# AdapThink: Adaptive Thinking Preferences for Reasoning Language Model

**Xu Wan**
Zhejiang University
wanxu@zju.edu.cn

**Wei Wang**
Alibaba DAMO Academy
richard.wang1@unimelb.edu.au

**Wenyue Xu**
Tongji University
2111141@tongji.edu.cn

**Wotao Yin**
Alibaba DAMO Academy
wotao.yin@alibaba-inc.com

**Jie Song**
Peking University
songjie@coe.pku.edu.cn

**Mingyang Sun**[*]
Peking University
smy@pku.edu.cn

## Abstract

Reinforcement Learning (RL)-based post-training has significantly advanced the complex reasoning capabilities of language models, fostering sophisticated self-reflection processes. However, this "slow thinking" paradigm presents a critical challenge to reasoning efficiency: models may expend excessive computation on simple questions and shift reasoning prematurely for complex ones. Previous mechanisms typically rely on static length budgets or predefined rules, lacking the adaptability for varying question complexities and models' evolving capabilities. To this end, we propose *AdapThink*, an *adap*tive post-training framework designed to induce more efficient *think*ing while maintaining the performance of reasoning language models. Specifically, AdapThink incorporates two key mechanisms: 1) A group-relative reward function that leverages model confidence and response's characteristic to dynamically adjust the preference of reflection-related transition words without resorting to a fixed length preference. 2) A diversity-aware sampling mechanism that balances the training group's solution accuracy with reasoning diversity via an entropy-guided score. Experiments on several mathematical reasoning datasets with DeepSeek-distilled models demonstrate AdapThink's advantages in enabling adaptive reasoning patterns and mitigating the inefficiencies.

## 1 Introduction

Recent breakthroughs in large language models such as OpenAI's o1 [16] and DeepSeek R1 [3] have demonstrated that reinforcement learning (RL)-based post-training methods can substantially improve their reasoning capabilities. This enhancement is primarily attributed to the emergence of models' sophisticated self-reflection behaviors [8, 7]. However, recent research has highlighted a significant inefficiency associated with this "*slow thinking*" pattern [14, 4, 2, 23, 1, 21]. Reasoning models frequently *overthink* simple problems, thereby spending unnecessary computational resources, while conversely *underthinking* complex challenges, leading to incomplete reasoning and incorrect answers. When presented with a simple problem (see Appendix A), a model can reach a correct answer with merely 479 tokens. Yet, its self-validation mechanisms—marked by phrases such as "Verify" and "Wait"—triggered unnecessary reflections, resulting in more than a quadrupled token consumption. In contrast, when attempting a complex problem, the same model exhibited frequent and unproductive shifts, marked by phrases "Alternatively" and "Another." Upon reaching its token limit, the model stopped and arrived at an incorrect answer.

---

[*]Coressponding author.

Therefore, an ideal chain of thought (CoT) would be capable of adjusting its self-reflection frequency and depth, adaptive to the problem difficulty and its level of confidence. Empirically, [11, 25] have demonstrated that employing reflection vocabularies does *not always* guarantee mathematical reasoning improvement. Drawing upon this insight, recent works explored direct control via modifying the input prompt [5, 14, 9, 4, 2] or providing indirect reasoning-length rewards [23, 1, 21]. A common limitation of these approaches is using token budgets governed by rules or offering rewards for adhering to a special length budget, overlooking the critical impact of the changes in models' capabilities and length preferences. Intuitively, models with limited reasoning abilities benefit from more extended CoT patterns, as redundant self-reflection could serendipitously contribute to reaching correct solutions. In contrast, high-performing models should aim to minimize token consumption to prevent overthinking and to maintain efficient CoT reasoning.

Motivated by these insights and our observations in Section 3, we propose **AdapThink**, an *RL post-training framework* for reasoning models. Our post-training empowers models to adjust their preference for reasoning depth, which is adaptive to their current operating capabilities. Instead of directly limiting the budget of reasoning length, our work adjusts the length preference through analyzing the distribution of diverse reasoning patterns observed in *groups of generated samples*. We find that reasoning length alone does not directly determine model reasoning performance. Instead, strategically regulating reasoning depth is a more effective way to mitigate overthinking and underthinking. Specifically, our main contributions are:

1. First, we introduce *a group-relative reward function*, a novel mechanism designed to adjust the model's current reasoning preferences. With it, the model is trained to determine the appropriate preferences for reflection based on intra-group accuracy of generated responses, while measuring reasoning efficiency *quantitatively* through counting key transition words that appear in these groups of training samples.

2. Furthermore, we propose a *diversity-aware sampling mechanism* to balance between achieving high reasoning accuracy and fostering rich diversity within the groups of training samples. To achieve this, we begin with an oversampling of reasoning instances. We then employ carefully defined *diversity metrics* to evaluate both the final answers and the intermediate steps of these instances, before applying diversity-aware downsampling to curate and enhance the overall quality of the instances used for RL post-training.

3. By post-training the DeepSeek-distilled Qwen model subject to a context length limit of only *2K tokens*, our method achieves superior performance compared to several length-control baselines when tested within an 8K-token limit, demonstrating AdapThink's effectiveness in developing powerful and efficient CoT models across multiple reasoning benchmarks.

## 2   Related Work

**Induce Longer CoT**   For inducing longer reasoning length, several works [24, 13, 19, 18, 5] have encouraged models to engage in deeper thinking through natural language feedback. For zero-shot CoT, [24, 13] stimulate model self-reflection by performing backward verification or multiple response voting. For few-shot CoT with demonstrations, [5] introduced five general standardized patterns to induce models to simulate human thinking and reshape the CoT. However, these one-size-fits-all approaches ignore the diversity of problem-solving paths, only have a monotonous processing mode. Similarly, [14] designed several budget forcing mechanisms to increase guidance in CoT, such as appending "Wait" multiple times when the model tries to end, forcing it to double-check. Likewise, [20] introduced special meta-action markers like <|continue|>, <|reflect|>, and <|explore|>, enabling the model to restart from intermediate steps, lengthen responses, and correct errors.

**Induce Shorter CoT**   Simple prompt methods are also effective for encouraging a shorter CoT [15, 27, 17]. For example, [14] added "Final answer" to terminate the model's thinking process. Besides, [5, 6] used stronger models to compress long CoTs semantically into shorter ones. [28, 21, 1, 29] design different length-budget signals from the perspective of rule-based reward design, encouraging models to balance accuracy and token efficiency during the thinking process. Although [21] already takes problem complexity and model confidence into account for the budget, its adoption of a uniformly shorter response preference is less suitable for models with weaker reasoning abilities. Moreover, while [29, 30] uses dynamic sampling to avoid the zero-advantage issue in group-relative
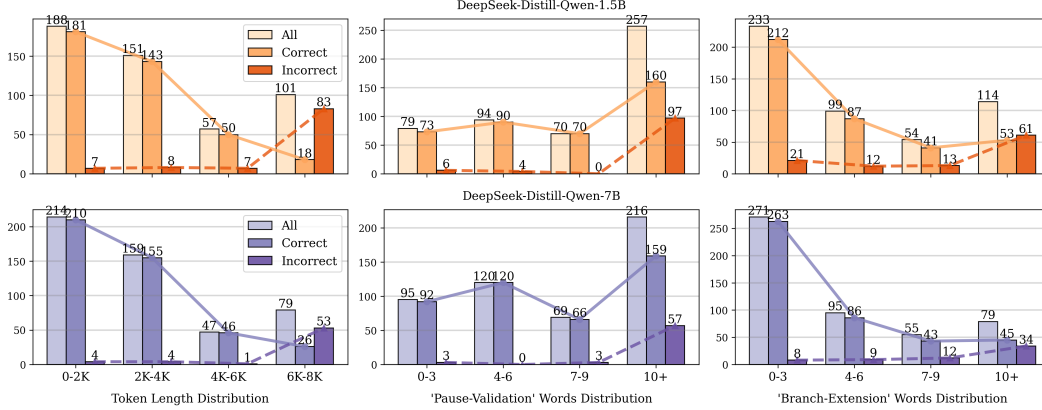
Figure 1: Distribution analysis for DeepSeek-Distill-Qwen 1.5B (top), 7B (bottom) on MATH-500. X-axes: token ranges (left), word frequency intervals (middle/right); Y-axes: sample counts.

advantage calculation, its sampling strategy focuses solely on outcome diversity metrics and neglects intra-group reasoning process diversity.

