

# Beyond Templates: Dynamic Adaptation of Reasoning Demonstrations via Feasibility-Aware Exploration

Anonymous ACL submission

## Abstract

Large language models (LLMs) have shown remarkable reasoning capabilities, yet aligning such abilities to small language models (SLMs) remains a challenge due to distributional mismatches and limited model capacity. Existing reasoning datasets, typically designed for powerful LLMs, often lead to degraded performance when directly applied to weaker models. In this work, we introduce Dynamic Adaptation of Reasoning Trajectories (DART), a novel data adaptation framework that bridges the capability gap between expert reasoning trajectories and diverse SLMs. Instead of uniformly imitating expert steps, DART employs a *selective imitation strategy* guided by step-wise adaptability estimation via solution simulation. When expert steps surpass the student’s capacity—signaled by an *imitation gap*—the student autonomously explores alternative reasoning paths, constrained by outcome consistency. We validate DART across multiple reasoning benchmarks and model scales, demonstrating that it significantly improves generalization and data efficiency over static fine-tuning. Our method enhances supervision quality by aligning training signals with the student’s reasoning capabilities, offering a scalable solution for reasoning alignment in resource-constrained models.

## 1 Introduction

Large language models (LLMs) have recently achieved remarkable performance in complex reasoning tasks such as mathematics and programming (OpenAI, 2024; Shao et al., 2024). A key insight from recent work (Zhou et al., 2024; Yue et al., 2024; Ye et al., 2025) is that small, high-quality instruction datasets are surprisingly effective at eliciting sophisticated reasoning abilities in large models. This discovery challenges traditional beliefs (Li et al., 2024; Yu et al., 2024) that complex cognitive skills necessarily require massive

supervised fine-tuning, opening promising avenues for data-efficient model alignment.

Despite the remarkable effectiveness of small, high-quality instruction datasets in eliciting sophisticated reasoning, mainstream approaches (Zhou et al., 2024; Ye et al., 2025; Muennighoff et al., 2025) remain reliant on **static, pre-collected** reasoning datasets. While effective in controlled environments, these datasets struggle to generalize across heterogeneous pretraining distributions, particularly for small language models (SLMs) with diverse training data and limited reasoning capabilities (Xu et al., 2024; Yeo et al., 2025). Disparities in model scale, reasoning proficiency, and training history exacerbate distributional mismatches, significantly hindering the activation of reasoning skills.

To address these challenges, we introduce **Dynamic Adaptation of Reasoning Trajectories (DART)**, a novel data adaptation framework designed to bridge the distribution gap between static reasoning datasets and diverse SLMs. Instead of enforcing uniform imitation of expert demonstrations, DART introduces a *selective imitation strategy* guided by *imitation feasibility estimate*. For each step provided by the expert, DART dynamically assesses the likelihood that the student model can successfully complete the reasoning process conditioned on adopting that step. When imitation is deemed infeasible, the student autonomously explores alternative trajectories while maintaining the consistency of the outcome with the objective of the original task. This approach enables DART to flexibly adapt high-quality reasoning datasets to heterogeneous model populations, significantly improving reasoning elicitation under distribution shift.

In summary, our contributions are as follows.

- We identify the critical limitations of applying static curated reasoning datasets to diverse

083 small language models and propose **DART**, a  
084 novel framework for adapted reasoning data  
085 guided by imitation feasibility.

- 086 • We introduce a Monte Carlo simulation-based  
087 method to estimate the feasibility of imita-  
088 tion per step, allowing selective supervision  
089 tailored to the student model capabilities.
- 090 • We develop an autonomous exploration  
091 mechanism that allows models to recover  
092 from infeasible supervision points, generat-  
093 ing outcome-consistent alternative reasoning  
094 paths.
- 095 • Through extensive experiments across differ-  
096 ent model scales and benchmarks, we demon-  
097 strate that DART substantially improves reason-  
098 ing performance over static fine-tuning,  
099 achieving superior data efficiency and gener-  
100 alization.

## 101 2 Preliminaries and Limitations of 102 Supervised Imitation on Expert 103 Trajectories

### 104 2.1 Problem Definition: Reasoning Capability 105 Elicitation via Minimal Demonstrations

106 We define the reasoning elicitation problem in the  
107 context of large language models (LLMs) with lat-  
108 ent pre-trained knowledge. Let  $\mathcal{Q}$  denote the space  
109 of reasoning problems,  $\mathcal{A}$  the space of answers, and  
110  $\mathcal{R}$  the space of reasoning chains, where each  $r \in \mathcal{R}$   
111 is a sequence of logical steps  $r = \{s_1, s_2, \dots, s_n\}$ .

112 The goal is to learn a reasoning function:

$$113 f : \mathcal{Q} \rightarrow \mathcal{R} \times \mathcal{A} \quad (1)$$

114 so that, given a question  $q \in \mathcal{Q}$ , the model gener-  
115 ates a logically valid reasoning chain  $r \in \mathcal{R}$  and  
116 a verifiable final answer  $a \in \mathcal{A}$ .

117 Prior work (e.g., (Ye et al., 2025), (Muennighoff  
118 et al., 2025)) suggests that reasoning competence  
119 in large language models (LLMs) can be elicited  
120 not by scale alone, but a small set of carefully  
121 crafted demonstrations that expose the underlying  
122 cognitive structure of reasoning. This paradigm as-  
123 sumes that latent reasoning skills embedded within  
124 pretrained models can be activated through appro-  
125 priately designed prompts in the form of explicit  
126 multi-step exemplars.

127 Let  $\mathcal{D} = \{(q_i, r_i, a_i)\}_{i=1}^N$  represent a compact  
128 yet high-quality dataset ( $N \ll |\mathcal{Q}|$ ), where each

tuple contains a question  $q_i$ , a structured reason-  
ing chain  $r_i$ , and its corresponding answer  $a_i$ .  
Each  $r_i$  serves as a **cognitive template**—an in-  
terpretable, step-wise reasoning demonstration de-  
signed to guide the model through logical steps  
with intermediate verification. Instead of introduc-  
ing new knowledge, these templates activate the  
model’s latent reasoning capabilities by leveraging  
structured prompting (Wei et al., 2022; Zhou et al.,  
2024; Ye et al., 2025).

### 139 2.2 Limitations of Supervised Imitation on 140 Expert Demonstrations

141 Despite its pedagogical appeal, supervised imita-  
142 tion over expert demonstrations exhibits critical  
143 limitations when applied to LLMs with diverse ca-  
144 pacity levels.

145 This paradigm (Wei et al., 2022; Ye et al., 2025)  
146 assumes that the model possesses sufficient latent  
147 competence to internalize and reproduce the reason-  
148 ing trajectory in each template. In practice, this  
149 assumption frequently fails. A template  $r_i$  may (i)  
150 over-challenge the model by invoking reasoning  
151 procedures not encoded in its weights, or (ii) mis-  
152 align with the model’s inductive biases, causing  
153 representational mismatch. We define a reasoning  
154 failure event  $\mathcal{F}$  as the inability of the model to em-  
155 ulate the intended behavior given an input-template  
156 pair:

$$157 \mathcal{F}(f; q, r, a) = \mathbb{I}[f(q) \not\approx (r, a)] \quad (2)$$

158 where  $\mathbb{I}[\cdot]$  is the indicator function. Such failures  
159 may arise from superficial imitation, incomplete  
160 reasoning chains, or insufficient justification for the  
161 final answer.

162 Compounding this challenge is the substantial  
163 cost associated with constructing template datasets  
164  $\mathcal{D}$  that satisfy the Cognitive Template Demonstra-  
165 tion criterion. Such templates demand meticulous  
166 logical decomposition, intermediate verification,  
167 and fine-grained pedagogical design. Furthermore,  
168 a template crafted for a specific model often fails to  
169 generalize to others due to differences in scale, pre-  
170 training corpus, or architectural inductive biases,  
171 resulting in distributional shifts. As highlighted  
172 in prior work on imitation learning (Pomerleau,  
173 1991; Ross et al., 2011), relying on static datasets  
174 for training can lead to a distribution mismatch  
175 between the output sequences encountered during  
176 training and those generated auto-regressively by  
177 the student at inference time, undermining general-  
178 ization and robustness.

