

# AdaR: A Framework for Equipping LLMs with Adaptive Reasoning

Anonymous ACL submission

## Abstract

Mathematical reasoning is a primary indicator of large language models (LLMs) intelligence. However, existing LLMs exhibit failures of robustness and generalization. This paper attributes these deficiencies to spurious reasoning, wherein generated reasoning traces bear negligible causal connection to answers. To address this challenge, we propose the AdaR framework to equip LLMs with adaptive reasoning, wherein models rely on problem-solving logic to produce answers. AdaR automatically synthesizes logically equivalent queries by varying variable values, and trains models with Reinforcement Learning with Verifiable Rewards (RLVR) on these data to penalize spurious logic while encouraging adaptive logic. To ensure data quality, we extract the problem-solving logic from the original query and generate the corresponding answer by code execution and then apply sanity check. Experimental results demonstrate that AdaR achieves substantial improvement in mathematical reasoning while maintaining high data efficiency. Furthermore, even for advanced LLMs, there still exists robustness and generalization deficiencies, which our work effectively mitigates.

## 1 Introduction

Large Language Models (LLMs) have demonstrated strong performance across various reasoning tasks (Wei et al., 2022a). Among these, mathematical reasoning serves as a crucial cognitive skill that supports problem-solving across tasks (Huang and Chang, 2023). Beyond early direct inference attempts (Liu et al., 2021; Brown et al., 2020), Chain-of-Thought (CoT) reasoning has been recognized as an effective approach to enhance mathematical reasoning (Wei et al., 2022b), as it breaks down complex problems into manageable steps (Chu et al., 2023).

However, existing mathematical LLMs still exhibit failures at two levels: (i) robustness on in-

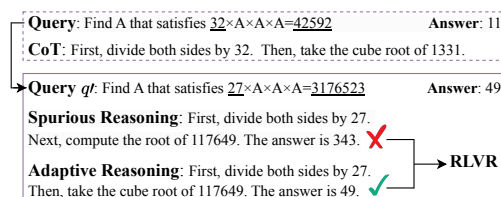


Figure 1: The examples of spurious and adaptive reasoning.

domain tasks (Mirzadeh et al., 2024); and (ii) generalization on out-of-domain tasks (Jahin et al., 2025). We attribute these deficiencies to **spurious reasoning**, the process by which LLMs derive gold answer  $y$  from superficial features but not the correct problem-solving logic  $L$ , therefore, producing reasoning trace (i.e. CoT) that is semantically similar to, yet logically inequivalent with correct traces. Consequently, even when the underlying problem-solving logic  $L$  remains unchanged, models relying on spurious reasoning fail to adapt to numerical changes of values in the variable set  $x$  and exhibit instability in performance. Meanwhile, spurious reasoning cannot be constructed from atom of thoughts (Teng et al., 2025). Therefore, it’s non-compositional along causal relations which in turn causes models to generalize ineffectively.

To overcome the mentioned weaknesses, the model is expected to understand the logic  $L$  as a function  $L(\cdot)$ , by correctly understanding the query  $q$  and disentangling  $x$  from  $q$ , which is an ability that reflects algebraic reasoning (Kieran, 2004). Applying the disentangled  $x$  to  $L(\cdot)$  may reveal the correct answer  $y$ . We introduce this process as **adaptive reasoning**. The query template is denoted as  $T$ , with  $x$  removed from  $q$ .

Learning adaptive reasoning from limited data is difficult. Thus previous work employ data synthesis (Yu et al., 2023; Lu et al., 2024). According to the above analysis, synthesizing mathematical data given an  $L$  have two essential dimensions: the

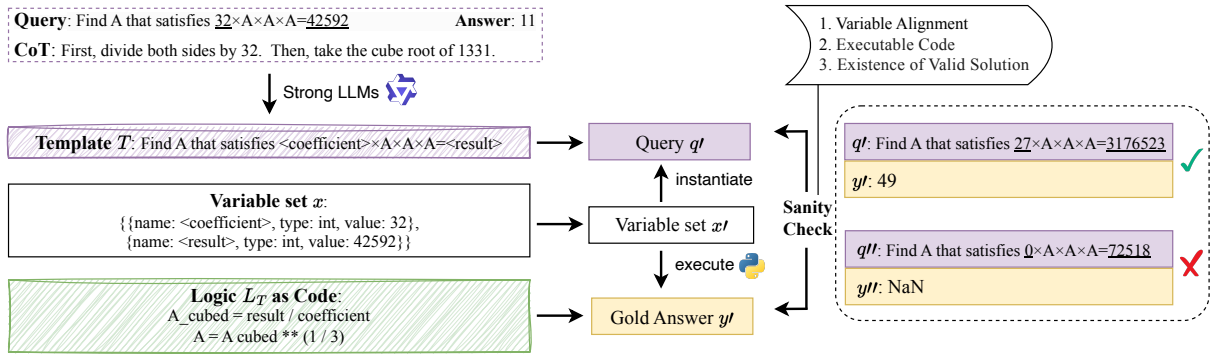


Figure 2: The pipeline of how we get the query-answer pairs by controllably perturbing variable values while preserving problem-solving logic and sanity.

query template  $T$ , and the variable set  $x$ . Yu et al. (2023) predominantly concentrated on perturbing  $T$  (e.g., paraphrasing) while keeping  $x$  fixed. However, training with varying  $T$  with fixed  $x$  may increase the risk of spurious reasoning overfitting to  $x$ . In contrast, Lu et al. (2024) attempted to prompt the LLM to modify  $x$ ; however, the employed strategy leaves  $x$  uncontrolled, and may also inadvertently alter  $L_T$ , making it difficult to ensure the correctness of the generated gold answer  $y$ .

In order to equip LLMs with **Adaptive Reasoning**, we propose the framework **AdaR**, which synthesizes data along both  $T$  and  $x$  dimensions and performs reinforcement learning with them. Crucially, AdaR introduces an automatic process for high-precision, controllable numerical perturbations, while ensuring verified  $y$ . To achieve this end, we decompose the complex perturbation task into the following manageable, verifiable sub-tasks. As shown in Figure 2, we first prompt an strong LLM to generate  $T$ ,  $L_T$  (as code, e.g. a Python program), and  $x$ . Subsequently, we *controllably perturb* values in  $x$  to predefined magnitudes and types. The perturbed variable set  $x'$  are then used to instantiate  $T$  to generate queries  $q'$ ; while be provided as input to  $L_T$ , which is executed to produce  $y'$ . Furthermore, we introduce *sanity check* to filter invalid instances.

AdaR then trains the model by Reinforcement Learning with Verifiable Rewards (RLVR) (Shao et al., 2024) to equip LLMs with adaptive reasoning. Notably, unlike in the single data situation, where outcome reward may inadvertently strengthen existing spurious reasoning, the correctness of outcomes on these synthetic data provides a reliable signal for inferring where their responses derive. In detail, responses that rely on spurious reasoning are more likely to produce incorrect answers on the perturbed synthetic data and are consequently

penalized in RLVR, thereby pushing the model to explore the adaptive problem-solving logic (Figure 1).

Extensive experiments across in-domain robust tasks and out-of-domain tasks demonstrate that AdaR achieves great gains (+8.50 points on average), with only 9K synthetic data on math-specialized and general base models, thereby demonstrating that our approach enhances the model’s robustness and generalization. Further analysis indicates that:

- Advanced LLMs still exhibit pronounced robustness and generalization deficiencies, which our work effectively mitigates.
- Perturbing the variable values is more effective than perturbing the query template for improving model’s performance.
- Evidence of heightened influence on logical order and improved algebraic thinking demonstrates that AdaR equips LLMs with adaptive reasoning

## 2 Method

In this section, we first present methods for synthesizing data that ensure controllability and sanity throughout the generation process. Following this, we introduce our training strategy, which is employed to more effectively integrate with synthetic data, preventing the model from learning spurious reasoning, thereby facilitating adaptive reasoning.

