

Thinking-Based Non-Thinking: Solving the Reward Hacking Problem in Training Hybrid Reasoning Models via Reinforcement Learning

Siyuan Gan¹, Jiaheng Liu¹, Boyan Wang¹, Tianpei Yang¹, Runqing Miao³, Yuyao Zhang³, Fanyu Meng³, Junlan Feng³, Linjian Meng^{†2}, Jing Huo^{†1}, Yang Gao¹

¹State Key Laboratory of Novel Software Technology, Nanjing University, Nanjing, China

²Shanghai Artificial Intelligence Laboratory, Shanghai, China

³Jiutian Research, Beijing, China

{gansiyuan, menglinjian}@smail.nju.edu.cn

{liujiaheng, boyanwang, tianpei.yang, huojing, gaoy}@nju.edu.cn

{miaorunqing, zhangyuyao, mengfanyu, fengjunlan}@chinamobile.com

Abstract

Large reasoning models (LRMs) have attracted much attention due to their exceptional performance. However, their performance mainly stems from thinking, a long Chain of Thought (CoT), which significantly increase computational overhead. To address this overthinking problem, existing work focuses on using reinforcement learning (RL) to train hybrid reasoning models that automatically decide whether to engage in thinking or not based on the complexity of the query. Unfortunately, using RL will suffer the the reward hacking problem, e.g., the model engages in thinking but is judged as not doing so, resulting in incorrect rewards. To mitigate this problem, existing works either employ supervised fine-tuning (SFT), which incurs high computational costs, or enforce uniform token limits on non-thinking responses, which yields limited mitigation of the problem. In this paper, we propose *Thinking-Based Non-Thinking* (TNT). It does not employ SFT, and sets different maximum token usage for responses not using thinking across various queries by leveraging information from the solution component of the responses using thinking. Experiments on five mathematical benchmarks demonstrate that TNT reduces token usage by around 50% compared to DeepSeek-R1-Distill-Qwen-1.5B/7B and DeepScaleR-1.5B, while significantly improving accuracy. In fact, TNT achieves the optimal trade-off between accuracy and efficiency among all tested methods. Additionally, the probability of reward hacking problem in TNT’s responses, which are classified as not using thinking, remains below 10% across all tested datasets.

1 Introduction

Large reasoning models (LRMs) (Xu et al., 2025), such as DeepSeek-R1 (Guo et al., 2025) and OpenAI o1 (Jaech et al., 2024), have recently become

a central research topic due to their superior capabilities. However, their performance gains rely on extended chain-of-thought (CoT) (Wei et al., 2022)—also referred to as *thinking*. This reliance leads to protracted and repetitive outputs, resulting in a substantial increase in inference overhead and latency (Li et al., 2025; Chen et al., 2024), a phenomenon referred to as the overthinking problem (Sui et al., 2025; Qu et al., 2025).

To tackle the overthinking problem, recent works increasingly explore hybrid reasoning models that can dynamically determine whether to activate the thinking mode (Zhang et al., 2025; Fang et al., 2025; Tu et al., 2025; Lou et al., 2025; Jiang et al., 2025; Chen et al., 2025; Zhan et al., 2025). The alternative, where the model directly outputs solution without thinking, is referred to as the non-thinking mode. Among these works, one of the most commonly used method involves employing reinforcement learning (RL), with a higher reward allocated to correct answer of the non-thinking mode than that of the thinking mode (Zhang et al., 2025; Fang et al., 2025; Tu et al., 2025). Unfortunately, the RL training method suffer from reward hacking problem that reduces the performance. Specially, models embed a significant portion of its reasoning process within the non-thinking response (Zhang et al., 2025; Tu et al., 2025). By doing so, the model leverages the thinking process to arrive at the correct answer while securing additional rewards for following the non-thinking format. This leads to increased rewards, but results in a significant increase in tokens for non-thinking mode responses.