**The Need for Imitation Feasibility-Aware Adaptation.** These limitations highlight the inadequacy of static demonstrations in addressing the diversity of model behaviors. We argue for a dynamic grounding mechanism that aligns template presentation with the target model’s internal capacity and abstraction level. Rather than treating  $\mathcal{D}$  as fixed input, the elicitation process should adaptively align the demonstrated reasoning path with the model’s own preferred or accessible inference trajectories, potentially reformulating how the reasoning unfolds to match internal representations. This motivates our central question:

*Can we design a dynamic adaptation mechanism that reliably anchors cognitive templates in model-specific latent space, enabling scalable and robust reasoning?*

In the following section, we instantiate this motivation via our proposed framework — **Dynamic Adaptation of Reasoning Trajectories (DART)**.

### 3 Methodology

In this section, we propose **Dynamic Adaptation of Reasoning Trajectories (DART)**, a capability-aware adaptation framework designed to align expert-level reasoning data with the capacity of small language models (SLMs). Instead of statically mimicking expert trajectories from the elicitation template set, DART introduces a selective imitation mechanism that dynamically adapts supervision signals based on the model’s reasoning proficiency. The framework comprises three key components: (1) step-wise adaptability estimation via solution simulation (Section 3.1), (2) imitation gap detection and adaptive path exploration (Section 3.2), and (3) learning from outcome-aligned adapted trajectories (Section 3.3). Figure 1 provides an overview of the pipeline.

#### 3.1 Step-wise Adaptability Estimation via Solution Simulation

To determine whether a given expert step is suitable for imitation, we introduce the concept of **adaptability**: the likelihood that a student model can reach the correct answer when conditioned on that step. This evaluation is conducted via solution simulation—akin to Monte Carlo Tree Search (Kocsis and Szepesvári, 2006; Silver et al., 2016; Świechowski et al., 2023)—by rolling out

multiple completions from partially constructed trajectories that incorporate the candidate step.

Let  $s_{<t} = \{s_0, s_1, \dots, s_{t-1}\}$  be the prefix of expert steps, and  $s_t$  the candidate step under evaluation. The adaptability score  $f_t$  is computed as:

$$f_t = Q(s_{<t}, s_t) = \frac{1}{N_{\text{sim}}} \sum_{i=1}^{N_{\text{sim}}} \mathbb{I}(a_i^{\text{final}} = a^*) \quad (3)$$

where  $N_{\text{sim}}$  denotes the total number of *rollouts* performed for each candidate step  $s_t$ , with each rollout simulating a complete reasoning trajectory conditioned on the prefix  $s_{<t}$  and the adoption of step  $s_t$ .

Empirically observed patterns (see Section 5.1) suggest that adaptability tends to rise in the early stages of expert trajectories, but drops sharply beyond a certain point. This non-monotonic behavior motivates our definition of the **imitation gap**, a regime in which continued imitation becomes counterproductive due to the increasing complexity of the remaining expert steps.

#### 3.2 Adaptive Path Exploration

To avoid overfitting to brittle expert demonstrations, we monitor the *adaptability score* throughout the trajectory and halt imitation once a significant drop is detected (see Equation 3). Motivated by the need to overcome low-adaptability segments that may hinder generalization, DART transitions to autonomous rollout beyond the gap, generating a continuation from the last high-adaptability prefix:

$$\tau_{\text{adapt}} = (s_0, s_1, \dots, s_{t-1}, s'_t, s'_{t+1}, \dots, s'_T), \quad (4)$$

where  $s'_t, \dots, s'_T$  are student-generated reasoning steps. Inspired by outcome-based learning strategies (DeepSeek-AI et al., 2025), we do not constrain this trajectory to mimic the expert’s form. Instead, we enforce an *outcome consistency* constraint to ensure semantic alignment, as described in Eq. equation 5, as we observe that process supervision (Lightman et al., 2024; Zhang et al., 2025), such as via a Process Reward Model (PRM), often encounters inherent ambiguities and standardization challenges in practice.

$$C(\tau_{\text{adapt}}, \tau_{\text{expert}}) = \begin{cases} 1, & \text{if } \mathcal{O}(\tau_{\text{adapt}}) = \mathcal{O}(\tau_{\text{expert}}), \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

Here,  $C \in \{0, 1\}$  denotes task-level agreement, with  $\mathcal{O}(\cdot)$  representing the final answer obtained by

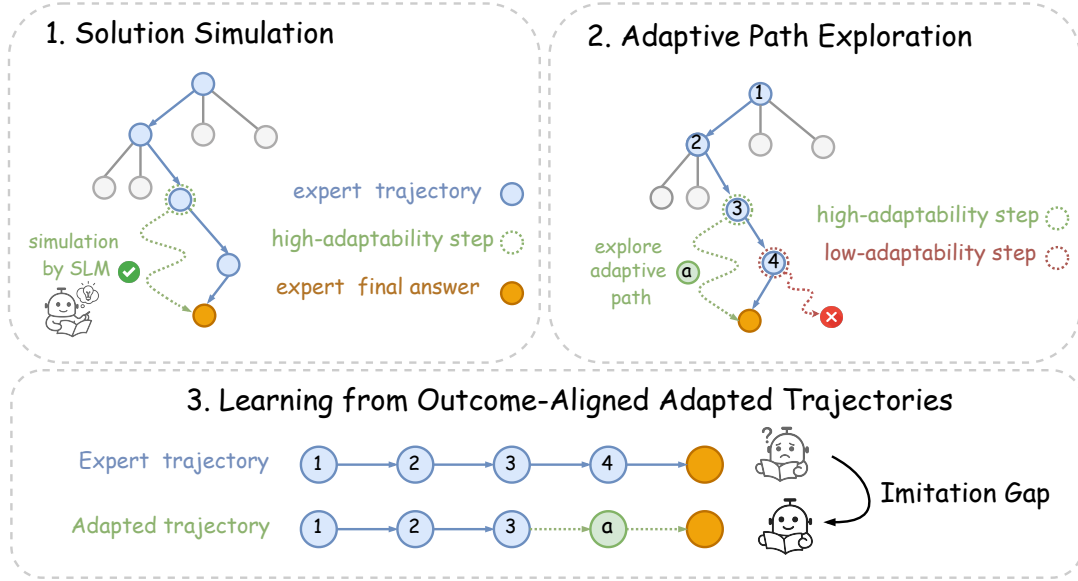


Figure 1: Overview of the DART framework.

executing a reasoning path. Specifically,  $\mathcal{O}(\tau_{\text{expert}})$  refers to the outcome of the expert demonstration, while  $\mathcal{O}(\tau_{\text{adapt}})$  captures the result of the student’s adapted trajectory. The constraint  $\mathcal{O}(\tau_{\text{adapt}}) = \mathcal{O}(\tau_{\text{expert}})$  ensures that, although the reasoning paths may differ, their semantic outcomes are equivalent. This outcome consistency criterion allows the student to depart from brittle expert traces while preserving task correctness.

This strategy empowers the student model to develop its own reasoning strategies beyond segments with low adaptability, guided solely by the correctness of the final outcome. By anchoring supervision at the outcome level rather than mimicking intermediate steps, we alleviate the brittleness of process-level imitation. This encourages robust generalization, reduces reliance on ambiguous or inconsistent expert demonstrations, and aligns with the broader goal of enabling flexible yet goal-directed reasoning.

### 3.3 Learning from Outcome-Aligned Adapted Trajectories

To effectively activate the student model’s own reasoning ability, we apply a standard cross-entropy loss on the outcome-aligned adapted trajectories generated during autonomous exploration. This training objective encourages the model to reinforce reasoning patterns that are not only aligned with the task goal but also feasible under its own capacity.

Training proceeds by distilling the adapted

trajectory  $\tau_{\text{adapt}}$  using a standard cross-entropy loss (Kim and Rush, 2016; Bengio et al., 2003):

$$L_{\text{DART}} = - \sum_{t=1}^T \mathbb{E}_{(s_{<t}, a_t) \sim \tau_{\text{adapt}}} [\log \pi_{\text{student}}(a_t | s_{<t})] \quad (6)$$

Here,  $s_{<t} = \{s_0, \dots, s_{t-1}\}$  denotes the contextual prefix consisting of all prior reasoning steps up to time  $t$ , and  $a_t$  is the corresponding next-step decision. This loss encourages the student model  $\pi_{\text{student}}$  to maximize the likelihood of producing  $a_t$  when conditioned on its own reasoning history.

By learning from outcome-aligned yet model-compatible trajectories, DART provides high-quality supervision that reflects the student’s actual competence. This approach decouples the training signal from rigid trajectory matching, improving both robustness and scalability across models with varying capacity.

## 4 Experiments

We evaluate DART across a series of mathematical reasoning benchmarks to assess its effectiveness in adapting expert data to student models of varying capacities.