### 2.1 Data Synthesis with Executable Code and Verifiable Answers

A straightforward method for perturbing values in queries and obtaining correct answers without human annotation is to prompt a LLM (Wang et al.,

2025). However, our preliminary experiments indicate that this approach offers neither explicit control over the magnitudes and types of perturbations nor any guarantee of answer correctness. To synthesize the desired query-answer pairs, AdaR constrains the model to identify only the metadata within the data and to translate a CoT into code; controllable perturbations are then applied externally, and code execution is used to ensure answer correctness. The AdaR framework accordingly decomposes it into a sequence of manageable, verifiable sub-tasks as follows.

**Convert Logic in Text to Logic in Code.** Executable problem-solving code can serve as a problem-solving logic provided its correctness is guaranteed. Crucially, such code generalizes effectively: by substituting the input variable set, it can solve any perturbed query of the original query, for its output can serve as reliable gold answer. When CoT and gold answer are provided, problem-solving code generation becomes a straightforward process of translating the textual logic (i.e. CoT) into executable code, thereby reducing the model’s reasoning burden and enabling direct verification of correctness against the gold answer. Building on these insights, we provide a query  $q$ , the corresponding CoT  $z$  and gold answer  $y$  to an open-source LLM to synthesize the logic  $L_T$  as code for producing gold answer  $y'$  of subsequent perturbed query  $q'$ . Because code generation requires abstracting the concrete numerical values in  $q$  into a mapping to the input variables, we further instruct the model to produce  $T$  as a byproduct by applying this mapping into the original query  $q$ . Appendix A.3 details the prompt employed.

**Controllable Perturbation.** By comparing  $q$  with  $T$ , we construct a variable set  $x$  that records each variable’s name, numeric type and its value as it appears in the original query. This variable set, enriched with metadata, enables direct numeric perturbations at a predefined magnitude while ensuring numerical validity, thereby avoiding the unpredictability of perturbations performed through an LLM. More specifically, we then apply independent perturbations by sampling each variable’s value within a range of  $\pm\alpha\%$  of its original value. The parameter  $\alpha$  allows the control of the perturbation magnitude. To ensure numerical validity, both the numeric type and the sign of every variable must remain unchanged before and after perturbation. Formally, for the  $i$ -th variable with original

value  $x_0^i$ , we draw:

$$x^i = x_0^i \times (1 + \Delta_i), \Delta_i \sim \text{Uniform}(-\alpha\%, \alpha\%)$$

$$\text{s.t. type}(x^i) = \text{type}(x_0^i), \text{sign}(x^i) = \text{sign}(x_0^i) \quad (1)$$

**Sanity Check.** After applying perturbations, we instantiate  $T$  with  $x'$  to obtain  $q'$ , and execute  $L_T$  as code with the same  $x'$  as input to generate  $y'$ . Although providing  $z$  and  $y$  has reduced task difficulty,  $T$  and  $L_T$  as code generated by LLMs remain uncertain and may contain errors. Furthermore, a meaningful query imposes inter-variable constraints on its variables; independent perturbations can violate these constraints and thereby introduce errors. To ensure that the perturbed data remain well-posed, we conduct a *sanity check* along the following aspects:

- **Variable Alignment (VA).** We compare variables referenced in  $T$  with those used in each  $L_T$  as code. Any mismatch indicates potential errors (e.g. hallucination) in the LLM’s output.
- **Executable Code (EC).** Since  $L_T$  as code is used to derive the  $y$ , executability is a critical requirement: (i) the code runs without runtime errors; and (ii) providing the original variable set  $x_0$  as input reproduces the gold answer  $y_0$ .
- **Existence of Valid Solution (EVS).** As perturbations do not incorporate inter-variable constraints,  $q'$  may be invalid in realistic scenarios (e.g., selecting 20 items from a set of 10). To provide evidence that a valid solution exists, we perform cross-validation by comparing the answer generated by the code with the output of a mathematical LLM under perturbed query input. Nevertheless, because post-training can introduce spurious reasoning, the model may produce an incorrect answer when conditioned solely on the perturbed query. To handle this, we supply the corresponding  $L_T$  as code in the context as a hint, which has been verified by EC, enabling the model to optionally ground its reasoning on  $L_T$ .

If a perturbed instance fail to pass the *sanity check*, we reattempt *controllable perturbation*. If the number of attempts exceeds  $\tau$  times, we conclude this instance likely involves complex inter-variable constraints and then discard it from synthesis.

**Queries Paraphrasing.** We adopt paraphrasing (Yu et al., 2023) as a data augmentation strategy to increase the diversity of  $T$ . This approach complements the *controllable perturbation* introduced earlier to increase the diversity of  $x$ . The impact of these scalable dimensions on model performance is further explored in Section 3.3 and Section 4.2.

**Remark.** *Our synthetic data contains no CoT, it consists solely of the query and the corresponding gold answer. As will be described in the next subsection, CoTs generated through either spurious reasoning or adaptive reasoning can be sampled from the target LLM.*

## 2.2 Training Strategy

The training objective of Supervised Fine-Tuning (SFT) and its variant Rejection sampling Fine-Tuning (RFT) (Yuan et al., 2023) (detailed in Appendix A.2) makes models prone to memorizing provided CoTs rather than developing adaptive reasoning (Chu et al., 2025), thereby causing them to converge to local optimum (Kang et al., 2025). RLVR has recently been widely adopted to enhance generalization (Guo et al., 2025). However, its outcome reward generated by spurious reasoning is indistinguishable from that generated by adaptive reasoning, which may inadvertently reinforce spurious reasoning.

To address this, we combine RLVR with our synthetic data. RLVR’s reward-driven enables the model to learn from comparison among rewards obtained by solving perturbed queries using different reasoning process. More specifically, when given the original query  $q$ , it’s hard to determine where the gold answer  $y$  is derived. However, when the model is evaluated on perturbed query  $q'$ , reliance on CoT from spurious reasoning is more likely to yield incorrect outcome, whereas reliance on CoT from adaptive reasoning  $L_T(\cdot)$  is more likely to yield correct outcome. All perturbed queries are placed into the same batch. Then, each reasoning process is more likely to receive appropriate feedback, thereby promoting adaptive reasoning.

## 3 Experiment

### 3.1 Experimental Setup

**Data synthesis.** We use the Qwen2.5-72B-Instruct (Yang et al., 2024) as the open-source LLM to generate the query templates and the problem-solving codes. We select 9K instances from ORCA-MATH (Mitra et al., 2024) as seed data

for data synthesis. Using these seed data, AdaR synthesizes instances with a predefined magnitude  $\alpha = 500$  and a maximum of  $\tau = 50$  attempts. For each data in seed data, we select one corresponding synthetic data to construct the “ORCA-AdaR-train”, which contains 9K instances. We select another 2.5K instances to form “ORCA-AdaR-test”, ensuring no overlap with “ORCA-AdaR-train”. The details are shown in Appendix A.4.1.

**Models.** To prove the adaptability of our framework, we conduct experiments on two categories of base models: math specialized base LLM, specifically Qwen2.5-Math-7B (Yang et al., 2024) and DeepSeekMath-7B (Shao et al., 2024), and 8B general base LLM, specifically LLaMA3-8B (Grattafiori et al., 2024). In analysis, we further explore the deployment of AdaR on Instruct Model.

**Evaluation.** For a comprehensive evaluation of mathematical reasoning, we adopt 7 benchmarks covering both in-domain and out-of-domain evaluation. Specifically, we evaluate in-domain mathematical competence using GSM8K (Cobbe et al., 2021) and evaluate in-domain robustness using ORCA-AdaR-test (ORCA-AdaR) and GSM-SYM (main/p1/p2) (Mirzadeh et al., 2024). Within GSM-SYM, the main, p1, and p2 subsets exhibit increasing difficulty. For out-of-domain evaluation, we use MATH (Hendrycks et al., 2021), CollegeMath (Tang et al., 2024), TheoremQA (Chen et al., 2023), and American Invitational Mathematics Examination (AIME) problems for 2025 to evaluate generalization. Results are reported as pass@1 for all datasets except AIME 2025, for which, in accordance with (Yu et al., 2025), we report avg@32. Further details about the evaluation setup and benchmarks are provided in the Appendix A.4.3.