One approach to mitigate this reward hacking problem is to use supervised fine-tuning (SFT) with a dataset that is significantly larger than the RL dataset, thereby fixing the model’s output, as did in Thinkless (Fang et al., 2025). However, this also incurs a substantial increase in computational overhead. To avoid such high computational overhead, AdaptThink (Zhang et al., 2025) set a smaller max-

† Corresponding author.

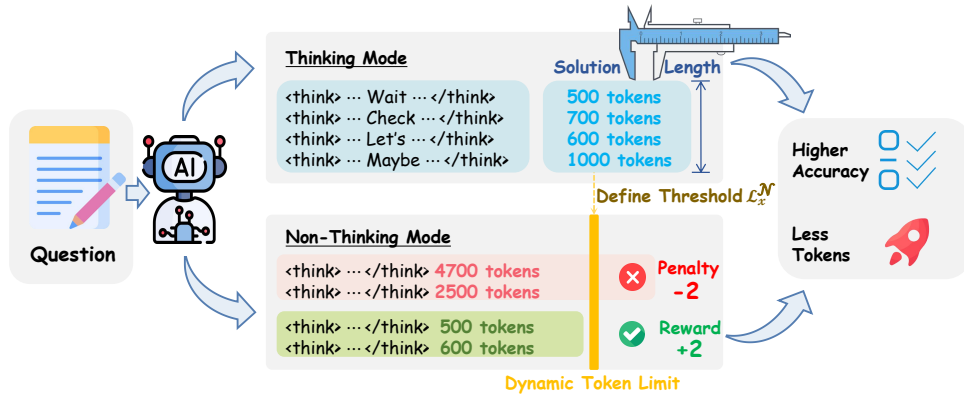


Figure 1: Overview of TNT.

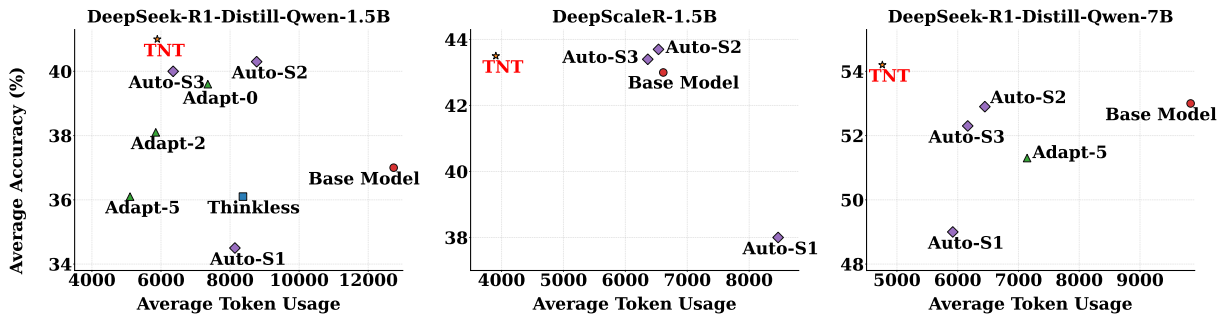


Figure 2: Average accuracy and token usage comparison across different hybrid reasoning model training methods on mathematical benchmarks. We only presented the evaluation results of their open-source checkpoints while some of these methods lack the trained checkpoints based on DeepScaleR-1.5B, and DeepSeek-R1-Distill-Qwen-7B. Adapt- x refers to AdaptThink with $\delta = x * 0.01$. Auto- S_x refers to AutoThink-Stage x . We also follow these abbreviation conventions in the table below.

imum token usage for the non-thinking mode than that for the thinking mode. This way, the response judged in the non-thinking mode cannot increase token usage to actually operate in thinking mode and obtain a higher reward. Unfortunately, Zhang et al. (2025) set the maximum token usage for the non-thinking mode to be uniform across queries, resulting in only limited mitigation of reward hacking problem. This is because applying such a uniform maximum may result in poor detection of reward hacking problem because thinking mode tokens for simple queries can be fewer than non-thinking mode tokens for complex ones, making uniform limits ineffective (see details in Section A.2).