### 4.1 Experimental Setup

**Adaptation Datasets.** We conduct adaptation experiments using two datasets. (1) **LIMO** dataset (Ye et al., 2025), a curated set of 817 high-quality math reasoning examples with multi-step CoT demonstrations tailored. We use the official

Table 1: Main results (%) on LIMO and Math-QwQ-32B across adaptation strategies and model sizes. **Static** overfits to noisy data, while **Adaptation-Full** improves results through exploration and filtering of low-adaptability segments.

Dataset	Method	GSM8K	GaoKao	Olympiad Bench	College Math	MMLU STEM	Avg.
Qwen2.5-0.5B-Instruct							
-	No-Tuning	49.1	30.4	9.3	28.9	36.7	30.9
Math-QwQ-32B	Static	39.8 <b>-9.3</b>	20.5 <b>-9.9</b>	5.9 <b>-3.4</b>	17.3 <b>-11.6</b>	27.9 <b>-8.8</b>	22.3 <b>-8.6</b>
Math-QwQ-32B	<b>Adaptation-Full</b>	49.6 <b>+0.5</b>	30.9 <b>+0.5</b>	9.3 <b>+0.0</b>	27.5 <b>-1.4</b>	37.5 <b>+0.8</b>	31.0 <b>+0.1</b>
LIMO	Static	49.6 <b>+0.5</b>	26.8 <b>-3.6</b>	7.7 <b>-1.6</b>	27.3 <b>-1.6</b>	32.9 <b>-3.8</b>	28.9 <b>-2.0</b>
LIMO	<b>Adaptation-Full</b>	52.2 <b>+3.1</b>	32.5 <b>+2.1</b>	9.8 <b>+0.5</b>	29.1 <b>+0.2</b>	37.2 <b>+0.5</b>	32.2 <b>+1.3</b>
Qwen2.5-3B-Instruct							
-	No-Tuning	87.0	56.6	27.3	39.9	47.6	51.7
Math-QwQ-32B	Static	82.0 <b>-5.0</b>	46.2 <b>-10.4</b>	20.1 <b>-7.2</b>	35.7 <b>-4.2</b>	51.3 <b>+3.7</b>	47.1 <b>-4.6</b>
Math-QwQ-32B	<b>Adaptation-Full</b>	86.6 <b>-0.4</b>	57.9 <b>+1.3</b>	29.0 <b>+1.7</b>	44.4 <b>+4.5</b>	50.6 <b>+3.0</b>	53.7 <b>+2.0</b>
LIMO	Static	85.4 <b>-1.6</b>	53.8 <b>-2.8</b>	25.2 <b>-2.1</b>	41.8 <b>+1.9</b>	54.8 <b>+7.2</b>	52.2 <b>+0.5</b>
LIMO	<b>Adaptation-Full</b>	87.2 <b>+0.2</b>	59.5 <b>+2.9</b>	30.4 <b>+3.1</b>	43.9 <b>+4.0</b>	62.7 <b>+15.1</b>	56.7 <b>+5.0</b>
LLaMA-3B-Instruct							
-	No-Tuning	38.4	21.3	10.7	16.0	48.0	26.9
Math-QwQ-32B	Static	67.7 <b>+29.3</b>	30.4 <b>+9.1</b>	10.2 <b>-0.5</b>	21.1 <b>+5.1</b>	39.2 <b>-8.8</b>	33.7 <b>+6.8</b>
Math-QwQ-32B	<b>Adaptation-Full</b>	72.3 <b>+33.9</b>	34.0 <b>+12.7</b>	10.2 <b>-0.5</b>	21.8 <b>+5.8</b>	38.3 <b>-9.7</b>	35.3 <b>+8.4</b>
LIMO	Static	28.1 <b>-10.3</b>	15.1 <b>-6.2</b>	3.6 <b>-7.1</b>	11.4 <b>-4.6</b>	48.6 <b>+0.6</b>	21.4 <b>-5.5</b>
LIMO	<b>Adaptation-Full</b>	47.2 <b>+8.8</b>	24.9 <b>+3.6</b>	6.4 <b>-4.3</b>	18.1 <b>+2.1</b>	47.2 <b>-0.8</b>	28.8 <b>+1.9</b>

328 filtered release<sup>1</sup>. (2) The **Math-QwQ-32B** dataset  
 329 is a synthetic dataset derived from the MATH  
 330 benchmark (Hendrycks et al., 2021b), where the  
 331 Qwen/QwQ-32B-Preview model<sup>2</sup> generates long-  
 332 form Chain-of-Thought (CoT) solutions for 5,383  
 333 problems in the training subset.

334 **Adaptation Strategies.** We evaluate three adap-  
 335 tation strategies to disentangle the effects of selec-  
 336 tive imitation and adaptive exploration. (1) *No-*  
 337 *Tuning* denotes direct zero-shot evaluation. (2)  
 338 *Static* reflects standard offline supervised fine-  
 339 tuning on the full set of expert trajectories, without  
 340 any filtering or adaptability mechanism. (3) The  
 341 *Adaptation-Full* strategy represents the complete  
 342 DART pipeline, integrating imitation gap detection  
 343 with outcome-consistent student exploration. This  
 344 approach empowers the model to autonomously  
 345 explore alternative reasoning paths when expert im-  
 346 itation becomes unreliable. If the model can’t find  
 347 a suitable alternative path, it discards that expert  
 348 example. We evaluate our method on Qwen2.5-

Instruct models at 0.5B, 1.5B, 3B and LLaMA-3B-  
 Instruct scales, covering a diverse range of SLMs.

351 **Benchmark Tasks.** We evaluate DART on seven  
 352 diverse benchmarks encompassing a broad spec-  
 353 trum of mathematical reasoning. These include  
 354 GSM8K (Cobbe et al., 2021), covering grade-  
 355 school to competition-level problems, To assess lin-  
 356 guistic and cultural generalization, we incorporate  
 357 GaoKao 2023 En (Liao et al., 2024), a Chinese na-  
 358 tional exam benchmark. OlympiadBench (He et al.,  
 359 2024) features high-difficulty, compositional prob-  
 360 lems from international math competitions. Col-  
 361 lege Math (Tang et al., 2024) probes undergraduate-  
 362 level topics in calculus, algebra, and discrete math.  
 363 MMLU-STEM (Hendrycks et al., 2021a) evaluates  
 364 STEM-focused reasoning breadth. Overall adap-  
 365 tation is quantified by the arithmetic mean (Avg.)  
 366 across all benchmarks.

367 **Training and Model Selection** To ensure the re-  
 368 liability of experimental results, we conducted sys-  
 369 tematic training and model selection for all models.  
 370 For both the *Static* and *Adaptation-Full* strategies,  
 371 we trained models at the 0.5B, 1.5B, and 3B scales

<sup>1</sup><https://huggingface.co/GAIR/LIMO>

<sup>2</sup><https://huggingface.co/Qwen/QwQ-32B-Preview>

for 15 epochs, saving a model checkpoint at the end of each epoch, resulting in 15 checkpoints per model. These checkpoints were evaluated on the validation sets, and the model with the best performance was selected as the final model. The training parameter settings were consistent with LIMO (Ye et al., 2025). All experiments are conducted on the same NVIDIA A100 GPU infrastructure. Additional implementation details, including configurations and setups, are provided in Appendix A.

## 4.2 Main Results

Table 1 reports performance across five mathematical reasoning benchmarks, demonstrating the effectiveness of the proposed DART framework in aligning expert reasoning with the capabilities of small language models (SLMs).

**Static Results** Static, which rigidly imitates expert trajectories without adaptation, exhibits clear limitations. On Qwen2.5-0.5B, it decreases accuracy by **8.6 points** on *Math-QwQ-32B* and by **2.0 points** on *LIMO*. On Qwen2.5-3B, Static reduces accuracy on *Math-QwQ-32B* by **4.6 points**. A similar degradation is observed on LLaMA-3B, where accuracy on *LIMO* drops by **5.5 points**. These results are obtained under our careful training and model selection protocol (see Sec. 4.1), ensuring that the observed degradation is not caused by insufficient training but rather reflects the inherent limitations of the Static strategy. These findings indicate that imitating expert demonstrations without adaptation not only constrains small models but can also undermine the performance of larger ones.

**Adaptation-Full Results** Adaptation-Full shows robust improvements over No-Tuning, consistently enhancing performance across datasets and model scales. For example, on the Qwen2.5-0.5B model, Adaptation-Full improves the average accuracy on *LIMO* by +1.3 points, while on the larger Qwen2.5-3B, it yields gains of +2.0 and +5.0 points on *Math-QwQ-32B* and *LIMO*, respectively. The effect is even more pronounced on LLaMA-3B, where Adaptation-Full boosts *Math-QwQ-32B* by +8.4 points and *LIMO* by +1.9 points. On average, Adaptation-Full achieves a **+4.9 point** improvement over No-Tuning, demonstrating its effectiveness in aligning reasoning trajectories with model capacity while maintaining stability across different architectures.