Moreover, some studies have further focused on addressing overthinking and underthinking in long CoTs. [23, 22] introduced thought switching penalties to influence the token decoding probability distribution early in CoT generation, reducing initial thought-switching. Similarly, [2] presented an efficiency metric to evaluate each token's contribution to accuracy and used length preference optimization to achieve more efficient CoT patterns. However, both approaches depend on auxiliary judgments from more powerful reasoning models, making the performance improvements inherently constrained by the capabilities of the reference models.

## 3 Observations

To investigate potential overthinking and underthinking issues in current language models, we first conducted a comprehensive analysis of the generation patterns from DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B models [3]. We examined the CoT generation efficiency from three perspectives:

- **Token Length Distribution**: Following the inference settings in [1], we set the maximum token limit to 8,192 during inference and segmented the model's response range into four equal intervals to analyze the distribution patterns.
- **"Pause-Validation" Words Distribution**: We analyzed the frequency of transition words in the responses, focusing on four key indicators of "Pause-Validation": "*wait*", "*hold on*", "*check*", and "*verify*". The distribution was categorized into four bins with an interval of 3.
- **"Branch-Extension" Words Distribution**: Similarly, we identified four representative markers of "Branch-Extension": "*alternatively*", "*however*", "*another*", and "*instead*". Their occurrences were also categorized into four bins with an interval of 3.

As shown in Figure 1, incorrect answers typically correlate with higher token consumption and increased usage of both "Pause-Validation" and "Branch-Extension" words. This suggests that when confronting challenging problems, models tend to exhibit uncertain behavior through excessive alternative explorations and repetitive self-verifications. Intriguingly, for correct answers, we observed that the responses predominantly fall within lower token ranges and show lower frequencies of "Branch-Extension" word occurrences. Notably, while correct answers demonstrate reduced exploration of alternative paths, they maintain consistent levels of "Pause-Validation" expressions, suggesting a potential distinction between the roles of different transition words in the reasoning process. These observations motivate us to explore effective ways to achieve more efficient reasoning while maintaining performance, with particular attention to the distinct roles of different transition words in the reasoning process.

# 4 Methodology

In this section, we introduce the proposed post-training framework AdapThink in detail. Unlike previous approaches that *directly impose length restrictions* or set *token budgets* [21, 28, 1, 14], we propose to modulate the reasoning behavior by regulating the preference of reflection words based on group-level response characteristics. To this end, AdapThink introduces two key innovations: (1) a group-relative reward function for controlling reasoning preferences, and (2) a diversity-aware sampling mechanism for ensuring intra-group diversity during training. The pseudo-code of AdapThink can be found in Appendix C.

## 4.1 Group-relative Reasoning Preference Reward

We begin with a pre-trained reasoning language model $\pi_\theta$ and a dataset $\mathcal{D} = \{(x^k, y^k)\}_{k=1}^N$, where each instance contains the input prompt $x$ and the final answer $y$. For each $x \in \mathcal{D}$, $\pi_\theta$ performs reasoning to generate $|\mathcal{G}|$ samples $\mathcal{G} := \{y_i\}_{i=1}^{|\mathcal{G}|}$. Next, we analyze these samples from two perspectives:

- **Model Confidence.** We represent the average correctness of the model's $|\mathcal{G}|$ samples within a group as the model confidence, denoted as $\varphi$.

$$\varphi = \frac{1}{|\mathcal{G}|} \sum_{i=1}^{|\mathcal{G}|} \mathbb{I}(y_i = y), \tag{1}$$

  where $\mathbb{I}(\cdot)$ is the indicator function . A larger $\varphi$ reflects more reliable reasoning patterns for the current question $x$. We set $\varphi_{\text{low}}$ and $\varphi_{\text{high}}$ as the boundaries for low and high confidence levels. $\varphi$ directly influences the p*reference direction* of the AdapThink's reasoning preference reward.

- **Response Characteristics.** Given the distinct reasoning patterns observed between correct and incorrect answers, we partition the model's responses into two groups based on their correctness: $\mathcal{G}_\mathcal{T}$ for correct answers and $\mathcal{G}_\mathcal{F}$ for incorrect ones. The distribution of transition words within these groups determines the *specific values* of the reasoning preference reward.

Then, based on the reflection words analysis in Section 3, we aim to control the usage of "Branch-Extension" transition words while promoting efficient output generation. Specifically, we introduce three reward components to guide the reasoning process: token consumption reward $\lambda_l$ that quantifies the response length, output completion reward $\lambda_o$ that tracks the presence of answer completion marker in the `</think>` tag, and "Branch-Extension" transition words reward $\lambda_b$ that measures the numbers of four switch words. Each component reward $\lambda_*$ is calculated as the normalized deviation from its corresponding group mean:

$$\lambda_*(y_j) = \frac{r_*(y_j)}{\mu_*(\mathcal{G}(y_j))} - 1, \quad * \in \{l, o, b\} \tag{2}$$

where $\mathcal{G}(y_j)$ maps response $y_j$ to either $\mathcal{G}_\mathcal{T}$ or $\mathcal{G}_\mathcal{F}$ depending on whether $y_i$ is correct or not, and $\mu_*(\mathcal{G})$ represents the mean measurement within group $\mathcal{G}$.

Then, we introduce the group-relative reasoning preference reward (GRPR) to adaptively enable different reward components based on the model's confidence $\varphi$ for the current question:

$$r(x, \mathcal{G}, \theta) = \text{clip}\Big(|\omega(\varphi)|(\lambda_o - \lambda_l) + \mathbb{I}(\omega(\varphi) < 0)\omega(\varphi)\lambda_b, r_{\min}, r_{\max}\Big) \tag{3}$$

where the clip function bounds the reward within $[r_{\min}, r_{\max}]$ to ensure stable training signals. Here, $\omega(\varphi)$ is a cosine interpolation function that smoothly transitions between reward preferences at confidence boundaries:

$$\omega(\varphi) = \begin{cases} +1 & \text{if } \varphi \leq \varphi_{\text{low}} \\ \cos\left(\frac{\varphi - \varphi_{\text{low}}}{\varphi_{\text{high}} - \varphi_{\text{low}}}\pi\right) & \text{if } \varphi_{\text{low}} < \varphi < \varphi_{\text{high}} \\ -1 & \text{if } \varphi \geq \varphi_{\text{high}} \end{cases} \tag{4}$$