**Baseline.** We primarily compare our framework, AdaR, against the following baselines: (1) The base models undergo SFT using 9K seed data—this same data is also leveraged for data synthesis. These models are referred to as the “Initial SFT” models. (2) In “Standard-RLVR” setting, the “Initial SFT” is subsequently trained with RLVR on another 9K subset of ORCA-MATH, with its templates being disjoint from the training set used for “Initial SFT”. (3) We also include other synthesis methods, such as MetaMATH (Yu et al., 2023), which improves query diversity through paraphrasing and self-verification; MathGenie (Lu et al.,

Method	In-Domain					Out-of-Domain				AVG
	GSM8K	ORCA-AdaR	main	p1	p2	MATH	College	Theorem	AIME	
<b>Qwen2.5-MATH (7B MATH-Specialized Base Model)</b>										
Initial SFT	81.43	77.12	74.26	64.48	52.44	42.84	28.55	14.88	4.27	48.92
Standard-RLVR	86.96	81.80	81.34	74.32	64.68	61.74	29.48	20.38	10.94	56.85
MetaMATH	84.61	82.44	79.78	71.24	63.28	66.92	39.25	26.13	9.38	58.11
MathGenie	87.64	84.96	80.40	72.60	63.56	55.94	21.11	20.05	10.31	55.17
AdaR	91.81	86.08	89.72	82.74	73.24	75.90	48.62	37.13	14.27	<b>66.61</b>
<b>DeepSeekMath (7B MATH-Specialized Base Model)</b>										
Initial SFT	70.28	65.48	61.70	50.26	29.36	28.44	22.32	11.50	0.21	37.73
Standard-RLVR	81.43	72.56	74.76	64.34	40.00	39.26	21.61	20.38	0.73	46.12
MetaMATH	80.89	70.28	75.14	64.54	39.56	38.02	19.13	22.63	0.21	45.60
MathGenie	80.13	73.44	74.18	64.02	40.04	39.16	22.39	23.75	3.33	46.72
AdaR	81.43	75.44	76.56	64.94	42.00	41.18	24.23	25.25	3.33	<b>48.26</b>
<b>Llama3 (8B General Base Model)</b>										
Initial SFT	67.17	57.84	59.98	47.30	24.92	18.20	9.33	8.13	0.00	32.54
Standard-RLVR	73.09	55.68	67.06	51.58	25.16	19.02	9.01	9.00	0.00	34.40
MetaMATH	74.52	59.96	70.74	56.58	26.72	19.68	9.04	9.38	0.00	36.29
MathGenie	72.47	58.90	67.26	53.72	25.38	18.30	8.92	9.00	0.00	34.88
AdaR	77.77	61.52	73.96	56.90	30.28	20.48	9.62	9.63	0.00	<b>37.80</b>

Table 1: Performance comparison across In-Domain and Out-of-Domain mathematical benchmarks. The datasets “main”, “p1”, and “p2” are from GSM-SYM. Best results are highlighted in bold.

2024), which generate questions by applying back-translation to paraphrased responses. For all RLVR training procedures, we adopt DAPO (Yu et al., 2025) for its faster convergence and the elimination of the need for a value model, which together reduce training time and computational resource requirements. Additional details of the training setup are provided in the Appendix A.4.2.

### 3.2 Main Results

The main results are summarized in Table 1. We highlight the following three observations:

**Observation 1: Our method with a small amount of synthetic data yields substantial performance gains.** Using only 9K synthetic data, our method surpasses other methods across all base models. Compared to other synthetic data methods, we improve over MetaMATH by 8.50 points and over MathGenie by 11.44 points on average.

**Observation 2: Our method significantly enhances model robustness and generalization.** The enhancement manifests across three levels: (i) perturbed variable values in seen queries during training (ORCA-AdaR-test); (ii) perturbed variable values in unseen queries during training (GSM-SYM main/p1/p2, where main, p1, and p2 represent increasing levels of difficulty); (iii) out-of-domain data. Across all base models, we observe

gains on the first two levels compared to the sub-optimal MetaMATH (+4.66 / +4.75 points). At a finer level of granularity, we observe that “Initial SFT” rarely produces correct answers across all perturbations of variable values within a specific query template. In contrast, AdaR significantly improves this capability. Besides, on the third level we obtain larger improvements: when using Qwen2.5-MATH-7B as the base model, the average gain reaches +8.56 points. These results indicate that the model generates CoTs from adaptive reasoning rather than spurious reasoning, which better supports robustness and generalization.

**Observation 3: The effectiveness of our method correlates with the mathematical reasoning ability of the base model.** The number of additional rollouts required to obtain both correct and incorrect answers follows the ordering Qwen2.5-MATH(458.00) < DeepSeekMath(563.20) < LLama3(942.83), which is inversely correlated with their respective performance gains (+17.69, +10.53, and +5.26 points). It implies possessing sufficient mathematics-related knowledge is a necessary condition for sampling both positive and negative responses, thereby facilitating the learning of adaptive reasoning. The result also suggests that AdaR is complementary to the mathematical pre-training with a large amount of real or synthetic mathematical data.

Method	Strategy	Sanity Check			Paraphrase	In-Domain				Out-of-Domain
		VA	EC	EVS		ORCA-AdaR	main	p1	p2	
Initial SFT	-	-	-	-	-	77.12	74.26	64.48	52.44	22.63
	RFT	✓	✓	✓	✓	82.32	76.54	67.66	55.48	23.82
	RLVR	✗	✓	✓	✓	85.36	87.98	80.34	69.12	39.04
AdaR	RLVR	✓	✗	✓	✓	-	-	-	-	-
	RLVR	✓	✓	✗	✓	85.48	89.06	81.90	71.88	40.17
	RLVR	✓	✓	✓	✗	<b>87.56</b>	89.63	80.96	71.88	41.31
	RLVR	✓	✓	✓	✓	86.08	<b>89.72</b>	<b>82.74</b>	<b>73.24</b>	<b>43.98</b>

Table 2: Ablation Study. Best results are highlighted in bold.

### 3.3 Ablation study

The AdaR framework employs techniques to ensure the model can master adaptive reasoning from controllable and diverse synthetic data: (i) deploying the RLVR training strategy; (ii) introducing the *sanity check*, comprising Variable alignment (VA), Executable Code (EC), and Existence of Valid Solution (EVS); (iii) incorporating the paraphrase to further enhance the diversity of query templates.

As shown in Table 2, each technique within the AdaR framework is essential to the final performance. Notably, without EC, directly applying subsequent EVS retains only 0.2% data, which is insufficient to support training; consequently, no result is reported for this configuration. Although this indicates that EVS functionally contains EC, the rule-based EC is more computationally efficient than the model-based SE, by  $218\times$  (as shown in Appendix A.6), and thus remains essential. Besides, the result proves that paraphrase is complementary to perturbation of variable values, since it enhances the ability to understand different templates of specific problem-solving logic.

To further examine the importance of RLVR, we evaluate reasoning ability by measuring the accuracy of problem-solving code generation and computational ability by measuring the accuracy when the model itself is instructed to execute the code (the full prompt is provided in Appendix A.3). As summarized in the Table 3, relative to “Initial-SFT”, “AdaR-RFT” improves computation while degrading reasoning, whereas “AdaR” with RLVR yields substantial improvements in both reasoning and computation ability. It is consistent with Chu et al. (2025): RFT’s objective tends to memorize superficial features, thereby exacerbating overfitting to computation patterns rather than the corresponding problem-solving logic, while RLVR improves both abilities by encouraging the maximization of the exploration reward.

Method	ORCA-AdaR-test	Reasoning	Computation
Initial-SFT	77.12	75.24	54.92
Standard-RLVR	81.80	69.52	68.52
AdaR-RFT	82.32	72.36	75.36
AdaR	<b>86.08</b>	<b>80.96</b>	<b>90.76</b>

Table 3: Evaluation of reasoning and computation ability on ORCA-AdaR.