## 5 Analysis

We further analyze the internal mechanisms of DART, aiming to understand why selective imitation and autonomous exploration improve reasoning capabilities.

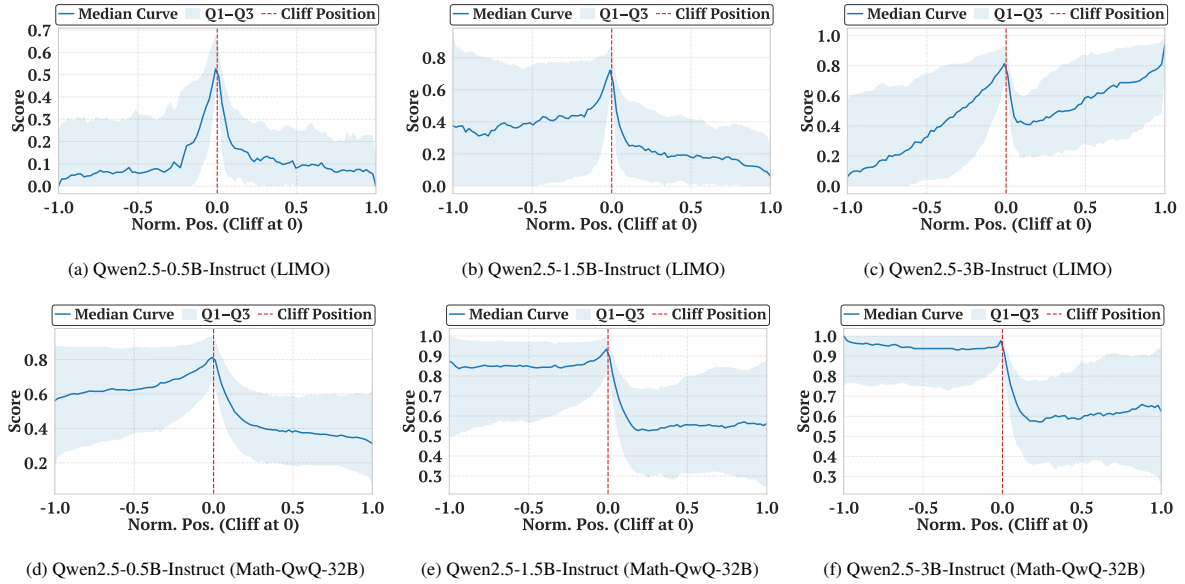
### 5.1 Step-wise Adaptability Reveals the Emergence of the Imitation Gap

To empirically validate the *imitation gap* hypothesis introduced in Section 3, we estimate the step-wise adaptability scores of Qwen2.5-Instruct student models across three parameter scales (0.5B, 1.5B, 3B) on two reasoning datasets (LIMO and Math-QwQ-32B). Each adaptability score quantifies the model’s average probability of reaching the correct final answer when conditioned on imitating a given intermediate step from the expert trajectory. To remove any confounding effect of unequal trajectory lengths, we length-normalize every trace and register the detected cliff at  $x = 0$ . The curve shows the median adaptability score, and the shaded area the interquartile range (Q1–Q3).

As shown in Figure 2, these curves reveal a consistent behavioral pattern: early in the reasoning path, student models exhibit increasing adaptability as they benefit from following expert steps. However, beyond a certain point, adaptability scores sharply decline—signaling that the student has encountered steps that exceed its reasoning capacity, leading to degraded rollout completions and a collapse in trajectory success.

This non-monotonic pattern reveals the **imitation gap**—a critical region where student models falter due to misalignment between their capabilities and the expert’s step distribution. This misalignment arises from distributional discrepancies, where expert trajectories include reasoning patterns outside the student’s abstraction space. Consequently, continued imitation in this zone not only fails to benefit learning but actively impairs performance. This phenomenon underscores our central argument: effective reasoning supervision must be dynamically aligned with model-specific capabilities, as realized in our DART framework. To elucidate the imitation gap’s impact on adaptability score decline and reasoning performance, we present a case study on the LIMO dataset, pinpointing its onset in a complex reasoning task through Qwen2.5-3B-Instruct’s step-wise adaptability scores (see Table 6).

Figure 2: Step-wise adaptability scores across expert trajectories for Qwen2.5-Instruct student models of varying sizes (0.5B, 1.5B, 3B parameters) under LIMO (top row) and Math-QwQ-32B dataset (bottom row) supervision. The emergence of the **Imitation Gap** is evident: initial steps yield positive adaptation, but continued step-by-step imitation can become harmful.



## 5.2 Impact of Search Restriction on Adaptation Strategies

To evaluate the impact of adaptation strategies without autonomous search, we designed two variants: **Adaptation-First** and **Adaptation-Gap**. **Adaptation-First** halts imitation once a feasible solution state is detected, whereas **Adaptation-Gap** monitors adaptability scores and terminates imitation when sharp declines occur, as described in Section 5.1. Table 2 presents the evaluation results on Math-QwQ-32B for 1.5B and 3B Qwen2.5-Instruct models. Both strategies exhibit performance degradation compared to **Adaptation-Full**, highlighting the critical role of autonomous search for recovery in complex reasoning paths. Notably, **Adaptation-Gap** consistently outperforms **Adaptation-First** across all benchmarks, with significant accuracy gains in the average performance (35.8% vs. 30.3% for 1.5B and 43.9% vs. 31.0% for 3B). This improvement stems from its capacity-aware truncation, which effectively filters out low-adaptability segments, preventing error propagation and enhancing stability.

## 5.3 Capacity-Aligned Lexical Dynamics Under Adaptation

To investigate how DART reshapes student model behavior at different scales, we analyze keyword frequency changes between static and adapted dataset. Table 3 lists the top 20 tokens with the largest shifts in the first sentence of each reasoning

Table 2: Accuracy (%) on the Math-QwQ-32B dataset for Qwen2.5-1.5B-Instruct and Qwen2.5-3B-Instruct models under different adaptation strategies. **Adaptation-First** performs early stopping at feasible solution states, while **Adaptation-Gap** selectively truncates imitation paths based on adaptability declines. **Adaptation-Full** integrates autonomous search, achieving the highest performance across benchmarks. **Bold** values indicate the best results in each group.

Method	GSM8K	GaoKao	Olympiad Bench	College Math	MMLU STEM	Avg.
Qwen2.5-1.5B-Instruct						
Adaptation-First	41.2	33.0	10.7	24.9	41.7	30.3
Adaptation-Gap	60.0	34.5	13.6	29.7	41.3	35.8
<b>Adaptation-Full</b>	<b>74.2</b>	<b>48.6</b>	<b>19.6</b>	<b>39.4</b>	<b>57.7</b>	<b>47.9</b>
Qwen2.5-3B-Instruct						
Adaptation-First	36.7	32.7	12.9	26.1	46.8	31.0
Adaptation-Gap	77.9	43.4	18.4	36.0	44.0	43.9
<b>Adaptation-Full</b>	<b>86.6</b>	<b>57.9</b>	<b>29.0</b>	<b>44.4</b>	<b>50.6</b>	<b>53.7</b>

step for the Qwen2.5-Instruct series 0.5B, 1.5B, and 3B models.