For low-confidence cases ($\omega(\varphi) > 0$), the GRPR reward emphasizes output completion while controlling length; for high-confidence cases ($\omega(\varphi) < 0$), it additionally controls the branch-extension exploration.
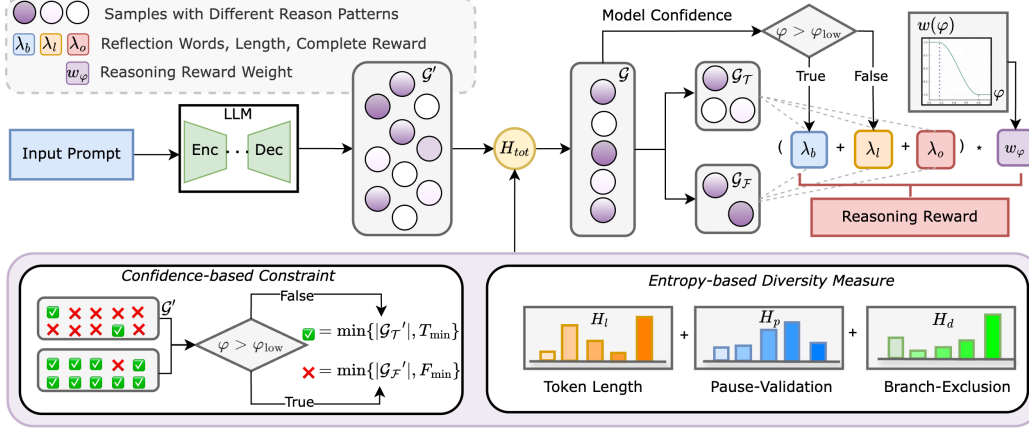
Figure 2: The AdapThink framework architecture, featuring confidence-guided sample selection and entropy-based diversity measurement for adaptive reasoning reward computation.

## 4.2 Diversity-aware Sampling

As shown in Equation 3, since GRPR incorporates group-specific distribution patterns from $\mathcal{G}_{\mathcal{T}}$ and $\mathcal{G}_{\mathcal{F}}$ into its reward computation, the diversity of reason patterns within each group significantly influences the effectiveness of GRPR-based post-training. While previous studies have addressed the zero-advantage problem through dynamic upsampling strategies [29] to avoid homogeneous output labels within groups, they haven't considered the diversity of reasoning processes.

To address this limitation, we propose a two-stage sampling strategy that enhances reasoning diversity while maintaining group balance:

**Stage 1: Oversampling and Diversity Measurement** First, we oversample by a factor of $K$ to expand the sample pool to $\mathcal{G}'$ candidates, allowing exploration of diverse reasoning styles within each group. To evaluate sample diversity, we introduce entropy-based measures $H_l$, $H_p$, and $H_b$ for token length,"Pause-Validation" words and "Branch-Extension" words, respectively:

$$H_*(\mathcal{G}') = -\frac{1}{\log |\mathcal{G}'|} \sum_{s \in \mathcal{S}} \frac{|s|}{|\mathcal{G}'|} \log \frac{|s|}{|\mathcal{G}'|}, \quad * \in \{l, b, p\} \tag{5}$$

where we partition each component's measurement range into four equal intervals, forming bins $\mathcal{S}$. The normalized entropy ensures $H_* \in [0, 1]$, with higher values indicating a more uniform distribution for reflection words count. The overall diversity score combines these dimensions:

$$H_{\text{tot}} = \alpha_l H_l + \alpha_p H_p + \alpha_b H_b \tag{6}$$

where $\alpha_l$, $\alpha_p$ and $\alpha_b$ denote the weights for length, "Pause-Validation" words and "Branch-Extension" words diversity score, respectively.

**Stage 2: Confidence-constrained Downsampling** To enhance RL training efficiency while maintaining sample diversity, we select $|\mathcal{G}|$ samples from the oversampled pool $\mathcal{G}'$, which consists of correct answers $\mathcal{G}_{\mathcal{T}}'$ and incorrect answers $\mathcal{G}_{\mathcal{F}}'$. To prevent model bias towards a single output label, we calculate the minimum number of samples required from each category as shown in Equation 8:

$$T = \min\{|\mathcal{G}_{\mathcal{T}}'|, T_{\min}\}, \quad F = \min\{|\mathcal{G}_{\mathcal{F}}'|, F_{\min}\} \tag{7}$$

where $T_{\min}$ and $F_{\min}$ represent the minimum requirements for correct and incorrect samples, respectively. Based on the model's confidence level $\varphi$, we adapt our selection strategy to prioritize either $T$ correct or $F$ incorrect samples while maximizing reasoning diversity. The final sample selection process adapts to the model's confidence level $\varphi$, as formulated in Equation 8:

$$\mathcal{G} = \begin{cases} \underset{\mathcal{G}_c \subseteq \mathcal{G}_{\mathcal{T}}', |\mathcal{G}_c|=T}{\arg\max} H_{\text{tot}}(\mathcal{G}_c) \cup \underset{\mathcal{G}_s \subseteq \mathcal{G}_{\text{remain}}', |\mathcal{G}_s|=|\mathcal{G}|-T}{\arg\max} H_{\text{tot}}(\mathcal{G}_s) & \text{if } \varphi \le \varphi_{\text{low}} \\ \underset{\mathcal{G}_c \subseteq \mathcal{G}_{\mathcal{F}}', |\mathcal{G}_c|=F}{\arg\max} H_{\text{tot}}(\mathcal{G}_c) \cup \underset{\mathcal{G}_s \subseteq \mathcal{G}_{\text{remain}}', |\mathcal{G}_s|=|\mathcal{G}|-F}{\arg\max} H_{\text{tot}}(\mathcal{G}_s) & \text{otherwise} \end{cases} \tag{8}$$

where $\mathcal{G}_{\mathcal{T}}'$ and $\mathcal{G}_{\mathcal{F}}'$ denote the pools of correct and incorrect answers from oversampling, respectively. When model confidence is low ($\varphi \le \varphi_{\text{low}}$), the algorithm first selects $T$ correct samples $\mathcal{G}_c$ from $\mathcal{G}_{\mathcal{T}}'$ that maximize diversity; otherwise, it prioritizes selecting $F$ incorrect samples from $\mathcal{G}_{\mathcal{F}}'$. The remaining $\mathcal{G}_s$ training samples are then selected from the candidate pool $\mathcal{G}_{\text{remain}}'$ to reach the target size $|\mathcal{G}|$, while maintaining maximum diversity measured by $H_{\text{tot}}$.

## 5 Experimental Setup

**Datasets and Pretrained LLMs** We conduct experiments on a curated lightweight mathematics dataset that spans various difficulty levels. This dataset combines queations from *DeepScaleR-Preview-Dataset* [10], including about 5K question-answer pairs sampled from AIME (1984-2023), AMC (prior to 2023), and MATH training sets. For baseline models, we employ *DeepSeek-R1-Distill-Qwen-1.5B*, which has demonstrated competitive performance against GPT-4-0513 and QwQ-32B-Preview on various mathematical reasoning benchmarks [3].

**Baselines** We conduct post-training on DeepSeek-R1-Distill-Qwen-1.5B using Group Relative Policy Optimization (GRPO) [3] as our base algorithm. GRPO optimizes policy $\pi_\theta$ by maximizing the following objective:

$$\mathcal{J}(\theta) = \mathbb{E}[x \sim P(x), \mathcal{G} \sim \pi_{\theta_{\text{old}}}(y|x)]$$

$$\frac{1}{|\mathcal{G}|} \sum_{i=1}^{|\mathcal{G}|} \left( \min\left( \frac{\pi_\theta(y_i|x)}{\pi_{\theta_{\text{old}}}(y_i|x)} A_i, \text{clip}\left( \frac{\pi_\theta(y_i|x)}{\pi_{\theta_{\text{old}}}(y_i|x)}, 1-\epsilon, 1+\epsilon \right) A_i \right) - \beta \mathbb{D}_{\text{KL}}(\pi_\theta || \pi_{\text{ref}}) \right) \tag{9}$$

where $A_i$ is the advantage computed using binary accuracy reward $r_i = \mathbb{I}(y_i = y)$ within $\mathcal{G}$:

$$A_i = \frac{r_i - \text{mean}(\{r_1, r_2, \ldots, r_{|\mathcal{G}|}\})}{\text{std}(\{r_1, r_2, \ldots, r_{|\mathcal{G}|}\})} \tag{10}$$

To evaluate different strategies for controlling reasoning length and efficiency, we implement the following baselines with different reward setting:

*(1) Length Controlled Policy Optimization (LCPO) [1]*: it extends GRPO by modifying its reward function to include a correctness reward and a length penalty:

$$r_i = \mathbb{I}(y_i = y) - \alpha \cdot |l_i - l^*| \tag{11}$$

where $l_i$ is the generated output length, and $\alpha$ is a scalar that regulates the trade-off between generating the correct answer and meeting the target length $l^*$, which we set to 2,048 tokens during training.