## 4 Analysis

In this subsection, we address the following Research Questions (RQs) regarding AdaR: (1) Does training with AdaR enable the model to master adaptive reasoning? (2) To what extent does the scaling along different dimensions (query template  $T$ , variable set  $x$ , and problem-solving logic  $L$ ) influence the model’s performance respectively? (3) During the RLVR stage, when employing synthetic data whose templates are unseen by the target LLM, does AdaR continue to exhibit strong performance? (4) To what extent does diversity in the computational difficulty levels of synthetic queries affect the model’s performance? (5) Is AdaR applicable to Instruct models? To answer these RQs (excluding RQ5), we conduct the following analysis using Qwen2.5-MATH-7B as the base model. Because the AIME evaluation metric differs from those of the other out-of-domain benchmarks, we exclude AIME from the analysis for better reporting.

### 4.1 AdaR Enables Adaptive Reasoning

**Spurious reasoning causes generalization deficiencies.** Spurious reasoning that relies on superficial features can derive answers without adhering to a correct logical order. By contrast, adaptive reasoning requires strict adherence to the correct logical order to derive answers step by step. To quantify this effect, we introduce the metric, Influence to Logical Order (ILO), which measures relative change rate in perplexity (PPL) of correct answers when the sentence order of the correspond-

ing CoT is shuffled.

$$\text{ILO} = \frac{1}{n} \sum_{i=1}^n \frac{|\text{PPL}(y | q, z) - \text{PPL}(y | q, R^i(z))|}{\text{PPL}(y | q, z)}, \quad (2)$$

where  $q$  denotes the query,  $z$  denotes the CoT,  $y$  denotes the answer,  $R$  denotes a random permutation function operating at the sentence level, and  $n$  denotes the number of random permutations (set to 5 to keep variance small). We find that the ILO of CoTs that are unable to generate correct answers across all perturbations of variable values within a specific query template is significantly lower than that of CoTs capable of producing correct answers across all such perturbations (114.24% vs. 221.87%). This result indicates that spurious reasoning (characterized by a lower ILO metric) is more likely to induce generalization deficiencies than adaptive reasoning (characterized by a higher ILO metric). Based on ILO, we show that AdaR significantly improves the adaptive reasoning capability over “Initial SFT” (150.49% vs. 119.22%).

**Enhancing algebraic thinking.** We further examine model outputs and observe that, despite no code data being provided during training, the proportion of CoTs containing structural code snippets (e.g. “... B + H = 157 - (23 + 41). Substituting the value of H into the equation, we get ...”) increases from 55% to 90% after training with AdaR (please refer to Appendix A.5 for more details), indicating the emergence of algebraic thinking, treating unknown and known variables on equal footing and solving queries via variable calculation. We then likewise evaluate reasoning ability by measuring the accuracy of problem-solving code generation which engages algebraic thinking. As shown in Table 3, although RLVR is employed in both, training on standard data (Standard-RLVR) decreases accuracy, whereas training on our synthetic data with queries which are semantically similar but instantiate different variables improves accuracy.

## 4.2 Effect of Scaling Different Dimensions

In AdaR, we modulate sampling frequency to scale data size along query template  $T$  and values in variable set  $x$ . In addition, we scale data size along problem-solving logic  $L$  by subsampling, which regulates diversity in common scenarios. Accordingly, we evaluate performance trends when scaling  $x$ ,  $T$  or  $L$  and show the result in Figure 3. We observe steady improvements in model performance

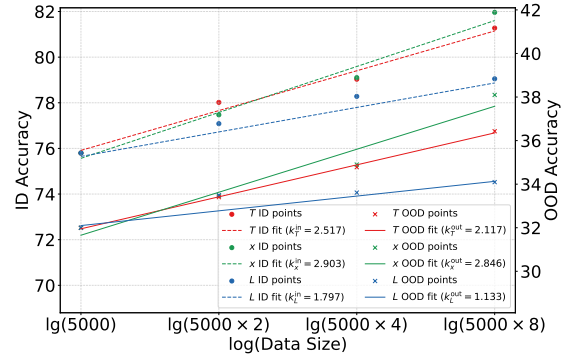


Figure 3: Performance of scaling.

with increasing data size in both in-domain (ID) and out-of-domain (OOD) settings. To characterize marginal returns with scale, we fit the scaling curves with a log-linear approximation. The fitted coefficients satisfy  $k_L < k_T < k_x$  throughout, indicating that scaling along the dimension  $x$  yields the greatest marginal returns. This is mainly because perturbing  $T$  facilitates query comprehension, whereas perturbing  $x$  is the essential driver of adaptive reasoning. We attribute the limited effectiveness of scaling along  $L$  to the simple subsampling, which causes data from same distribution constrains the control of diversity. We leave more comprehensive exploration of  $L$  to future work.

## 4.3 Effect of Seen Perturbed Queries

AdaR learns adaptive reasoning by comparing feedback on responses to perturbed queries whose problem-solving logic are consistent; consequently, the model should be exposed to these queries during the RLVR stage. We posit two situations where this comparison arises.

First, during SFT, memorization of superficial features associated with a given query induces spurious reasoning. When solving a perturbed query derived from that query in SFT, the model relies on spurious reasoning, producing unstable rollouts and thus yielding different feedback for comparison. It’s consistent with result shown in Table 4, training AdaR with perturbed data conditioned on SFT seed data leads to higher performance than “Standard-RLVR” which is trained on normal data.

Second, during RLVR, presenting the model with multiple perturbed queries that share consistent problem-solving logic but instantiate different valuable values can likewise induce comparison: if the model depends on spurious reasoning, it will receive different feedback across responses to these queries. To validate this, we additionally sample a 2.25K subset from the ORCA-MATH dataset with

Method	# Samples	In-Domain	Out-of-Domain
Standard-RLVR	9K	77.82	37.20
AdaR-Lite $\times 4$	2.25K $\times 4$	83.42	49.90
AdaR	9K	84.72	53.88

Table 4: Effect of seen perturbed queries.

no overlap with the SFT training data and apply 4 rounds of *controllable perturbation* for comparison. We denote this setting as “AdaR-Lite $\times 4$ ”. As reported in Table 4, “AdaR-Lite $\times 4$ ” outperforms “Standard-RLVR”, further proving this hypothesis.

#### 4.4 Effect of Perturbation Magnitude

We investigate the influence of perturbation magnitude and present the results in Figure 4. We observe that the performance increases with larger perturbation in variable values up to a point and then exhibits a slight decline. Larger perturbation magnitude enables the model to learn previously unseen numerical calculations, and thereby improves performance. However, when the generated queries contain an excessively high proportion of invalid data (e.g. selecting 2 items from a set of 10 is the original query, the maximum  $\alpha$  can be set is 500) before to the *sanity check* (refer to Section 2.1), it leads to more noise that ultimately degrades model performance.

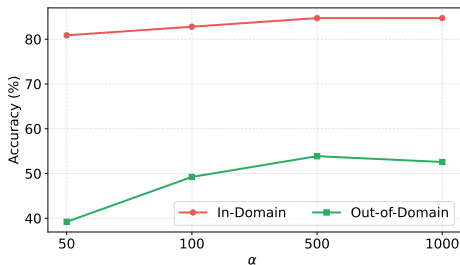


Figure 4: Influence of perturbation magnitude.

#### 4.5 AdaR is Applicable to Instruct Model

We further investigate using Qwen3-4B-Instruct-2507 (Yang et al., 2025) as the base model. AdaR still yields substantial performance gains (+5.86) and effectively mitigates robustness failures that persist under test-time scaling (+1.51). We additionally observe that the model’s rollout time on seed dataset is shorter than that of AdaR ( $\approx 0.36\times$ ), suggesting that generating correct answers becomes more difficult when the model is exposed to mild numerical perturbations, revealing pronounced brittleness. Taken together, these findings highlight a pronounced limitation in advanced LLMs: strong performance on conventional test sets often obscures fundamental deficiencies in

generalization and robustness. In contrast, AdaR effectively exposes and mitigates such latent deficiencies, demonstrating its broad applicability. Additional details are provided in Appendix A.7.