Adaptation reduces exploratory terms like *but*, *wait*, and *alternatively*, while amplifying goal-oriented expressions such as *step*, *solve*, *find*, and *need*. In the 1.5B model, *but* and *wait* drop by 0.36% and 0.20% percentage points, while *find* and *need* rise by 0.13% and 0.14% points. This shift reflects a transition from hesitant exploration to decisive, solution-driven reasoning. These changes reduce uncertainty and digression—traits often seen in expert trajectories but burdensome for smaller

Table 3: Top 20 Keyword Frequency Changes Across Model Sizes

Keyword	0.5B (%)			1.5B (%)			3B (%)		
	Static	Adapted	$\Delta$	Static	Adapted	$\Delta$	Static	Adapted	$\Delta$
but	2.73	2.59	-0.14	2.73	2.37	-0.36	2.73	2.27	-0.46
alternatively	0.86	0.79	-0.07	0.86	0.72	-0.14	0.86	0.71	-0.15
wait	2.30	2.23	-0.07	2.30	2.10	-0.20	2.30	2.00	-0.30
therefore	1.55	1.50	-0.05	1.55	1.43	-0.13	1.55	1.40	-0.15
check	0.51	0.47	-0.04	0.51	0.40	-0.11	0.51	0.33	-0.18
another	0.29	0.26	-0.03	0.29	0.20	-0.09	0.29	0.17	-0.12
then	0.97	0.94	-0.03	-	-	-	-	-	-
pi	0.11	0.09	-0.02	-	-	-	-	-	-
perhaps	0.55	0.53	-0.02	-	-	-	-	-	-
length	0.22	0.24	+0.02	-	-	-	-	-	-
step	0.27	0.30	+0.02	0.27	0.37	+0.10	0.27	0.41	+0.14
now	0.43	0.46	+0.03	0.43	0.49	+0.06	0.43	0.50	+0.07
first	0.89	0.92	+0.03	0.89	0.96	+0.07	0.89	0.96	+0.07
since	0.80	0.83	+0.03	-	-	-	-	-	-
have	0.64	0.67	+0.03	0.64	0.69	+0.05	-	-	-
let	1.83	1.86	+0.04	1.83	1.88	+0.05	-	-	-
need	0.54	0.58	+0.04	0.54	0.68	+0.14	0.54	0.69	+0.16
find	0.37	0.42	+0.04	0.37	0.50	+0.13	0.37	0.54	+0.16
newline	-	-	-	0.00	0.05	+0.05	-	-	-
equation	-	-	-	0.76	0.81	+0.05	0.76	0.85	+0.09

models. In static supervision, such expressions appear frequently, straining low-capacity models and widening the *Imitation Gap* (Sec. 3.1), where expert strategies exceed model capabilities.

DART bridges this gap by replacing brittle reasoning paths with model-originated decision traces. This adaptation maintains task objectives while restructuring execution to fit model capacity, leading to stable and efficient reasoning.

## 6 Related Work

**Chain-of-Thought Reasoning** Early work on chain-of-thought reasoning (CoT) (Wei et al., 2022) primarily focused on *short CoT*, where models generate concise reasoning paths to solve problems. Recent advances (Chen et al., 2025) have shifted towards *long CoT prompting*, encouraging more elaborate reasoning chains that enable systematic exploration of multiple paths (*branching*) and backtracking when errors are detected. While techniques like knowledge distillation (Hinton et al., 2015; Luo et al., 2025) and reinforcement learning (Hou et al., 2025) have been used to equip large language models (LLMs) with long CoT capabilities, these efforts remain largely confined to models with substantial parameter sizes. In contrast, our work specifically addresses the unique challenges associated with training smaller-scale models for complex reasoning tasks.

**Data-Efficient Reasoning Elicitation** A related line of work investigates how minimal supervision can elicit latent reasoning abilities in pretrained

models (Ye et al., 2025; Muennighoff et al., 2025). These methods rely on a few carefully designed *cognitive templates*, to guide reasoning, but often assume that models possess the necessary prior knowledge. This assumption makes the templates brittle when cognitive demands exceed model capacity. To address this limitation, we propose a feasibility-aware adaptation framework that dynamically adjusts supervision to model ability, enabling robust reasoning across diverse capacity profiles.

## 7 Conclusion

We propose Dynamic Adaptation of Reasoning Trajectories (DART), a data adaptation framework designed to improve reasoning elicitation for small language models. By introducing adaptability-based selective imitation and outcome-consistent exploration, our method aims to better align expert demonstrations with model capabilities. Experimental results across several benchmarks show that DART can improve reasoning performance compared to static fine-tuning. We hope this work provides a step toward more flexible and model-aware data alignment strategies for reasoning tasks.

## 8 Limitations

Our framework is effective for structured reasoning tasks with verifiable outcomes. However, its extension to open-ended tasks with inherent output uncertainty remains limited, suggesting the need for refined supervision mechanisms and evaluation metrics to ensure outcome consistency.

## References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. A neural probabilistic c model. *J. Mach. Learn. Res.*, 3(null):1137–1155.
- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. [Towards reasoning era: A survey of long chain-of-thought for reasoning large language models](#). *Preprint*, arXiv:2503.09567.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168. CoRR, abs/2110.14168.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, and 1 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, Bangkok, Thailand. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. [Measuring mathematical problem solving with the math dataset](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#). *Preprint*, arXiv:1503.02531.
- Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu, Kaizhi Qian, Jacob Andreas, and Shiyu Chang. 2025. [Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning](#). *arXiv preprint arXiv:2504.01296*.
- Yoon Kim and Alexander M. Rush. 2016. [Sequence-level knowledge distillation](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas. Association for Computational Linguistics.
- Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer.
- Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, and 1 others. 2024. [Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions](#). *Hugging Face repository*.
- Wendi Li and Yixuan Li. [Process reward model with q-value rankings](#). In *The Thirteenth International Conference on Learning Representations*.
- Minpeng Liao, Chengxi Li, Wei Luo, Jing Wu, and Kai Fan. 2024. [MARIO: math reasoning with code interpreter output - a reproducible pipeline](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 905–924. Association for Computational Linguistics.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. [Let’s verify step by step](#). In *The Twelfth International Conference on Learning Representations*.
- Yijia Luo, Yulin Song, Xingyao Zhang, Jiaheng Liu, Weixun Wang, GengRu Chen, Wenbo Su, and Bo Zheng. 2025. [Deconstructing long chain-of-thought: A structured reasoning optimization framework for long cot distillation](#). *Preprint*, arXiv:2503.16385.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. [s1: Simple test-time scaling](#). *Preprint*, arXiv:2501.19393.
- OpenAI. 2024. [Openai o1 system card](#). *Preprint*, arXiv:2412.16720.
- Dean A. Pomerleau. 1991. [Efficient training of artificial neural networks for autonomous navigation](#). *Neural Computation*, 3(1):88–97.
- Stephane Ross, Geoffrey Gordon, and Drew Bagnell. 2011. [A reduction of imitation learning and structured prediction to no-regret online learning](#). In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 627–635, Fort Lauderdale, FL, USA. PMLR.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *Preprint*, arXiv:2402.03300.

683 David Silver, Aja Huang, Chris J Maddison, Arthur  
684 Guez, Laurent Sifre, George Van Den Driessche, Ju-  
685 lian Schrittwieser, Ioannis Antonoglou, Veda Pan-  
686 neershelvam, Marc Lanctot, and 1 others. 2016. Mas-  
687 tering the game of go with deep neural networks and  
688 tree search. *nature*, 529(7587):484–489.

689 Maciej Świechowski, Konrad Godlewski, Bartosz Saw-  
690 icki, and Jacek Mańdziuk. 2023. Monte carlo tree  
691 search: A review of recent modifications and appli-  
692 cations. *Artificial Intelligence Review*, 56(3):2497–  
693 2562.

694 Zhengyang Tang, Xingxing Zhang, Benyou Wang, and  
695 Furu Wei. 2024. [Mathscale: Scaling instruction tun-  
696 ing for mathematical reasoning](#). In *Proceedings of  
697 the 41st International Conference on Machine Learn-  
698 ing (ICML)*.

699 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
700 Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le,  
701 and Denny Zhou. 2022. [Chain-of-thought prompt-  
702 ing elicits reasoning in large language models](#). In  
703 *Advances in Neural Information Processing Systems*,  
704 volume 35, pages 24824–24837. Curran Associates,  
705 Inc.

706 Zhangchen Xu, Fengqing Jiang, Luyao Niu, Bill Yuchen  
707 Lin, and Radha Poovendran. 2024. [Stronger mod-  
708 els are not stronger teachers for instruction tuning](#).  
709 *Preprint*, arXiv:2411.07133.

710 Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie  
711 Xia, and Pengfei Liu. 2025. [Limo: Less is more for  
712 reasoning](#). *Preprint*, arXiv:2502.03387.

713 Edward Yeo, Yuxuan Tong, Morry Niu, Graham  
714 Neubig, and Xiang Yue. 2025. [Demystifying  
715 long chain-of-thought reasoning in llms](#). *Preprint*,  
716 arXiv:2502.03373.

717 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu,  
718 Zhengying Liu, Yu Zhang, James T. Kwok, Zhen-  
719 guo Li, Adrian Weller, and Weiyang Liu. 2024.  
720 [Metamath: Bootstrap your own mathematical  
721 questions for large language models](#). *Preprint*,  
722 arXiv:2309.12284.

723 Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen.  
724 2024. [Mammoth2: Scaling instructions from the  
725 web](#). *Preprint*, arXiv:2405.03548.

726 Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen  
727 Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jin-  
728 gren Zhou, and Junyang Lin. 2025. [The lessons of  
729 developing process reward models in mathematical  
730 reasoning](#). *Preprint*, arXiv:2501.07301.