*(2) Token Length Budget (TLB) [21]*: it establishes a token length budget metric $l_b$ based on sampling accuracy and length statistics, and designs a difficulty-adaptive reward as follows:

$$r_i = \begin{cases} \max(-0.5\lambda + 0.5, 0.1) & \text{if } y_i = y \\ \min(0.9\lambda - 0.1, -0.1) & \text{otherwise} \end{cases} \tag{12}$$

where $\lambda = l_i/l_b - 1$ denotes the relative length reward.

*(3) Cosine Reward Function (CosFn) [28]*: it introduces a cosine length scaling reward:

$$r_i = r_{\min} + \frac{1}{2}(r_{\max} - r_{\min})(1 + \cos(\frac{l_i \, \pi}{l_{\max}})) \tag{13}$$

where different pairs of $(r_{\min}, r_{\max})$ are adopted from [28]: $(0.5, 1.0)$ for correct answers and $(-1, -0.5)$ for incorrect ones, with $l_{\max}$ set to 2,048 tokens.

**Evaluation Protocol**    We evaluate our method from three perspectives:

*(1) Math Accuracy.* We adopt PASS@1 $= \frac{1}{|\mathcal{G}|} \sum_{i=1}^{|\mathcal{G}|} p_i$ as our primary accuracy metric following [3], where $p_i$ indicates the correctness of the $i$-th response.

*(2) Reasoning Efficiency.* We measure efficiency through two metrics: (1) the average token length $\overline{L} = \frac{1}{|\mathcal{G}|} \sum_{i=1}^{|\mathcal{G}|} l_i$, and (2) the average frequency of reflection words $\overline{n}_* = \frac{1}{|\mathcal{G}|} \sum_{i=1}^{|\mathcal{G}|} n_*(y_i)$ in $|\mathcal{G}|$ responses, where $n_* := \{n_p, n_b\}$ represent the counts of "Pause-Validation" and "Branch-Extension" words, respectively.

*(3) Group Diversity.* We quantify the diversity of $\mathcal{G}$ using the proposed entropy-based measures $H_l$, $H_p$, and $H_b$, which evaluate the distribution uniformity of response length, "Pause-Validation" words, and "Branch-Extension" words respectively.

**Training Details and Hyperparameters**    For AdapThink training, we build upon the PEFT framework [12] using LoRA adaptation. Training uses a batch size of 32 and runs for 5 epochs with bfloat16 precision. The maximum token limit is set to 2K tokens during training, it is extended to 8K tokens for evaluation. For sampling, we generate 16 rollouts per question using temperature 0.7 and top-p 0.95, then select 8 samples based on our diversity-aware strategy with equal weights $\alpha_l = \alpha_b = \alpha_d = 1$. The confidence thresholds are set to $\varphi_{\text{low}} = 0.15$ and $\varphi_{\text{high}} = 0.5$. Following [26]'s observation that total incorrect samples degrade GRPO training performance, we set a large $T_{\min} = 3$ and $F_{\min} = 1$ to preserve more correct samples for hard questions.

# 6    Results

In this section, we evaluate and analyze the performance of AdapThink across various settings through three key aspects: answer accuracy, reasoning efficiency, and group diversity.

## 6.1    Overall Comparison

As shown in Figure 3 and Table 1, *AdapThink significantly improves accuracy while reducing response length, demonstrating distinct patterns for different types of reflection words.* During training, AdapThink exhibits a notable decrease in switch words, while check words follow a pattern of initial increase followed by decline. This control over key reasoning words effectively generalizes to testing scenarios, resulting in a 27% average performance improvement and 15.6% reduction in response length compared to the base model.

For LCPO, we observe that the late-stage increase in check words does not correspond to improvements in completion or accuracy rewards. Interestingly, LCPO's strict length control shows a preference for shorter responses but fails to generalize to effective reasoning patterns, leading to an average 8.2% lower accuracy compared to AdapThink during an 8K token limit.

For the GRPO and TLB methods, they demonstrated consistent minimal changes in reflection words during training. During the testing phase, their responses consumed an average of 215 and 488 more tokens per question respectively compared to AdapThink. While CosFn's separate length control for correct and incorrect answers leads to improved accuracy, its reasoning efficiency falls short of AdapThink.

## 6.2    Ablation Studies

To verify the effectiveness of each component in the AdapThink framework, we conducted three ablation experiments, and the evaluation results are shown in Table 2.

**Reward Components.**    First, we performed ablation studies on the reward components $\lambda_l$ and $\lambda_b$ in GRPR. Without length control reward $\lambda_l$, we observed a substantial increase in pause-validation words, rising by 244.8% for correct samples and 173.2% for incorrect ones. The outputs became significantly longer, with $L^{\mathcal{T}}$ increasing by 17.7% and a 2.1 percentage point accuracy drop. Similarly, removing $\lambda_b$ resulted in less efficient reasoning patterns with pause-validation frequency increasing by 65.5% and branch-extension usage rising by 55.0%, and a 2.5 percentage point drop in accuracy.
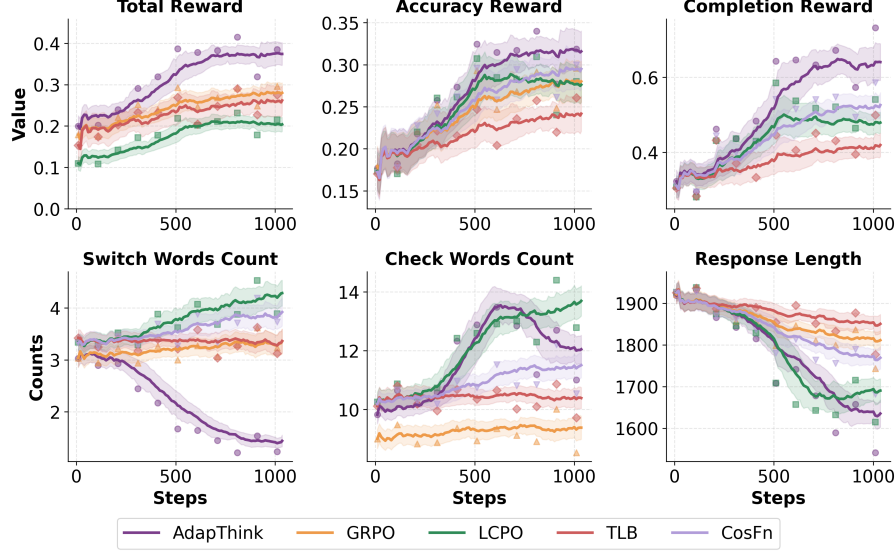
7

Figure 3: Training comparison among five post-training frameworks on DeepSeek-Distill-Qwen-1.5B with a 2K token length constraint.