## 5 Related Work

### Analysis of Robustness and Generalization.

Wang et al. (2023) discovered that even invalid CoT demonstrations achieve comparable performance to valid ones. Similarly, we demonstrate that CoT arising from spurious reasoning is less sensitive to random variations in logical order. LLMs exhibit brittleness when facing problem variations: the GSM-SYM benchmark (Mirzadeh et al., 2024) demonstrates poor robustness on in-domain tasks with altered numerical values, while Jahin et al. (2025) shows limited generalization to out-of-domain problems, indicating models may memorize patterns rather than genuinely understand reasoning. The generation of GSM-SYM relies on human effort. In contrast, AdaR is fully automatic.

### Approaches to Improving Reasoning Ability.

Data synthesis has emerged as a promising solution (Wang et al., 2025), with Li et al. (2023) paraphrasing questions to diversify templates, though showing limited improvement in adaptive reasoning due to unchanged mathematical structures. Lu et al. (2024) advances this by augmenting datasets through LLM-based generation and verification, helping models identify shared logical patterns through contrasting perturbed queries, but remains constrained by simple perturbations and error-prone verification. More importantly, we propose adaptive reasoning, which can explain why the aforementioned methods are ineffective.

## 6 Conclusion

Robustness and generalization remain central challenges for LLMs when solving mathematical problems. We attribute these failures to spurious reasoning that relies on superficial features and encourage adaptive reasoning that can adapt to varying variable values. Therefore, we propose AdaR, a framework that enables adaptive reasoning and comprises a data synthesis component and a model training component. Experimental results demonstrate substantial performance improvements with a small amount of data on both in-domain and out-of-domain tasks. Further analyses indicate that AdaR indeed facilitates adaptive reasoning and is a scalable and broadly applicable framework.

## 640 Limitations

641 Our framework is subject to several limitations.  
642 Firstly, although Section 4.2 underscores the impor-  
643 tance of scaling  $L$ , AdaR currently enhances only  
644  $T$  and  $x$ . We therefore identify scaling the diversity  
645 of  $L$  as a promising direction for future work. Sec-  
646 ondly, AdaR requires the entire problem-solving  
647 process of each seed queries to be expressible as a  
648 fully executable program, and it assumes the pres-  
649 ence of easily perturbed input variables. These  
650 requirements implies applying AdaR in more gener-  
651 al settings (e.g., theorem proving) would require  
652 more sophisticated designs. Thirdly, due to com-  
653 putational constraints, we do not investigate large  
654 reasoning models or models with larger parameter  
655 counts as base.

## 656 References

657 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
658 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
659 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
660 Askell, and 1 others. 2020. Language models are  
661 few-shot learners. *Advances in neural information  
662 processing systems*, 33:1877–1901.

663 Wenhui Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan,  
664 Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony  
665 Xia. 2023. Theoremqa: A theorem-driven question  
666 answering dataset. *arXiv preprint arXiv:2305.12524*.

667 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Sheng-  
668 bang Tong, Saining Xie, Dale Schuurmans, Quoc V  
669 Le, Sergey Levine, and Yi Ma. 2025. Sft mem-  
670 orizes, rl generalizes: A comparative study of  
671 foundation model post-training. *arXiv preprint  
672 arXiv:2501.17161*.

673 Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang  
674 Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu,  
675 Bing Qin, and Ting Liu. 2023. Navigate through enig-  
676 matic labyrinth a survey of chain of thought reason-  
677 ing: Advances, frontiers and future. *arXiv preprint  
678 arXiv:2309.15402*.

679 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
680 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
681 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
682 Nakano, and 1 others. 2021. Training verifiers  
683 to solve math word problems. *arXiv preprint  
684 arXiv:2110.14168*.

685 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
686 Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
687 Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,  
688 Alex Vaughan, and 1 others. 2024. The llama 3 herd  
689 of models. *arXiv preprint arXiv:2407.21783*.

690 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao  
691 Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-  
692 rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.

Deepseek-r1: Incentivizing reasoning capability in  
llms via reinforcement learning. *arXiv preprint  
arXiv:2501.12948*. 693  
694  
695

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul  
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-  
cob Steinhardt. 2021. Measuring mathematical prob-  
lem solving with the math dataset. *arXiv preprint  
arXiv:2103.03874*. 696  
697  
698  
699  
700

Jie Huang and Kevin Chen-Chuan Chang. 2023. To-  
wards reasoning in large language models: A survey.  
In *61st Annual Meeting of the Association for Com-  
putational Linguistics, ACL 2023*, pages 1049–1065.  
Association for Computational Linguistics (ACL). 701  
702  
703  
704  
705

Afrar Jahin, Arif Hassan Zidan, Wei Zhang, Yu Bao, and  
Tianming Liu. 2025. Evaluating mathematical rea-  
soning across large language models: A fine-grained  
approach. *arXiv preprint arXiv:2503.10573*. 706  
707  
708  
709

Feiyang Kang, Michael Kuchnik, Karthik Padthe, Marin  
Vlastelica, Ruoxi Jia, Carole-Jean Wu, and Newsha  
Ardalani. 2025. Quagmires in sft-rl post-training:  
When high sft scores mislead and what to use instead.  
*arXiv preprint arXiv:2510.01624*. 710  
711  
712  
713  
714

Carolyn Kieran. 2004. Algebraic thinking in the early  
grades: What is it. *The mathematics educator*,  
8(1):139–151. 715  
716  
717

Chengpeng Li, Zheng Yuan, Hongyi Yuan, Guanting  
Dong, Keming Lu, Jiancan Wu, Chuanqi Tan, Xi-  
ang Wang, and Chang Zhou. 2023. Mugglemath:  
Assessing the impact of query and response aug-  
mentation on math reasoning. *arXiv preprint  
arXiv:2310.05506*. 718  
719  
720  
721  
722  
723

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,  
Hiroaki Hayashi, and Graham Neubig. 2021. **Pre-  
train, prompt, and predict: A systematic survey of  
prompting methods in natural language processing.**  
*Preprint*, arXiv:2107.13586. 724  
725  
726  
727  
728

Ilya Loshchilov and Frank Hutter. 2017. Decou-  
pled weight decay regularization. *arXiv preprint  
arXiv:1711.05101*. 729  
730  
731

Zimu Lu, Aojun Zhou, Houxing Ren, Ke Wang,  
Weikang Shi, Junting Pan, Mingjie Zhan, and Hong-  
sheng Li. 2024. Mathgenie: Generating synthetic  
data with question back-translation for enhancing  
mathematical reasoning of llms. *arXiv preprint  
arXiv:2402.16352*. 732  
733  
734  
735  
736  
737

Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi,  
Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar.  
2024. Gsm-symbolic: Understanding the limitations  
of mathematical reasoning in large language models.  
*arXiv preprint arXiv:2410.05229*. 738  
739  
740  
741  
742

Arindam Mitra, Hamed Khanpour, Corby Rosset, and  
Ahmed Awadallah. 2024. Orca-math: Unlocking  
the potential of slms in grade school math. *arXiv  
preprint arXiv:2402.14830*. 743  
744  
745  
746