731 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer,  
732 Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping  
733 Yu, Lili Yu, and 1 others. 2024. [Lima: Less is more  
734 for alignment](#). *Advances in Neural Information Pro-  
735 cessing Systems*, 36.

## Appendix 736

### A Implementation Details 737

738 In this section, we provide a detailed account of  
739 the experimental configurations and setups to en-  
740 sure the transparency and reproducibility of our  
741 research. We introduce the prompt designs used  
742 in our experiments. Furthermore, we elaborate  
743 on the parameter configurations for simulation ex-  
744 periments and adaptive path exploration. These  
745 configurations are designed to balance computa-  
746 tional efficiency and response diversity, ensuring  
747 the stability and adaptability of the model across  
748 various tasks.

#### A.1 Experiment Prompts 749

750 In our simulation experiments, we employed a  
751 structured prompting approach to guide the lan-  
752 guage model through multi-step reasoning tasks.  
753 The primary simulation prompt used in our study  
754 is defined as follows:

#### Simulation Prompt

```

Problem: data["question"]
Existing reasoning path:
data["answer"]
Guidelines for continuing the
reasoning:
1. Understand the existing
path: Carefully analyze the
existing reasoning path and
understand the logic and basis
of each step.
2. Identify the next step:
Based on the last step of the
existing path, determine the
possible directions for the
next step of reasoning.
3. Reason step-by-step:
Start from the last step
of the existing path and
proceed with the reasoning
step-by-step.
4. Final conclusion: When
the reasoning is complete,
put your final answer within
boxed{}.
Continue reasoning step by
step, and put your final
answer within \boxed{}.
```

### Exploration Prompt

```
Problem: {data["question"]}
Existing reasoning (read-only;
cite only key points when
anchoring, do NOT restate the
whole text): {data["answer"]}
[Guidelines (strictly follow)]
1. Role & Boundaries
- Continue only from the
last step of the existing
reasoning; do not restate or
rewrite prior content.
- If new symbols/variables are
needed, first define their
meaning and domain in one
sentence, then use them.
2. Anchoring & Continuation
- Use one line to anchor the
key equation/state of the
"last step" (key points only;
do not restate the full text).
- If you can determine
the next step number from
previous steps, continue
that numbering; if not, do
not number--start reasoning
directly.
3. Explore
- Following your own reasoning
style and anchored to what has
been established, carry the
reasoning forward from here.
Final conclusion: When the
reasoning is complete, put
your final answer within
\boxed{}
```

This prompt encourages the model to decompose the problem into intermediate steps and to clearly indicate the final answer using LaTeX-style boxed notation. This formatting ensures consistency across outputs and facilitates automated evaluation of results.

In addition to standard simulation prompting, we introduce a dedicated exploration prompt tailored for the adaptive trajectory rollout described in Section 3.2. This prompt is activated once a low-adaptability segment is detected and aims to continue reasoning beyond the imitation gap. It conditions the model on the prefix of high-adaptability reasoning steps and allows for autonomous contin-

uation constrained only by outcome correctness:

This exploration prompt encourages the model to develop its own reasoning path from the last trustworthy segment, fostering flexible generalization while maintaining semantic alignment with the expert outcome.

### A.2 Parameter Configuration for Simulation

The simulation procedure in Algorithm 3.1 adopts stochastic decoding to explore alternative reasoning paths beyond expert demonstrations. We sample  $N = 4$  candidate continuations per step, corresponding to the adaptation simulation count  $N_{\text{sim}}$ .

Each trajectory is generated with a maximum length of  $\text{MAX\_NEW\_TOKENS} = 4000$ . To promote determinism while retaining minimal stochasticity, we set the sampling temperature to  $\text{TEMPERATURE} = 0.1$ . Decoding is performed in batches of  $\text{BATCH\_SIZE} = 32$  to enable efficient parallel inference under hardware constraints. These settings ensure stable simulation rollouts with low-variance outputs, suitable for evaluating adaptability under controlled decoding conditions.

Regarding the computational cost of our simulation, we employ SGLang as the inference deployment framework for small-scale models. Given that our approach primarily focuses on adapting reasoning templates to small-scale models, it encounters challenges stemming from distributional mismatches and the limited capacity of small language models (SLMs) compared to large language models (LLMs). Consequently, we are able to deploy our model on a single GPU. To enhance simulation efficiency, we utilize a distributed rollout engine, with multiple SGLang workers managed by an SGLang router to achieve load balancing. Within our code framework, for a 3B model with an estimated four simulations per step, processing the LIMO dataset on a single node equipped with eight A100 GPUs requires approximately six hours.

### A.3 Parameter Configuration for Adaptive Path Exploration

To support the adaptive rollout mechanism described in Section 3.2, we configured the EXPLORE phase with carefully selected hyperparameters to balance computational efficiency and response diversity. The sampling procedure was executed with a candidate beam size of  $\text{NUM\_SAMPLES} = 8$ , meaning that at each decision step, eight reasoning continuations were generated for evaluation based on

821 the adaptability score.

822 We set the maximum generation length to  
823 `MAX_NEW_TOKENS = 2000` to allow sufficient space  
824 for multi-step reasoning without premature trunca-  
825 tion. A temperature of `TEMPERATURE = 0.7` was  
826 employed to introduce moderate randomness in  
827 token sampling, facilitating the exploration of alter-  
828 native reasoning paths while retaining coherence.

829 Batch inference was performed with a  
830 `BATCH_SIZE = 64` to utilize GPU resources effi-  
831 ciently during large-scale rollouts. The underlying  
832 language model was run using half-precision  
833 arithmetic (`DTYPE = float16`), which reduced  
834 memory footprint and improved throughput  
835 without compromising output quality.

836 Additionally, the maximum number of concu-  
837 rent sequences handled by the inference engine  
838 (VLLM) was set to `MAX_NUM_SEQS = 512`, en-  
839 abling high-throughput parallel generation during  
840 exploration. These settings ensured scalable, stable,  
841 and semantically diverse adaptation rollouts  
842 that align with the outcome consistency constraint  
843 described in Equation equation 5.

## 844 B Impact of Search Path Quality on Model 845 Performance

846 To investigate the impact of search quality on  
847 model performance, we conducted a comparative  
848 experiment (see Table 4). After completing the  
849 adaptation path search, we removed paths exhibit-  
850 ing severe repetition phenomena. As illustrated in  
851 Figure 3, the proportion of repeated paths during  
852 exploration decreases progressively with increas-  
853 ing model parameter size, indicating that improve-  
854 ments in the model’s generative capability and con-  
855 textual memory effectively reduce repetition.

856 We refer to the results after removing such re-  
857 peated paths as *Adaptation-Cleaned* and systemati-  
858 cally evaluated these against the complete search  
859 results without removing repeated paths, denoted  
860 as *Adaptation-Raw*. Experimental results demon-  
861 strate that filtering out repeated paths leads to sig-  
862 nificant performance gains, further highlighting the  
863 critical role of search path quality in overall model  
864 performance.

## 865 C Comparative Analysis of Truncation 866 Methods under Search Constraints

867 In our previous section (see Section 5.2), we  
868 investigate two truncation methods under differ-  
869 ent search constraints. Specifically, we designed  
870 two variants: **Adaptation-First** and **Adaptation-**

**Gap**. The **Adaptation-First** method halts imi- 871  
872 tation once a feasible solution state is detected,  
873 whereas **Adaptation-Gap** monitors adaptability  
874 scores and terminates imitation when sharp de-  
875 clines occur, as detailed in Section 5.1.

876 We compare the truncation positions of the two  
877 methods across different datasets and model sizes.  
878 Our analysis indicates that on more challenging  
879 datasets, or when the model capacity is limited  
880 (e.g., results on the 0.5B models for both datasets),  
881 the truncation points identified by **Adaptation-**  
882 **First** and **Adaptation-Gap** are largely consistent.  
883 This can be attributed to the complexity of the rea-  
884 soning cognitive templates in these datasets relative  
885 to the model’s capabilities: once the model identi-  
886 fies a path leading to a feasible solution, continued  
887 imitation often ventures into regions that are diffi-  
888 cult to adapt to, typically accompanied by a sharp  
889 decline in adaptability scores. Consequently, the  
890 truncation positions under both **Adaptation-First**  
891 and **Adaptation-Gap** modes are generally aligned.