Table 1: Evaluation performance comparison across mathematical reasoning benchmarks with an 8K token length constraint.

| Model | Method | AIME2025 | | AIME2024 | | MATH500 | | AMC | |
|---|---|---|---|---|---|---|---|---|---|
| | | PASS@1 | Token | PASS@1 | Token | PASS@1 | Token | PASS@1 | Token |
| DeepSeek Distill Qwen -1.5B-2K | Baseline | 17.92 | 7242 | 20.00 | 7451 | 77.25 | 3736 | 52.26 | 5743 |
| | +GRPO | 24.17↑6.25 | 6763↓479 | 26.67↑6.67 | 6994↓457 | 81.40↑4.15 | 3090↓646 | **62.95**↑10.69 | 4839↓904 |
| | +TLB | 20.83↑2.91 | 6984↓258 | 22.50↑2.50 | 7200↓251 | 80.40↑3.15 | 3352↓384 | 58.13↑5.87 | 5239↓504 |
| | +LCPO | 20.83↑2.91 | 6724↓518 | 29.17↑9.17 | 6367↓1084 | 78.60↑1.35 | 2523↓1213 | 61.14↑8.88 | 4649↓1094 |
| | +CosFn | 21.67↑3.75 | 6956↓286 | 28.33↑8.33 | 7175↓276 | 82.70↑5.45 | 2729↓1007 | 61.14↑8.88 | 4847↓896 |
| | +AdapThink | **25.42**↑7.50 | 6455↓787 | **31.67**↑11.67 | 6812↓639 | **82.80**↑5.55 | 2742↓994 | 62.24↑9.98 | **4814**↓929 |
| FastCuRL-1.5B-24K[‡] | | 25.42 | 7085 | 31.67 | 7117 | 86.40 | 3430 | 62.80 | 5440 |
| DeepScaler-1.5B-24K[‡] | | 25.42 | 6397 | 37.50 | 6423 | 88.00 | 2831 | 66.87 | 4893 |

[‡]Originally trained with 24K token limit. We evaluate their official models with an 8K token limit.

**Diversity-aware Selection.** Subsequently, we investigated the effectiveness of the diversity-aware selection mechanism by setting different oversampling ratios $K$. Our experiments revealed that $K = 2$ yielded better performance across all metrics, achieving 2.9% higher accuracy with less token consumption. The higher diversity scores $H_l$ and $H_d$ suggest that the selection strategy promotes varied reasoning patterns in the training phase and accelerates the efficient exploration of RL.

**Curriculum Learning.** Inspired by [10], we implemented a curriculum learning for AdapThink. Specifically, we conducted secondary training on the best checkpoint from the 2K token limit to 4K. For comparison, we also trained a model directly with 4K token limit using AdapThink. Our results showed that progressive training from 2K to 4K token limit outperformed direct 4K training, achieving the highest accuracy while maintaining the most efficient reasoning patterns.

## 7 Discussion

### 7.1 How Do Different Transition Words and Preferences Influence the Model's Performance?

As shown in Figure 4, we conducted a comparative analysis of three AdapThink variants by implementing control over different reflection words: AdapThink(check) calculates $\lambda_p$ in Equation 3 using the number of "Pause-Validation" words $n_p$, AdapThink(switch) controls "Branch-Extension" words $n_b$, and AdapThink(all) regulates both types simultaneously. Our experimental results demonstrate that reducing $n_p$ significantly influences learning efficiency, highlighting the crucial role of self-

Table 2: Ablation studies of AdapThink. $L$ represents the maximum response token limit (in thousand) in training phase. The diversity metrics are from training phase while other metrics are from testing on AIME 2025. Superscripts $\mathcal{T}$ and $\mathcal{F}$ indicate correct and incorrect answers respectively.

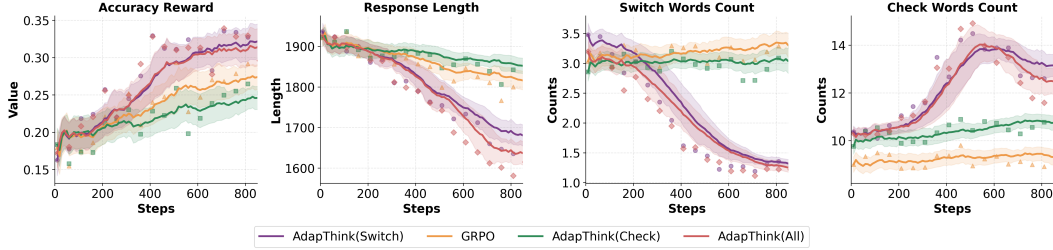| | | Setting | | | | Accuracy | | | Efficiency | | | | Diversity | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\lambda_l$ | $\lambda_o$ | $\lambda_b$ | $K$ | $L$ | PASS@1 | $\overline{L}^{\mathcal{T}}$ | $\overline{L}^{\mathcal{F}}$ | $n_p^{\mathcal{T}}$ | $n_p^{\mathcal{F}}$ | $n_b^{\mathcal{T}}$ | $n_b^{\mathcal{F}}$ | $H_l$ | $H_b$ | $H_d$ |
| ✓ | ✓ | ✓ | 2 | 2 | 25.4 ↑7.5 | 3561 ↓580 | 7442 ↓478 | 16.5 ↓4.6 | **55.7** ↓9.4 | 2.0 ↓1.4 | 13.3 ↑1.1 | 25.7 ↑3.2 | 28.6 ↓5.3 | 35.5 ↑1.7 |
| ✗ | ✓ | ✓ | 2 | 2 | 23.3 ↑5.4 | 4772 ↑630 | 8154 ↑235 | 56.9 ↑35.8 | 152.2 ↑87.1 | 3.4 ↓0.0 | **8.3** ↓3.9 | 22.2 ↓0.3 | 25.3 ↓8.6 | **38.9** ↑5.1 |
| ✓ | ✗ | ✓ | 2 | 2 | 22.5 ↑4.6 | 3775 ↓366 | 7853 ↑66 | 18.9 ↓2.2 | 71.4 ↑6.3 | 2.2 ↓1.2 | 12.2 ↓0.0 | 24.2 ↑1.7 | 27.8 ↓6.2 | 33.5 ↓0.4 |
| ✓ | ✓ | ✗ | 2 | 2 | 22.9 ↑5.0 | 4395 ↓253 | 7872 ↑48 | 27.3 ↑6.2 | 70.7 ↑5.6 | 3.1 ↑0.3 | 11.1 ↓1.1 | 23.5 ↑1.0 | 24.8 ↓9.2 | 33.8 ↓0.1 |
| ✓ | ✓ | ✓ | 1 | 2 | 22.5 ↑4.6 | 3947 ↓195 | 7793 ↑126 | 20.3 ↓0.8 | 71.1 ↑6.0 | 1.9 ↓1.5 | 10.9 ↓1.3 | 21.9 ↓0.6 | 26.5 ↓7.5 | 32.6 ↓1.2 |
| ✓ | ✓ | ✓ | 1.75 | 2 | 25.0 ↑7.1 | 3803 ↓339 | 7579 ↓340 | 20.1 ↓1.0 | 58.8 ↓6.3 | 2.3 ↓1.1 | 12.6 ↑0.4 | 23.2 ↑0.8 | 32.6 ↑1.4 | 34.3 ↑0.5 |
| ✓ | ✓ | ✓ | 2 | 4 | 24.6 ↑6.7 | 3861 ↓281 | 7909 ↑10 | 23.9 ↑2.8 | 78.1 ↑13.0 | 3.4 ↓0.0 | 14.1 ↑1.9 | 24.5 ↑2.1 | 27.9 ↓6.1 | 26.8 ↓7.0 |
| ✓ | ✓ | ✓ | 2 | 2→4 | **26.0** ↑8.1 | **3339** ↓802 | **6886** ↓1034 | **13.8** ↓7.3 | 63.2 ↑1.9 | **1.7** ↓1.7 | 13.9 ↑1.7 | **31.6** ↑9.1 | **32.7** ↓1.3 | 27.9 ↓5.9 |



Figure 4: Training comparison of AdapThink variants.

reflection processes even within short CoT reasoning. Additionally, AdapThink(all) exhibits a notable reduction in check words than AdapThink(switch), resulting in slightly degraded performance.