747	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	802
748	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	803
749	Zhang, YK Li, Yang Wu, and 1 others. 2024.	Gao, Chengen Huang, Chenxu Lv, and 1 others.	804
750	Deepseekmath: Pushing the limits of mathematical	2025. Qwen3 technical report. <i>arXiv preprint</i>	805
751	reasoning in open language models. <i>arXiv preprint</i>	<i>arXiv:2505.09388</i> .	806
752	<i>arXiv:2402.03300</i> .		
753	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,	807
754	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong	808
755	Lin, and Chuan Wu. 2024. Hybridflow: A flexible	Tu, Jingren Zhou, Junyang Lin, and 1 others. 2024.	809
756	and efficient rlhf framework. <i>arXiv preprint</i>	Qwen2. 5-math technical report: Toward mathe-	810
757	<i>arXiv:2409.19256</i> .	matical expert model via self-improvement. <i>arXiv</i>	811
		<i>preprint arXiv:2409.12122</i> .	812
758	Zhengyang Tang, Xingxing Zhang, Benyou Wang, and	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu,	813
759	Furu Wei. 2024. Mathsacle: Scaling instruction	Zhengying Liu, Yu Zhang, James T Kwok, Zhen-	814
760	tuning for mathematical reasoning. <i>arXiv preprint</i>	guo Li, Adrian Weller, and Weiyang Liu. 2023.	815
761	<i>arXiv:2403.02884</i> .	Metamath: Bootstrap your own mathematical ques-	816
		tions for large language models. <i>arXiv preprint</i>	817
762	Fengwei Teng, Quan Shi, Zhaoyang Yu, Jiayi Zhang,	<i>arXiv:2309.12284</i> .	818
763	Yuyu Luo, Chenglin Wu, and Zhijiang Guo. 2025.	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	819
764	Atom of thoughts for markov llm test-time scaling.	Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,	820
765	<i>arXiv preprint arXiv:2502.12018</i> .	Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo:	821
766	Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu,	An open-source llm reinforcement learning system	822
767	and Junxian He. 2024. <a href="#">Dart-math: Difficulty-aware</a>	at scale. <i>arXiv preprint arXiv:2503.14476</i> .	823
768	<a href="#">rejection tuning for mathematical problem-solving.</a>		
769	<i>Preprint</i> , arXiv:2407.13690.	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting	824
		Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and	825
770	Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen,	Jingren Zhou. 2023. <a href="#">Scaling relationship on learning</a>	826
771	You Wu, Luke Zettlemoyer, and Huan Sun. 2023.	<a href="#">mathematical reasoning with large language models.</a>	827
772	<a href="#">Towards understanding chain-of-thought prompting:</a>	<i>Preprint</i> , arXiv:2308.01825.	828
773	<a href="#">An empirical study of what matters.</a> In <i>Proceedings</i>	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	829
774	<i>of the 61st Annual Meeting of the Association for</i>	Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.	830
775	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	2024. <a href="#">Llamafactory: Unified efficient fine-tuning</a>	831
776	pages 2717–2739, Toronto, Canada. Association for	<a href="#">of 100+ language models.</a> In <i>Proceedings of the</i>	832
777	Computational Linguistics.	<i>62nd Annual Meeting of the Association for Computa-</i>	833
778	Zaitian Wang, Pengfei Wang, Kunpeng Liu, Pengyang	<i>tional Linguistics (Volume 3: System Demonstra-</i>	834
779	Wang, Yanjie Fu, Chang-Tien Lu, Charu C. Aggar-	<i>tions)</i> , Bangkok, Thailand. Association for Computa-	835
780	wal, Jian Pei, and Yuanchun Zhou. 2025. <a href="#">A com-</a>	tional Linguistics.	836
781	<a href="#">prehensive survey on data augmentation.</a> <i>Preprint</i> ,		
782	arXiv:2405.09591.	<b>A Appendix</b>	837
783	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<b>A.1 The use of LLMs</b>	838
784	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	We use LLMs (e.g., GPT-5) only to polish writing.	839
785	and 1 others. 2022a. Chain-of-thought prompting	More specifically, their application focuses on two	840
786	elicits reasoning in large language models. <i>Advances</i>	key areas: correcting grammatical errors and sug-	841
787	<i>in neural information processing systems</i> , 35:24824–	gesting more appropriate word choices to enhance	842
788	24837.	expression. Additionally, we conduct a thorough	843
789	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	double check of all content refined by LLMs. This	844
790	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	verification process is critical to preventing the in-	845
791	and 1 others. 2022b. Chain-of-thought prompting	clusion of harmful information and ensuring the	846
792	elicits reasoning in large language models. <i>Advances</i>	overall accuracy, reliability, and appropriateness of	847
793	<i>in neural information processing systems</i> , 35:24824–	the final content.	848
794	24837.	<b>A.2 Training Objective</b>	849
795	Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An,	<b>A.2.1 SFT &amp; RFT</b>	850
796	Zhenyu Duan, Yimin Du, Junchen Liu, Tanglifu Tan-	Supervised fine-tuning (SFT) is a mainstream post-	851
797	glifu, Xiaowei Lv, and 1 others. 2025. <a href="#">Light-rl:</a>	training strategy for eliciting step-by-step reason-	852
798	<a href="#">Curriculum sft, dpo and rl for long cot from scratch</a>	ing. Let $\pi_\theta$ denote a policy over responses given	853
799	and beyond. In <i>Proceedings of the 63rd Annual Meet-</i>		
800	<i>ing of the Association for Computational Linguistics</i>		
801	<i>(Volume 6: Industry Track)</i> , pages 318–327.		

854 queries, parameterized by  $\theta$ . Given a dataset  $\mathcal{D}$   
855 of query–response pairs  $(q, r)$ , SFT minimizes the  
856 negative log-likelihood objective

$$857 \quad \mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(q,r) \sim \mathcal{D}} [\log \pi_{\theta}(r | q)], \quad (3)$$

858 Reject sampling Fine-Tuning (RFT) constructs an  
859 SFT dataset by sampling from the model to be  
860 trained and retaining high-scoring outputs.

### 861 A.2.2 RLVR

862 Reinforcement Learning with Verifiable Reward  
863 (RLVR) is a reinforcement learning method that  
864 saves video memory usage. Given a query  $q$ ,  
865 the model samples  $r \sim \pi_{\theta}(\cdot | q)$  and a verifier  
866  $v(q, r) \in [0, 1]$  evaluates it. RLVR maximizes

$$867 \quad J(\theta) = \mathbb{E}_{q \sim \mathcal{D}, r \sim \pi(\cdot | q)} [v(q, r)] \quad (4)$$

868 In mathematical reasoning,  $v$  typically checks  
869 whether the predicted answer exactly matches the  
870 gold answer.

### 871 A.3 Prompts

872 We show the prompts used for generating the query  
873 template and problem-solving code in Prompt 1,  
874 Solution Existence in Prompt 2, code generation  
875 which evaluates reasoning ability in Prompt 3, code  
876 execution which evaluates computational ability in  
877 Prompt 4.

## 878 A.4 General Settings

### 879 A.4.1 Data Synthesis

880 Given that most publicly available math word  
881 problem datasets are constructed based on  
882 GSM8K (Cobbe et al., 2021), and considering  
883 that existing LLMs may have been trained on sub-  
884 stantial variations of this dataset, which alleviates  
885 to some extent the model’s shortcut learning is-  
886 sues, we have opted for a newer and larger dataset  
887 ORCA-MATH (Mitra et al., 2024). We sample the  
888 data for synthesis by uniformly selecting using 42  
889 as the random seed.

890 We use Qwen2.5-72B for data synthesis, includ-  
891 ing problem-solving codes, query templates, the  
892 existence of valid solutions, and implementations  
893 of other baselines. The temperature is set to 0.7, the  
894 top p is set to 0.95, and the maximum generation  
895 length is 4096 tokens.

896 In addition, we experiment with using Qwen2.5-  
897 7B for data synthesis under the same parameter  
898 settings, with the results reported in the Table 7.

899 Although employing a weaker model for data syn-  
900 thesis leads to a reduction in overall model perfor-  
901 mance, we observe that by scaling data size along  
902 the variable value  $x$ , the performance gap can be  
903 reduced from 8.14 to 3.09.

### 904 A.4.2 Training

905 All baselines are first trained with the same SFT  
906 procedure, after which RLVR training is performed  
907 using their respective data. For SFT, we use  
908 LLaMA-Factory (Zheng et al., 2024). All mod-  
909 els are fine-tuned for 3 epochs with a batch size of  
910 128 on 4 NVIDIA A100 GPUs. The peak learn-  
911 ing rate is 1e-5 with a linear warm-up over the  
912 first 3% of training steps, followed by cosine decay  
913 to a minimum of 1e-7. The maximum generation  
914 length is set to 4096 tokens. For DAPO, we use  
915 verRL (Sheng et al., 2024), adopting the training  
916 setup of Yu et al. (2025). We utilize the AdamW  
917 optimizer (Loshchilov and Hutter, 2017) with a  
918 constant learning rate of 1e-6 and a linear warm-  
919 up over 20 rollout steps. For rollout, the prompt  
920 batch size is 256 and we sample 16 responses per  
921 prompt. For training, the mini-batch size is 32,  
922 corresponding to 128 gradient updates per rollout  
923 step. Overall, SFT training requires approximately  
924 12 hours and RL training requires approximately 2  
925 days.