892 Conversely, on the Math-Qwen dataset, notable  
893 differences in truncation positions emerge. Many  
894 models, after reaching the step at which the fi-  
895 nal answer can be searched, continue to utilize  
896 subsequent adaptable path segments. Thus, the  
897 **Adaptation-Gap** method is able to detect and lever-  
898 age a greater number of these usable step frag-  
899 ments, resulting in more substantial performance im-  
900 provements, as reported in Table 2.

## 901 D Proof of Existence of Imitation Gap

902 To rigorously establish the existence of the imita-  
903 tion gap in behavioral cloning for reasoning tasks,  
904 we model the process as a deterministic Markov De-  
905 cision Process (MDP)  $M = (\mathcal{S}, \mathcal{A}, \mathcal{T}, r, \rho, T)$  (Li  
906 and Li), where:

- 907 •  $\mathcal{S}$ : state space of reasoning prefixes including  
908 the initial instruction  $x$ ;
- 909 •  $\mathcal{A}$ : action space of reasoning steps  $s_t$ ;
- 910 •  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ : deterministic transition  
911 appending  $s_t$  to the prefix;
- 912 •  $r : \mathcal{S} \rightarrow \mathbb{R}$ : reward function, with  $r(s_T) = 1$   
913 if the trajectory yields the correct answer  $a^*$ ,  
914 and 0 otherwise;
- 915 •  $\rho$ : initial distribution over instructions ( $s_0 \sim$   
916  $\rho$ );
- 917 •  $T$ : maximum trajectory length (horizon).

Table 4: Comparison of accuracy (%) on the Math-QwQ-32B dataset for 0.5B and 1.5B models under different adaptation strategies. The table contrasts the performance between *Adaptation-Raw* (without removing repeated paths) and *Adaptation-Cleaned* (with repeated paths removed). Columns for MATH and Minerva Math are excluded, and the average is computed over the remaining datasets. **Bold** values indicate the best results.

Model	Method	GSM8K	GaoKao	Olympiad Bench	College Math	MMLU STEM	Avg.
<b>Math-QwQ-32B Dataset</b>							
Qwen2.5-0.5B-Instruct	Adaptation-Raw	47.5	29.1	9.3	26.8	28.2	28.2
	Adaptation-Cleaned	<b>49.6</b>	<b>30.9</b>	<b>9.3</b>	<b>27.5</b>	<b>37.5</b>	<b>31.0</b>
Qwen2.5-1.5B-Instruct	Adaptation-Raw	70.8	44.2	18.8	38.3	44.6	43.3
	Adaptation-Cleaned	<b>74.2</b>	<b>48.6</b>	<b>19.6</b>	<b>39.4</b>	<b>57.7</b>	<b>47.9</b>

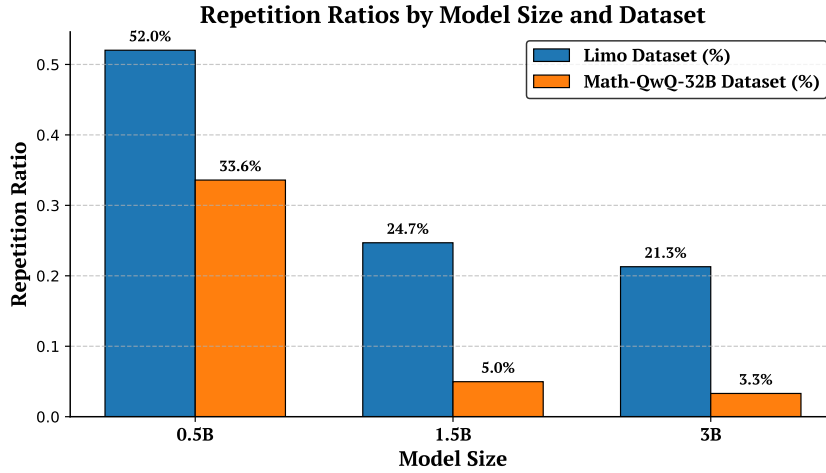


Figure 3: Repetition ratios(%) in search paths across different model sizes and datasets. Smaller models tend to have higher repetition ratios, particularly on the Limo dataset.

Consider an expert trajectory  $\tau_{\text{expert}} = (s_0, s_1, \dots, s_T)$  generated by a strong policy  $\pi_E$ , assumed to produce near-optimal steps. The student policy  $\pi_S$ , trained via behavioral cloning (BC) on expert demonstrations, minimizes the loss  $\mathbb{E}_{\tau \sim d^{\pi_E}} \left[ \sum_{t=1}^T -\log \pi_S(s_t | s_{<t}) \right]$ , where  $s_{<t} = (s_0, \dots, s_{t-1})$  is the prefix, and  $d^{\pi_E}$  is the expert state distribution.

Define the Q-value under  $\pi_S$  for appending the expert action  $s_t$  at prefix  $s_{<t}$ :

$$f_t = Q^{\pi_S}(s_{<t}, s_t) = \mathbb{E}_{s_{t+1:T} \sim \pi_S(\cdot | s_{\leq t})} [\mathbb{I}(\mathcal{O}(\tau) = a^*)]$$

where  $\tau = (s_0, \dots, s_T)$ ,  $\mathcal{O}(\tau)$  extracts the final answer,  $a^*$  is the ground truth,  $\mathbb{I}$  is the indicator function, and  $s_{\leq t} = (s_0, \dots, s_t)$ . Following (Li and Li), we use the sigmoid-transformed Q-value for probability interpretations:

$$f_t^\sigma = \sigma(f_t) = \mathbb{P}^{\pi_S}(\mathcal{O}(\tau) = a^* | s_{\leq t}).$$

**Lemma 1** (Existence of Imitation Gap). *There exists a step  $t_{\text{gap}} \in [1, T]$  such that the sequence of  $f_t$*

*values satisfies  $f_1 < f_2 < \dots < f_{t_{\text{gap}}-1}$ , followed by a sharp drop  $f_{t_{\text{gap}}} \ll f_{t_{\text{gap}}-1}$ .*

*Proof.* The proof is structured in three parts, leveraging Q-value rankings from process reward models (Li and Li) and the impact of distribution mismatch on the student policy.

**Part 1: Pre-gap monotonic increase.** For  $t < t_{\text{gap}}$ , the prefixes  $s_{<t}$  remain aligned with  $d^{\pi_E}$ , as the student policy  $\pi_S$  closely approximates  $\pi_E$ . Since the expert actions  $s_t$  are correct, we apply Lemma 3.3 from (Li and Li): for two correct steps  $s_n, s_m$  in a solution  $\tau$  with  $n < m$ , we have:

$$Q^*(s_{<n}, s_n) < Q^*(s_{<m}, s_m).$$

The proof, adapted to our student policy:

$$\begin{aligned} f_n^\sigma - f_m^\sigma &= \mathcal{P}^{\pi_S}(s_m | s_{<n}) \mathcal{P}^{\pi_S}(\tau | s_{\leq m}) \\ &\quad + \mathcal{P}^{\pi_S}(\overline{s}_m | s_{<n}) \mathcal{P}^{\pi_S}(\tau | \overline{s}_{\leq m}) \\ &\quad - \mathcal{P}^{\pi_S}(\tau | s_{\leq m}) \end{aligned}$$

$$= \mathcal{P}^{\pi_S}(\overline{s}_m | s_{<n}) [\mathcal{P}^{\pi_S}(\tau | \overline{s}_{\leq m}) - \mathcal{P}^{\pi_S}(\tau | s_{\leq m})],$$

Table 5: Comparison of truncation positions between **Adaptation-First** and **Adaptation-Gap** methods across datasets and model sizes. The relative localization difference represents the absolute difference between the relative truncation positions of these two methods. Higher differences are highlighted with deeper red.