## 7.2 Does Length-control Mechanisms Exhibit Reward Hacking Behavior?

To further investigate potential reward hacking introduced by length-control mechanisms, we employed N-gram (N=40) repetition rate metrics [28] to quantify repetitive patterns in model responses for the MATH-500 test dataset, as shown in Table 3. Our analysis reveals that LCPO exhibits severe repetitive patterns with an average N-gram repetition rate of 10.8%, especially 19.8% in incorrect answers, indicating potential reward hacking in its length control mechanism. While TLB and CosFn show moderate repetition, GRPO achieves the lowest overall repetition rate of 0.3% as it avoids introducing dense rewards for length control. Notably, AdapThink maintains consistently

Table 3: N-gram repetition analysis across different post-training methods.

| Algorithm | N-grams (%) | | |
|---|---|---|---|
| | **Total** | $\mathcal{G}_{\mathcal{T}}$ | $\mathcal{G}_{\mathcal{F}}$ |
| GRPO | **0.3** | **0.3** | **0.2** |
| TLB | 3.0 | 3.5 | 2.2 |
| LCPO | 10.8 | 5.9 | 19.8 |
| CosFn | 1.4 | 2.0 | **0.2** |
| AdapThink | 0.7 | 0.6 | 0.9 |

low repetition rates across both correct and incorrect responses, suggesting its effectiveness in promoting rational reasoning patterns.

## 8 Conclusion

In this paper, we presented AdapThink, a novel RL post-training framework that adaptively modulates reasoning preferences based on model capabilities. We highlight two key insights from our study: (1) the effectiveness of dynamic reasoning control through group-relative rewards based on reflection words, completion markers and length, and (2) the importance of maintaining diversity in reasoning patterns during GRPR-based training. Notably, with only a 2K token limit during training, our method outperformed existing baselines in both mathematical accuracy and token efficiency across various mathematical reasoning benchmarks. One limitation and potential future direction lies in exploring semantic-level thought transitions control beyond current word-level control. We will further conduct experiments investigating the impact of non-reflection-related vocabulary on the reasoning process.

# References

[1] Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.

[2] Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, et al. Do not think that much for 2+ 3=? on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187*, 2024.

[3] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[4] Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. Token-budget-aware llm reasoning. *arXiv preprint arXiv:2412.18547*, 2024.

[5] Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang, and Mengnan Du. The impact of reasoning step length on large language models. In *ACL (Findings)*, 2024.

[6] Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. C3ot: Generating shorter chain-of-thought without compromising effectiveness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 24312–24320, 2025.

[7] Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy, Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment. *arXiv preprint arXiv:2410.01679*, 2024.

[8] Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip HS Torr, Fahad Shahbaz Khan, and Salman Khan. Llm post-training: A deep dive into reasoning large language models. *arXiv preprint arXiv:2502.21321*, 2025.

[9] Yuliang Liu, Junjie Lu, Zhaoling Chen, Chaofeng Qu, Jason Klein Liu, Chonghan Liu, Zefan Cai, Yunhui Xia, Li Zhao, Jiang Bian, et al. Adaptivestep: Automatically dividing reasoning step through model confidence. *arXiv preprint arXiv:2502.13943*, 2025.

[10] Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. Notion Blog.

[11] Yan Ma, Steffi Chern, Xuyang Shen, Yiran Zhong, and Pengfei Liu. Rethinking rl scaling for vision language models: A transparent, from-scratch framework and comprehensive evaluation scheme. *arXiv preprint arXiv:2504.02587*, 2025.

[12] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. `https://github.com/huggingface/peft`, 2022.

[13] Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. *arXiv preprint arXiv:2308.00436*, 2023.

[14] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

[15] Sania Nayab, Giulio Rossolini, Marco Simoni, Andrea Saracino, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. Concise thoughts: Impact of output length on llm reasoning and cost. *arXiv preprint arXiv:2407.19825*, 2024.

[16] OpenAI. Learning to reason with LLMs. `https://openai.com/index/learning-to-reason-with-llms`, 2024. Accessed: 2024.

[17] Matthew Renze and Erhan Guven. The benefits of a concise chain of thought on problem-solving in large language models. In *2024 2nd International Conference on Foundation and Large Language Models (FLLM)*, pages 476–483. IEEE, 2024.

[18] Matthew Renze and Erhan Guven. Self-reflection in llm agents: Effects on problem-solving performance. *arXiv preprint arXiv:2405.06682*, 2024.

[19] William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*, 2022.

[20] Maohao Shen, Guangtao Zeng, Zhenting Qi, Zhang-Wei Hong, Zhenfang Chen, Wei Lu, Gregory Wornell, Subhro Das, David Cox, and Chuang Gan. Satori: Reinforcement learning with chain-of-action-thought enhances llm reasoning via autoregressive search. *arXiv preprint arXiv:2502.02508*, 2025.

[21] Yi Shen, Jian Zhang, Jieyun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai Wang, and Shiguo Lian. Dast: Difficulty-adaptive slow-thinking for large reasoning models. *arXiv preprint arXiv:2503.04472*, 2025.

[22] Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, et al. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.

[23] Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, et al. Thoughts are all over the place: On the underthinking of o1-like llms. *arXiv preprint arXiv:2501.18585*, 2025.

[24] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2550–2575, 2023.

[25] Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning. *arXiv preprint arXiv:2502.14768*, 2025.

[26] Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, et al. A minimalist approach to llm reasoning: from rejection sampling to reinforce. *arXiv preprint arXiv:2504.11343*, 2025.

[27] Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. Chain of draft: Thinking faster by writing less. *arXiv preprint arXiv:2502.18600*, 2025.

[28] Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.

[29] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

[30] Xiaojiang Zhang, Jinghui Wang, Zifei Cheng, Wenhao Zhuang, Zheng Lin, Minglei Zhang, Shaojie Wang, Yinghan Cui, Chao Wang, Junyi Peng, et al. Srpo: A cross-domain implementation of large-scale reinforcement learning on llm. *arXiv preprint arXiv:2504.14286*, 2025.

# A Examples of Model's Thinking Process

## A.1 Example 1: Overthinking on a Simple Problem

In this example, the model was presented with a straightforward probability problem. While it reached the correct answer early in its reasoning process (using only 479 tokens), it continued to reflect unnecessarily. The model's self-validation tendencies, marked by frequent use of phrases like "Wait", "check", and "verify", led to multiple alternative approaches and verifications. This resulted in consuming 2,613 tokens—more than five times the tokens needed for the initial correct solution.

---

**Overthinking Example**

**Question**

A bag contains 4 red balls and 3 blue balls. Two balls are drawn without replacement. What is the probability that both balls are of different colors?

**Model Response**

> **Thinking Process**
>
> Okay, let's see. Well, I can think ... **Alternatively**, maybe I should ... **Wait**, is it right? ... **Alternatively**, maybe it is easier to ... So, the probability is 4/7. But ... **Alternatively** ... But let me **check** both approaches ... **Alternatively** ... Approach 1 ... Wait ... Let me **check** ... Approach 2 ... **Alternatively** ... **Wait**, but let me just **verify** ... So all methods lead to 4/7.
> **Final Answer**:
>
> The probability that both balls are of different colors is $\boxed{\dfrac{4}{7}}$.

All methods confirm the probability that both balls are of different colors is $\boxed{\dfrac{4}{7}}$.

**Correctness**

Correct

**Tokens Consumption**

2613 (479 for first generating the correct answer.)