### 926 A.4.3 Evaluation

927 We compare AdaR with baselines on the following  
928 7 benchmarks:

- 929 • **GSM8K** [MIT] (Cobbe et al., 2021): The test  
930 set comprises 1,319 high-quality grade-school  
931 mathematics word problems, each requiring  
932 between 2 and 8 reasoning steps.
- 933 • **ORCA-AdaR-test**: The test set consists of  
934 2,500 high-quality grade-school mathematics  
935 word problems synthesized by AdaR. Seed  
936 problems are selected from ORCA [MIT Li-  
937 cense] (Mitra et al., 2024) and subjected to  
938 1-4 numerical perturbations to assess models’  
939 robustness in mathematical reasoning.
- 940 • **GSM-SYM** [CC-BY-4.0] (Mirzadeh et al.,  
941 2024): The test set includes three sub-  
942 sets—*main*, *p1*, and *p2*—with 5,000, 5,000,  
943 and 2,500 instances, respectively. Starting  
944 from the GSM8K test problems, a symbolic  
945 template is constructed for each problem; each  
946 template yields 50 instances. The subsets aug-  
947 ment the original logical structure by adding

0 (*main*), 1 (*p1*), or 2 (*p2*) additional reasoning clauses, thereby enabling a more rigorous evaluation of mathematical reasoning robustness.

- **MATH [MIT]** (Hendrycks et al., 2021): The test set comprises 5,000 problems drawn from high-school mathematics competitions. Problems are categorized into seven types (Prealgebra, Intermediate Algebra, Algebra, Precalculus, Geometry, Counting & Probability, and Number Theory) and five difficulty levels.
- **CollegeMath** (Tang et al., 2024): The test set contains 2,818 college-level problems curated from nine college-level mathematics textbooks, covering seven key disciplines: Algebra, Precalculus, Calculus, VectorCalculus, Probability, LinearAlgebra, and Differential Equations.
- **TheoremQA** (Chen et al., 2023): A theorem-driven question-answering benchmark containing 800 problems grounded in 350 theorems, designed to evaluate LLMs’ ability to apply domain-specific theorems across Mathematics, Physics, Electrical Engineering, Computer Science, and Finance.
- **AIME 2025**: A test set containing 30 problems from the 2025 American Invitational Mathematics Examination (AIME), curated to evaluate LLMs on challenging, Olympiad-level high-school mathematics across Algebra, Geometry, Number Theory, and Combinatorics.

We adopt the evaluation pipeline of Tong et al. (2024) with some modifications. Unless otherwise noted, we set the sampling temperature to 0.7 and the nucleus parameter  $\text{top}_p$  to 0.9.

Regarding evaluation metrics,  $\text{pass}@1$  is defined as the accuracy of the first sampled output, whereas  $\text{avg}@32$  is the mean accuracy computed over 32 sampled outputs. As shown in Table 5, when the test set is sufficiently large (approximately 1,000 samples), the performance differences between  $\text{pass}@1$  and  $\text{avg}@32$  remain within 1 percentage point and the  $\text{avg}@32$  results fully consistent with the conclusions drawn from the  $\text{avg}@1$  results.

### A.5 Structural Code Snippet

To calculate the frequency of structural code snippet in the responses, we inspect 20 responses to

Method	GSM8K	GSM-SYM-main	MATH
Initial SFT	81.43 / 80.42	74.26 / 73.89	42.84 / 42.57
Standard-RLVR	86.96 / 86.88	81.34 / 81.51	61.71 / 61.29
MetaMATH	84.61 / 83.80	79.78 / 78.88	66.92 / 66.29
AdaR	91.81 / 92.32	89.72 / 89.39	75.90 / 76.00

Table 5:  $\text{pass}@1$  /  $\text{avg}@32$  results across different benchmarks.

queries from GSM-SYM. We present cases of structural code snippets from “Initial-SFT” and AdaR in Figure 8.

### A.6 Generation Cost

The generation costs for all aspects included in the *sanity check* are reported in Table 6. This experiment is conducted on a single compute node.

### A.7 Performance on Instruct Models

**Dataset.** To evaluate whether AdaR remains effective on more challenging datasets, beyond the grade-school mathematics problems used in earlier experiments (i.e., ORCA dataset), we adopt the Light-R1 dataset (Wen et al., 2025). From this dataset, we sample 7.5K seed data and apply numerical perturbations 1–4 times, resulting in a 17K-instance training set, denoted as “R1-AdaR-train”. To perform a fine-grained analysis of whether AdaR improves model robustness and to avoid potential of data leakage, we construct two test sets, each containing 1K instances: “seen”, whose query templates appear in “R1-AdaR-train”, and “unseen”, whose templates do not appear in the “R1-AdaR-train”.

**Baselines.** In this experiment, we introduce two strong baselines: “Qwen3-4B-Instruct-2507” and the test-time-scaling-enabled “Qwen3-4B-Thinking-2507”. Because the Instruct model already possesses strong mathematical reasoning capabilities, we skip the SFT stage and directly perform RLVR training. The model trained using RLVR on the 7.5K seed data is referred to as “Standard-RLVR-Instruct”, while the model trained with AdaR synthetic data is denoted “AdaR-Instruct”.

Table 8 shows the full results.

### A.8 Training Dynamic Analysis

As shown in Figure 5 and Figure 6, during RLVR training, as the training steps increase, the model’s generated CoT sequences grow progressively longer. In the early stage, although the sequence length steadily increases, the model rapidly

1037 acquires stable short reasoning patterns, often ac-  
1038 companied by certain forms of spurious reasoning,  
1039 resulting in a decrease in entropy. As training con-  
1040 tinues and CoT sequences further lengthen, the  
1041 model begins to develop more adaptive reasoning  
1042 due to exposure to perturbed data. This transition  
1043 increases the model’s predictive uncertainty, lead-  
1044 ing to the subsequent rise in entropy.

1045 In parallel, Figure 7 shows that the AIME 2024  
1046 scores exhibits a rapid improvement during the  
1047 initial training stage and then plateaus with minor  
1048 fluctuations in later steps. This trajectory aligns  
1049 with entropy dynamics: the model first stabilizes  
1050 around short reasoning patterns that yield quick  
1051 performance gains, and later explores a broader  
1052 reasoning space as CoT sequences grow, leading  
1053 to small oscillations around a higher performance  
1054 level.

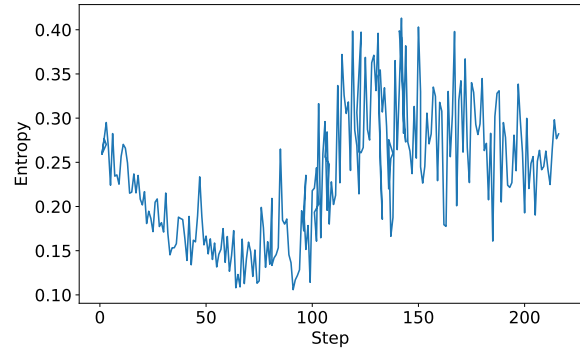


Figure 5: Entropy of actor model.

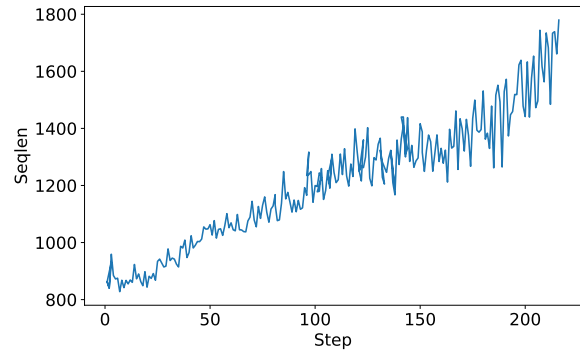


Figure 6: Generated sequence length of actor model.

