize Independent Variables:

- 1093 - Solve variables with direct numerical assignments first (e.g., constants like 'X = 5').
- 1094 - Substitute these values into the dependent equations immediately.

1095 3. **Build Dependency Graphs:**

- 1096 - For each unsolved variable, list all equations where it appears.
- 1097 - **Identify the simplest equation to resolve next (e.g., least dependencies).**

1098 4. **Solve Step-by-Step:**

- 1099 - Proceed incrementally, substituting known values into dependent equations.

1100 5. **Validate Intermediate Results:**

- 1101 - **Recompute critical steps to verify consistency.**
- 1102 - Flag contradictions for re-evaluation.

1103 6. **Target-Focused Resolution:**

- 1104 - **Clearly state the goal variable.**
- 1105 - Back-substitute from the goal to ensure all required variables are resolved.

1106 7. Final Answer:

- 1107 - Conclude with a boxed numerical answer ('boxed{N}') supported by the validated rea-
 1108 soning chain.

1109 Note: All variables are solvable.

1110 **iGSM Example prompt 2:**

1111 Instructions for Solving the Problem:

1112 1. Define All Variables Explicitly:

- 1113 - Assign variables to each entity and relationship mentioned, **including composite terms**
 1114 **(e.g., "Enclosure" = sum of sub-components).**
- 1115 - Label all variables clearly (e.g., (X.Store) for "X at Store").

1116 2. Translate Statements into Equations:

- 1117 - Convert every given statement into a mathematical equation, **even if it seems redundant.**
- 1118 - Preserve all operations (sums, differences, multipliers) exactly as stated.

1119 3. Substitute Known Values Early:

- 1120 - Substitute fixed values (e.g., "equals 9") immediately to simplify equations.
- 1121 - Do not assume unconstrained variables are zero unless explicitly stated.

1122 4. Explore Indirect Relationships:

- 1123 - Track how variables influence others indirectly (e.g., if (A = B + C) and (B) de-
 1124 pends on (D), express (A) in terms of (D)).
- 1125 - If a variable's value is unresolved, check if it can be expressed in terms of other known
 1126 variables.

1127 5. Self-Correct for Consistency:

- 1128 - If equations lead to contradictions (e.g., (0 = 1)), revisit earlier steps to identify
 1129 missed relationships or incorrect assumptions.
- 1130 - **Verify that all given statements are fully utilized in the solution.**

1131 6. Conclude Rigorously:

- 1132 - Ensure every step logically follows from the previous ones.
- 1133 - **If ambiguity remains, exhaustively test plausible interpretations to understand the prob-
 lem and conditions (e.g., "Enclosure" as a sum) to arrive at a valid answer.**

1134 C USAGE OF LLMS
1135

1136 We primarily used publicly available LLMs to assist with proofreading and grammar refinement of
1137 the paper draft. All technical content was verified by the authors. Our research also involves LLMs
1138 as a core component: we distill LLMs to curate datasets and also fine-tune LLMs. The authors take
1139 full responsibility for all research ideas, technical contributions, and conclusions.
1140

1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187