**Reflection Words Count**

19

---

## A.2 Example 2: Underthinking on a Complex Problem

The second example shows the model attempting a more challenging modular arithmetic problem. Here, the model exhibited frequent but unproductive shifts in thinking, marked by numerous "Wait"

and "Alternatively" statements. Unlike the first example where excessive reflection was merely inefficient, it reached the maximum token limit (4,096 token) before completing its calculation, resulting in an incorrect answer.

---

**Underthinking Example**

**Question**

Let $a_n = 6^n + 8^n$. Determine the remainder upon dividing $a_{83}$ by $49$?

**Model Response**

*Thinking Process*

Okay, so I need to find the remainder ... Let me think ... First ... But **wait** ... **Alternatively** ... **Alternatively** ... **Wait** ... Let's **check**... **Wait**, perhaps there's a better way ... Oh! **Wait** a second ... **Wait** ... Okay, so that's the first part ... **Wait** ... **Wait** first ... **Wait**, perhaps ... **Wait**, but just to make sure ... Similarly, $8^7 \equiv 1 \pmod{13}$, 83 divided by 7 gives 11*7=

**Correctness**

Incorrect

**Tokens Consumption**

4096 (max token limit)

**Reflection Words Count**

35

---

## B  Distribution of Reflection Words in Model Responses

To provide a comprehensive understanding of our choice of reflection words in Section 3, we conducted a detailed frequency analysis of various transition words in Deepseek-distilled Qwen 1.5B model's responses. Figure 5 presents the complete distribution of both "Pause-Validation" and "Branch-Extension" words.

## C  AdapThink Algorithm Details

As introduced in Section 4, AdapThink proposes a novel post-training framework that adaptively controls model's reasoning behavior through group-level response characteristics. Here we provide the complete algorithm details of AdapThink (Algorithm 1).
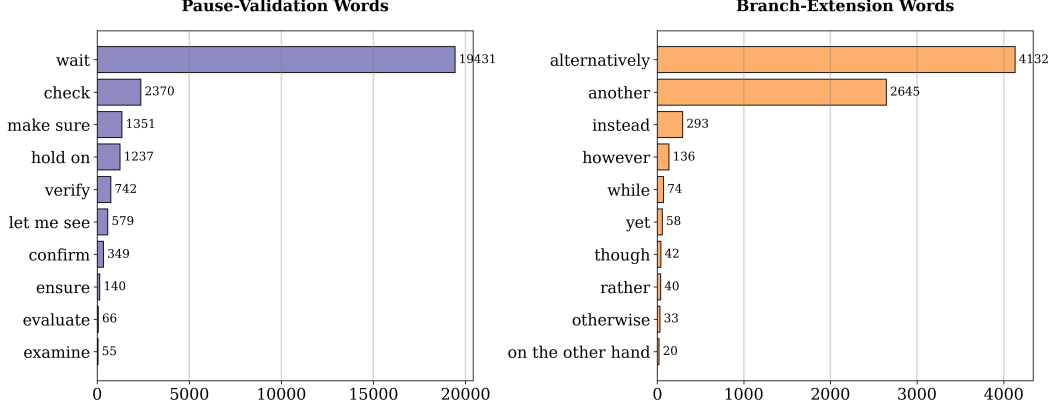
Figure 5: Frequency distribution of reflection words in base model's responses.

---

**Algorithm 1** AdapThink: Adaptive Thinking Post-trainig Framework

---

**Require:** Pre-trained model $\pi_\theta$, Dataset $\mathcal{D}$, Oversample factor $K$; Confidence thresholds $\varphi_{\text{low}}$, $\varphi_{\text{high}}$;Minimum correct/incorrect samples $T_{\min}$, $F_{\min}$
1: **for** each $(x, y) \in \mathcal{D}$ **do**
2:     Generate $K|\mathcal{G}|$ samples $\mathcal{G}'$ using $\pi_\theta$
3:     Compute model confidence $\varphi$ using Equation 1
4:     Partition $\mathcal{G}'$ into $\mathcal{G}_\mathcal{T}'$ and $\mathcal{G}_\mathcal{F}'$ based on correctness
5:     $T \leftarrow \min\{|\mathcal{G}_\mathcal{T}'|, T_{\min}\}$, $F \leftarrow \min\{|\mathcal{G}_\mathcal{F}'|, F_{\min}\}$
6:     **if** $\varphi \leq \varphi_{\text{low}}$ **then**
7:         Select $T$ correct samples from $\mathcal{G}_\mathcal{T}'$ maximizing $H_{\text{tot}}$
8:         Select remaining samples from $\mathcal{G}' \setminus \mathcal{G}_\mathcal{T}'$ maximizing $H_{\text{tot}}$
9:     **else**
10:         Select $F$ incorrect samples from $\mathcal{G}_\mathcal{F}'$ maximizing $H_{\text{tot}}$
11:         Select remaining samples from $\mathcal{G}' \setminus \mathcal{G}_\mathcal{F}'$ maximizing $H_{\text{tot}}$
12:     **end if**
13:     Compute component rewards $\lambda_l, \lambda_b, \lambda_o$ using Equation 2
14:     Calculate adaptive weight $\omega(\varphi)$ using Equation 4
15:     Update $\pi_\theta$ using Equation 9 with GRPR reward (Equation 3) and accuracy reward
16: **end for**
17: **return** Updated model $\pi_\theta$

---

# D   Hyperparameter Configuration

The core hyperparameters used in AdapThink are summarized in Table 4. The training is on 8 NVIDIA H100 GPUs, taking approximately 22 hours to complete 1000 steps.

# E   More Results

## E.1   Analysis of Detailed Reflection Words Distribution

To provide deeper insights into *how different post-training methods affect the model's reasoning behavior*, we analyze the distribution of reflection words across correct and incorrect responses, as shown in Figure 6.

For most pause-validation words and branch-extension words, we observe that AdapThink significantly reduces the frequency of validation words in both correct and incorrect responses compared to the base model, suggesting a more efficient reasoning process. In contrast, GRPO and TLB exhibit less pronounced changes in these reflection words; LCPO shows elevated counts of pause words, particularly in incorrect responses.

14

Figure 6: Distribution of reflection words across different post-training methods for correct and incorrect responses in MATH-500 datasets.

| Category | Parameter | Value |
|---|---|---|
| Model Settings | Base Model | DeepSeek-R1-Distill-Qwen-1.5B |
| | Torch Dtype | bfloat16 |
| | Max Sequence Length | 1024 |
| | Max New Tokens | 2048 |
| Training Settings | Learning Rate | 2e-6 |
| | Number of Epochs | 5 |
| | Batch Size | 8 |
| | Gradient Accumulation Steps | 8 |
| | Warmup Steps | 20 |
| LoRA Settings | LoRA Rank (r) | 32 |
| | LoRA Alpha | 32 |
| RL Settings | Algorithm | GRPO |
| | KL Coefficient | 0.15 |
| | Number of Generations | 16 |
| | Target Generations | 8 |
| | Min Correct Generations | 3 |
| | Min Incorrect Generations | 1 |
| Generation Settings | Temperature | 0.7 |
| | Top-p | 0.95 |

Table 4: Core hyperparameters for AdapThink post-training configuration

Our analysis provides insights into the relationship between reflection words and model performance. While LCPO achieves the shortest response lengths in Table. 1, its elevated use of pause words, especially in incorrect responses, suggests redundant reasoning patterns that may explain its relatively lower accuracy improvements. In contrast, AdapThink's systematic reduction in reflection words contributes to more efficient token usage while maintaining effective reasoning paths, resulting in consistent accuracy improvements across all benchmarks within token constraints.