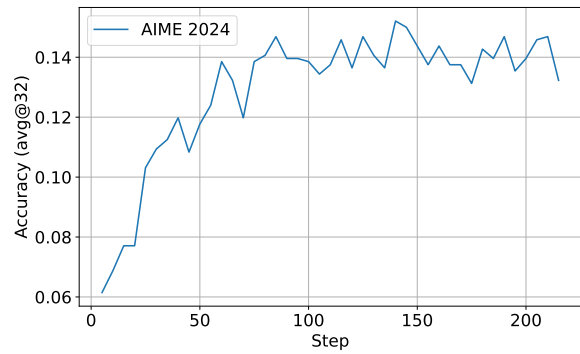


Figure 7: AIME 2024 scores on AdaR during training.

<p><math>q_i</math>: A group of 157 students has various hobbies. 23 like to play video games, 41 like to play volleyball, and the rest like to either hike or paint. How many like to hike if the number that like to paint is twice the number that prefer playing volleyball?</p> <p><math>y_i</math>: 11</p>	
<p><math>z_i</math> from "Inital-SFT":</p> <p>... we multiply the number of volleyball players by 2:</p> <p>Number of volleyball players = 41  Number of students who like to paint = <math>2 * 41 = 82</math></p> <p>...</p> <p>We subtract the sum of these from the total number of students to find out how many like to hike</p> <p>...</p> <p>Number of students who like to hike = <math>157 - (23 + 41 + 82)</math></p> <p>...</p> <p>Number of students who like to hike = 11...</p>	<p><math>z_i</math> from AdaR:</p> <p>Let's denote the number of students who like to paint as <math>(P)</math> and the number of students who like to hike as <math>(H)</math></p> <p>...</p> <p>The number of students who like to hike is twice the number that prefer playing volleyball, so <math>(P = 2 \times 41 = 82)</math></p> <p>...</p> <p>Therefore, the number of students who like to hike and paint is the total number of students minus the number of students who like to play video games and play volleyball.</p> <p>So, we can write the equation:</p> <p><math>(H + P = 157 - (23 + 41))</math></p> <p>Substituting the value of <math>(P)</math> into the equation, we get:  <math>(H + 82 = 157 - (23 + 41))</math>  <math>(H + 82 = 157 - 64)</math>  <math>(H + 82 = 93)</math></p> <p>Now, we solve for <math>(H)</math>:</p> <p><math>(H = 93 - 82)</math>  <math>(H = 11)</math>...</p>
<p><math>q_j</math>: A group of 183 students has various hobbies. 42 like to play video games, 25 like to play volleyball, and the rest like to either hike or paint. How many like to hike if the number that like to paint is twice the number that prefer playing volleyball?</p> <p><math>y_j</math>: 66</p>	
<p><math>z_j</math> from "Inital-SFT":</p> <p>... Since the number of students who like to paint is twice the number that prefer playing volleyball, we can calculate</p> <p>...</p> <p>Number of students who like to paint = <math>2 * 25 = 50</math></p> <p>...</p> <p>We can subtract these from the total to find out how many like to hike</p> <p>...</p> <p>Number of students who like to hike = <math>183 - (42 + 50)</math></p> <p>...</p> <p>Number of students who like to hike = 91...</p>	<p><math>z_j</math> from AdaR:</p> <p>Let's denote the number of students who like to paint as <math>(P)</math> and the number of students who like to hike as <math>(H)</math></p> <p>...</p> <p>The number of students who like to paint is twice the number that prefer playing volleyball, so <math>(P = 2 \times 25 = 50)</math></p> <p>...</p> <p>Therefore, the number of students who like to hike or play volleyball is the total number of students minus the number of students who like to play video games and the number of students who like to play volleyball. This can be expressed as:</p> <p><math>(H + P = 183 - 42 - 25)</math></p> <p>Substituting the value of <math>(P)</math> into the equation, we get:  <math>(H + 50 = 183 - 42 - 25)</math></p> <p>Now, let's calculate the right-hand side of the equation:  <math>(183 - 42 = 141)</math>  <math>(141 - 25 = 116)</math></p> <p>So, we have:  <math>(H + 50 = 116)</math></p> <p>To find <math>(H)</math>, we subtract 50 from both sides of the equation:  <math>(H = 116 - 50)</math>  <math>(H = 66)</math>...</p>

Figure 8: The case study of structural text in outputs. The green background indicates the correct reasoning step. The red background indicates the wrong reasoning step.

Sanity Check	Generation Times (s per 1K samples)	Retention rate (%)
Variable Alignment	2.90	67.70
Executable Code	2.80	70.70
Solution Existence	612	95.17

Table 6: Generation cost of component in *sanity check*.

Method	# Samples	In-Domain					Out-of-Domain				AVG
		GSM8K	ORCA-AdaR	main	p1	p2	MATH	College	Theorem	AIME	
Initial SFT	0K	81.43	77.12	74.26	64.48	52.44	42.84	28.55	14.88	4.27	48.92
Qwen2.5-7B	9K	87.34	83.72	81.02	73.84	62.08	64.36	35.42	25.13	13.33	58.47
	9K × 4	89.31	85.36	88.50	80.72	70.76	70.70	41.34	30.13	14.90	63.52
Qwen2.5-72B	9K	91.81	86.08	89.72	82.74	73.24	75.90	48.62	37.13	14.27	66.61
	9K × 4	92.72	87.16	90.32	83.50	73.80	76.20	49.96	37.50	16.46	67.51

Table 7: Performance comparison when employing different models for data synthesis.

Method	Robustness			Generalization					AVG	CoT Length
	seen	unseen	GSM-SYM	GSM8K	MATH	College	Theorem	AIME		
Qwen3-4B-Instruct-2507	60.99	56.62	85.50	92.27	91.12	50.85	47.50	46.46	66.41	1672
Qwen3-4B-Thinking-2507	66.32	63.42	91.19	93.33	<b>95.38</b>	<b>53.97</b>	<b>58.13</b>	<b>75.83</b>	<b>74.70</b>	4784
Standard-RLVR-Instruct	62.86	57.72	90.47	93.78	91.52	52.09	48.25	52.50	68.65	2355
AdaR-Instruct	<b>68.42</b>	<b>65.17</b>	<b>91.88</b>	<b>94.69</b>	90.83	53.30	55.00	58.83	72.27	3435

Table 8: Comparison with the Instruct model and Thinking model.

## Prompt 1: Generate Template and Code

### Task Description:

You are given a natural language query and its chain-of-thought response.

Your task is to: Generate a Query Template by abstracting specific values into variables.

Generate Python Code that executes the logic described in the COT response using the abstracted variables.

### Input Format:

Query: Original query with specific values

Response: Chain-of-thought reasoning that leads to the answer

### Output Requirements:

#### Query Template:

Replace only concrete values in the query with angle-bracketed placeholders like <variable\_name>. Do not replace names or general nouns (e.g., do not change “Jungkook” to <person\_name>). Preserve the original wording and structure of the query as much as possible.

#### Python Code:

Begin by defining variables that correspond to the placeholders in the template. Translate the logic in the response into executable Python code. The code should end with a print() statement that prints only the final result. Do not include comments with explanations or reasoning. Use the same variable names as in the template for consistency.

```
=== START EXAMPLE ===
```

```
{example}
```

```
=== END EXAMPLE ===
```

```
### Query:
```

```
{query}
```

```
### Response:
```

```
{response}
```

### Prompt 2: Existence of Valid Solution

#### Task Description:

Your task is to generate a Chain-of-Thought (CoT) explanation that answers the user's question by reasoning through the logic implied in a provided Python script. Use the script to inform your explanation, but do not output or reproduce any code.

#### Input Format:

Query: A question involving specific values or conditions.

Python Code: A script that solves the query or provides a key computational procedure.

#### Output Requirements:

Start by interpreting the question clearly. Reason through the problem step by step, using the Python code as a guide to inform your logic. Refer to relevant steps in the code as part of your reasoning. Do not output or reference the code in any form. Explicitly state the final answer after the final step within `\boxed{ }`.

```
### Query:  
{query}
```

```
### Python Code:  
{code}
```

1056

### Prompt 3: Code Generation

Please write a Python code to solve the following problem. Just give me the code, no explanation, no comments, no input statements. The code should be runnable and print the answer in the end.

```
### Query:  
{query}
```

```
### Python Code:
```

1057

### Prompt 4: Code Execution

Please help me run the following Python code and return its output result instead of the code itself:

```
{code}
```

1058