Dataset	Model Size	First Position	Gap Position	Relative Localization Difference
Limo Dataset	0.5B	0.7901	0.7707	0.0194
	1.5B	0.6785	0.6985	0.0200
	3B	0.5718	0.5983	0.0265
Math-QwQ-32B Dataset	0.5B	0.5055	0.5244	0.0189
	1.5B	0.2113	0.3664	0.1551
	3B	0.0840	0.3239	0.2399

where the first equality uses the Q-function definition, and the second uses  $\mathcal{P}^{\pi_S}(s_m | s_{<n}) + \mathcal{P}^{\pi_S}(\overline{s_m} | s_{<n}) = 1$ . Under Assumption 3.1,  $\mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq m}}) - \mathcal{P}^{\pi_S}(\tau | s_{\leq m}) < 0$ , since the correct step  $s_m$  has a higher probability of leading to a correct outcome. Thus, for  $n < m$ ,  $f_n^\sigma < f_m^\sigma$ , implying  $f_n < f_m$ . Since  $\pi_S \approx \pi_E$  for early steps, this holds for all  $t < t_{\text{gap}}$ , yielding:

$$f_1 < f_2 < \dots < f_{t_{\text{gap}}-1}.$$

**Part 2: Distribution mismatch and emergence of non-optimal step.** Due to differences in model capacity (e.g., the student being a smaller model), the expert data distribution  $d^{\pi_E}$  and the student model distribution  $d^{\pi_S}$  are inconsistent. As the number of steps increases, the prefixes  $s_{<t}$  grow increasingly complex, becoming likely to fall outside the training distribution of  $\pi_S$ . Consequently, the state observed by  $\pi_S$  at step  $t$  diverges from that of  $\pi_E$ , such that the expert action  $s_t$ , optimal under  $\pi_E$ , is not necessarily optimal under  $\pi_S$ . This distribution mismatch leads to a critical step  $t_{\text{gap}}$  where the expert action  $s_{t_{\text{gap}}} = s_E$  is non-optimal for  $\pi_S$ , as it does not maximize the expected reward under the student’s policy:

$$Q^{\pi_S}(s_{<t_{\text{gap}}}, s_E) < \max_{s \in \mathcal{A}} Q^{\pi_S}(s_{<t_{\text{gap}}}, s).$$

This non-optimality arises because the OOD prefix  $s_{<t_{\text{gap}}}$  causes  $\pi_S$  to misjudge the value of  $s_E$ , favoring an alternative action that aligns better with its biased distribution, analogous to selecting an incorrect step from a correct prefix.

**Part 3: Sharp drop behavior.** At  $t_{\text{gap}}$ , appending the non-optimal expert action  $s_{t_{\text{gap}}}$  produces an OOD state  $s_{\leq t_{\text{gap}}}$ , significantly reducing the probability of correct completion. We compare the Q-value of the correct prefix at  $t_{\text{gap}} - 1$  to the non-optimal step at  $t_{\text{gap}}$ . For the correct prefix at  $t_{\text{gap}} - 1$ ,

let  $s_{t_{\text{gap}}-1}$  be correct, and for the non-optimal step  $s_{t_{\text{gap}}}$ , we have:

$$f_{t_{\text{gap}}-1}^\sigma - \mathcal{V}^{\pi_S}(x) = \mathcal{P}^{\pi_S}(\overline{s_{t_{\text{gap}}-1}} | x) \left( \mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}-1}) - \mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}-1}}) \right).$$

$$f_{t_{\text{gap}}}^\sigma - \mathcal{V}^{\pi_S}(x) = \mathcal{P}^{\pi_S}(s_{t_{\text{gap}}} | x) \left( \mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}}}) - \mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}}) \right).$$

where  $\mathcal{V}^{\pi_S}(x) = \mathbb{P}^{\pi_S}(\mathcal{O}(\tau) = a^* | x)$ . Under Assumption 3.1,  $\mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}-1}) > \mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}-1}})$ , so the first difference is positive, implying  $f_{t_{\text{gap}}-1}^\sigma > \mathcal{V}^{\pi_S}(x)$ . For the non-optimal step,  $\mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}}}) < \mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}})$ , and since  $s_{t_{\text{gap}}}$  is non-optimal due to distribution mismatch,  $\mathcal{P}^{\pi_S}(s_{t_{\text{gap}}} | x) \gg \mathcal{P}^{\pi_S}(\overline{s_{t_{\text{gap}}} | x})$ , amplifying the negative difference. Thus:

$$f_{t_{\text{gap}}}^\sigma < \mathcal{V}^{\pi_S}(x) < f_{t_{\text{gap}}-1}^\sigma,$$

implying  $f_{t_{\text{gap}}} \ll f_{t_{\text{gap}}-1}$ , as the non-optimal step’s Q-value is significantly lower due to the low probability of recovery from incorrect branches.

The key size relation for the drop is:

$$f_{t_{\text{gap}}}^\sigma - f_{t_{\text{gap}}-1}^\sigma = \mathcal{P}^{\pi_S}(s_{t_{\text{gap}}} | s_{<t_{\text{gap}}}) \left[ \mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}}) - \mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}}}) \right] + \mathcal{P}^{\pi_S}(\overline{s_{t_{\text{gap}}} | s_{<t_{\text{gap}}})} \left[ \mathcal{P}^{\pi_S}(\tau | \overline{s_{\leq t_{\text{gap}}}}) - \mathcal{P}^{\pi_S}(\tau | s_{\leq t_{\text{gap}}}) \right] < 0.$$

where the negative term dominates under Assumption 3.1, ensuring  $f_{t_{\text{gap}}} \ll f_{t_{\text{gap}}-1}$ .  $\square$

## E Case Study

To elucidate the imitation gap’s impact on adaptability score decline and reasoning performance,

1014 we present a case study on the LIMO dataset, pin-  
1015 pointing its onset in a complex reasoning task  
1016 through Qwen2.5-3B-Instruct’s step-wise adapt-  
1017 ability scores (see Table 6).

## 1018 **F Description of Large Language Model** 1019 **Usage**

1020 In the preparation of this manuscript, we lever-  
1021 aged a large language model (LLM), specifically  
1022 Grok developed by xAI, to facilitate specific as-  
1023 pects of the writing process. The LLM was em-  
1024 ployed primarily for linguistic refinement, encom-  
1025 passing tasks such as enhancing sentence coher-  
1026 ence, improving syntactic clarity, and elevating the  
1027 overall readability of the text, while preserving  
1028 the integrity of the scientific content, methodolo-  
1029 gies, and findings. The rationale for this approach  
1030 was to optimize the communicative efficacy of the  
1031 manuscript, ensuring that intricate technical con-  
1032 cepts are articulated with precision and accessibil-  
1033 ity for a diverse academic readership. All outputs  
1034 generated by the LLM were subjected to rigorous  
1035 scrutiny, validation, and, where necessary, revision  
1036 by the authors to uphold the principles of accuracy,  
1037 originality, and academic rigor. Notably, the LLM  
1038 was not utilized for the generation of novel intellec-  
1039 tual contributions, experimental frameworks, data  
1040 analyses, or conclusions, which were exclusively  
1041 derived from human expertise. This judicious ap-  
1042 plication of LLMs adheres to established ethical  
1043 standards for AI-assisted academic writing, balanc-  
1044 ing the enhancement of textual quality with a com-  
1045 mitment to transparency and scholarly integrity.

Table 6: Case Study on the LIMO Dataset: Solutions are generated by Qwen2.5-3B-Instruct. Step-wise adaptability scores are reported for each reasoning step, with steps subsequent to the Imitation Gap (starting from Step 151) highlighted in gray.

Q: The vertices of a regular nonagon (9-sided polygon) are to be labeled with the digits 1 through 9 in such a way that the sum of the numbers on every three consecutive vertices is a multiple of 3. Two acceptable arrangements are considered to be indistinguishable if one can be obtained from the other by rotating the nonagon in the plane. Find the number of distinguishable acceptable arrangements.	Scores
...	
Step 145: Therefore, except in the case of $d=1$ , which would require all digits being same, but which is impossible with unique digits, otherwise, no fix. Therefore, fixed labelings only for identity, which are all our 1296 acceptable labelings.	0.750
Step 146: Therefore, the number of distinct arrangements is 1296 divided by 9, which is 144.	0.938
Step 147: So, is 144 the answer? But wait, let me cross-verify.	0.875
Step 148: Wait, though: in our prior thought process, we had 6 residue assignments, each giving 216 colorings. So $6 \times 216 = 1296$ .	0.875
Step 149: If we consider the condition that rotating must not create overcounts, then perhaps we had overcounted by factor 9, so we have to divide by 9. $1296 / 9 = 144$ .	1.000
Step 150: Alternatively, to think in terms of necklace counting: the number of distinct necklaces with beads labeled 1-9, under rotation, with a constraint that any three consecutive numbers sum to multiple of 3. So, the formula, via Burnside, gives 144. Does that make sense?	0.938
Step 151: But, another route. Let me think, maybe I have miscalculated for fixed labelings. Wait, notice that the fixed labelings under rotation by $k=3$ positions?	0.438
Step 152: Wait, because $d = \gcd(3,9) = 3$ , so the cycles decompose into 3 cycles, each of length 3. So, fixed labelings here would have to assign the same digit to each cycle, but since labels must be unique, is it possible?	0.688
Step 153: Wait, meaning that for each cycle (1,4,7), (2,5,8), (3,6,9); all three digits in cycle (1,4,7) must be the same, similarly for others, which is impossible with unique digits.	0.938