## E.2 Analysis of More Ablation

We conducted additional ablation studies to explore two variants:

- *AdapThink(Sequential)* - where we calculate $\lambda_b$ based on the number of sequential words ("first","then", "next", "finally", "therefore", "so", "thus") within each group.
- *AdapThink(No weight)* - where we remove the weighting mechanism for rewards, consistently favoring shorter responses with minimal reflection.

We present comparative results between these variants and the original AdapThink model, both during training comparison in Figure 7 and testing on the AIME2025 benchmark in Figure 7.
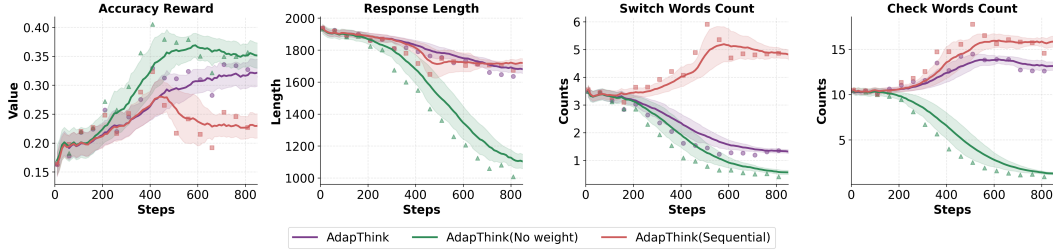


Figure 7: Training dynamics comparison among three AdapThink variants.
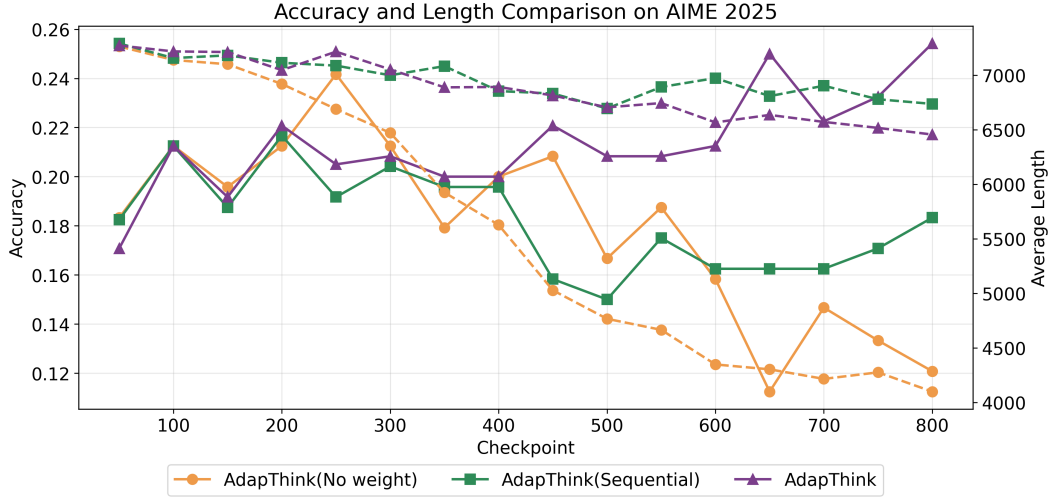


Figure 8: Test comparison among three AdapThink variants across checkpoints at 50-step intervals.

Interestingly, we observed two key findings:

(1) AdapThink(No weight) significantly reduced response lengths, which initially appeared advantageous in the early stages of post-training due to the inherent redundancy of reflection words in the base model. However, this eventually led to decreased accuracy rewards on the training set, with even more pronounced deterioration on the more challenging AIME 2025 dataset. As shown in Figure 8, we found that the reduced length in later stages of AdapThink(No weight) corresponded to lower PASS@1, highlighting a limitation of existing length control algorithms. Adapthink's dynamic weighting mechanism based on model capability mitigates this issue.

(2) AdapThink(Sequential) showed no significant changes in either performance or length. Our observations during both training and testing revealed that it consistently produced longer outputs compared to the original AdapThink algorithm, suggesting that controlling sequential words does not effectively address the overthinking/underthinking phenomenon.

## E.3 Example of AdapThink and Base Model

---

**AdapThink Example (from AIME 2025)**

**Question**

Find the number of ordered pairs $(x, y)$, where both $x$ and $y$ are integers between $-100$ and $100$, inclusive, such that $12x^2 - xy - 6y^2 = 0$.

**Model Response**

> **Thinking Process**
>
> Okay, so I need to find the number of ordered pairs $(x, y)$ where both $x$ and $y$ are integers between -100 and 100, inclusive. The condition given is $12x^2 - xy - 6y^2 = 0$. Hmm, let me try to figure this out step by step.
> First, looking at the equation: $12x^2 - xy - 6y^2 = 0$. It seems like a quadratic in terms of x and y. Maybe I can factor it or rearrange it somehow. ... Let's see:
>
> $$12x^2 - xy - 6y^2 = 0$$
>
> Divide both sides by 3:
>
> $$4x^2 - \frac{xy}{3} - 2y^2 = 0$$
>
> **Wait**, that doesn't help much because we have a fraction. Maybe factoring differently. ...
> So, seems consistent. Therefore, the number of ordered pairs is 117. **Final Answer**: The number of ordered pairs is $\boxed{117}$.

To find the number of ordered pairs $(x, y)$ where both $x$ and $y$ are integers between -100 and 100, inclusive, such that $12x^2 - xy - 6y^2 = 0$, we start by ... Thus, we subtract 1 to avoid double-counting. The total number of ordered pairs is:

$$51 + 67 - 1 = 117$$

Thus, the number of ordered pairs is $\boxed{117}$.

---

**Correctness**

Correct

---

**Tokens Consumption**

**4502**

---

**Reflection Words Count**

**19** (**1** "Branch-Extension" word and **18** "Pause-Validation" Words)

---

**DeepSeek-distilled Qwen 1.5B Example (from AIME 2025)**

### Question

Find the number of ordered pairs $(x, y)$, where both $x$ and $y$ are integers between $-100$ and $100$, inclusive, such that $12x^2 - xy - 6y^2 = 0$.

### Model Response

**Thinking Process**

Alright, so I have this problem here: I need to find the number of ordered pairs $(x, y)$ where both $x$ and $y$ are integers between $-100$ and $100$, inclusive. The equation given is $12x^2 - xy - 6y^2 = 0$. Hmm, okay. Let me try to figure this out step by step.

First, I notice that this is a quadratic equation in terms of x and y. Maybe I can factor it or rewrite it in a way that makes it easier to solve. Let me write it down again:

$$12x^2 - xy - 6y^2 = 0$$

**Alternatively**, maybe factor the quadratic expression as $(ax+by)(cx+dy) = 0$. Let me try that...
**Alternatively**, perhaps I can write the equation as: $12x^2 - xy - 6y^2 = 0$...
**Alternatively**, maybe factor it as $(ax + by)(cx + dy) = 0$...
**Alternatively**, perhaps factor out a 3:

$$3(4x^2 - \frac{1}{3}xy - 2y^2) = 0$$

But that doesn't seem helpful.
**Alternatively**, maybe factor as $(4x + ay)(3x + by) = 0$ ...
So, I think I can confidently say the answer is 117. **Final Answer**: $\boxed{117}$.

We need to find the number of ordered pairs $(x, y)$ where both $x$ and $y$ are integers between $-100$ and $100$, inclusive, such that $12x^2 - xy - 6y^2 = 0$ ... The total number of solutions is:

$$67 + 51 - 1 = 117$$

Thus, the number of ordered pairs $(x, y)$ is: $\boxed{117}$,

### Correctness

Correct

### Tokens Consumption

**7377**

### Reflection Words Count

**46** (**10** "Branch-Extension" words and **36** "Pause-Validation" Words)