To set the maximum token usage across different queries, we propose *Thinking-Based Non-Thinking* (TNT). As demonstrated in Figure 1, TNT employs the solution component of the thinking mode responses (i.e., the tokens following `</think>` are highlighted in blue as shown in Figure 1) to set the maximum token usage for the non-thinking mode. Specifically, LRMs’ thinking mode has been trained to ensure the solution component of the thinking mode responses does not involve thinking, such as DeepSeek-R1. Therefore, this component will not differ significantly from the output of the

non-thinking mode as the definition of the latter is that outputting solution without thinking. We can use the token usage of this component to set the maximum token usage for the non-thinking mode.

Consistent with previous hybrid reasoning model training methods (Zhang et al., 2025; Tu et al., 2025; Fang et al., 2025), we employ DeepSeek-R1-Distill-Qwen-1.5B/7B and DeepScaleR-1.5B as base models to evaluate TNT on five standard mathematical benchmarks: AIME24, AIME25, Minerva, AMC23, and Olympiad. The results show that TNT consistently reduces token usage by approximately 50% relative to base models, while substantially improving accuracy. Compared to other hybrid reasoning model training methods, TNT achieves the optimal balance between accuracy and token usage as Figure 2. Furthermore, the probability of reward hacking problem in TNT remains below 10% across all tested datasets.

2 Preliminary

Now, we introduce the LRMs. Consider a LRM parameterized by the parameter θ and denoted by π_θ . Given a prompt $x = [x_1, \dots, x_n, \langle \text{think} \rangle]$, where $[x_1, \dots, x_n]$ represents the query and the

special token `<think>` means the start of thinking. For this prompt, this LRM generates a response $y = [y_1, \dots, y_\tau, \text{</think>}, y_{\tau+2}, \dots, y_m]$. Here, $[y_1, \dots, y_\tau]$ corresponds to thinking, which is a long CoT consisting of constant exploration, reflection, and self-verification. The token `</think>` marks the end of thinking. The remaining sequence, $[y_{\tau+2}, \dots, y_m]$, denotes the final solution, which only includes the correct steps to solve the problem and the final answer. From the perspective of probability theory, the response y is a sample drawn from the conditional probability distribution $\pi_\theta(\cdot|x)$. Since y is generated in an auto-regressive way, the conditional probability $\pi_\theta(y|x)$ is:

$$\pi_\theta(y|x) = \prod_{t=1}^m \pi_\theta(y_t|x, y_{<t}). \quad (1)$$

In this paper, we follow the approach in Zhang et al. (2025); Tu et al. (2025); Lou et al. (2025)¹, referring to $[y_1, \dots, y_\tau] \neq \emptyset$ as the thinking mode and $[y_1, \dots, y_\tau] = \emptyset$ as the non-thinking mode. We use this approach to distinguish between thinking and non-thinking modes because it does not require modifications to the tokenizer, whereas the methods used in other works do require modifications to the tokenizer (Fang et al., 2025; Jiang et al., 2025; Zhan et al., 2025). For example, Fang et al. (2025) determine whether the response is thinking mode or non-thinking mode based on whether the first token of the response y is `<short>`. Unfortunately, `<short>` is usually not in the tokenizer of LRMs, so we need to modify the tokenizer.

3 Motivation

We now provide a detailed introduction to the reward hacking problem in training hybrid reasoning models by using RL with a higher reward allocated to the correct answer of the non-thinking mode compared to that of the thinking mode (Zhang et al., 2025; Fang et al., 2025; Tu et al., 2025).

As mentioned in Section 2, Zhang et al. (2025); Tu et al. (2025) and Fang et al. (2025) determine whether the response is thinking or non-thinking mode based on whether the first token in the response y is `</think>` or `<short>`, respectively. However, the model might generate `</think>` or `<short>` as the first token, indicating the non-thinking mode. Nonetheless, subsequent tokens conform to the thinking mode, producing a long CoT, as demon-

¹In Tu et al. (2025), they use the *ellipsis prompt* to make $x = [x_1, \dots, x_n, \text{<think>}, \backslash n, \dots, \backslash n]$, which enables sampling the non-thinking mode without using the off-policy sampling.

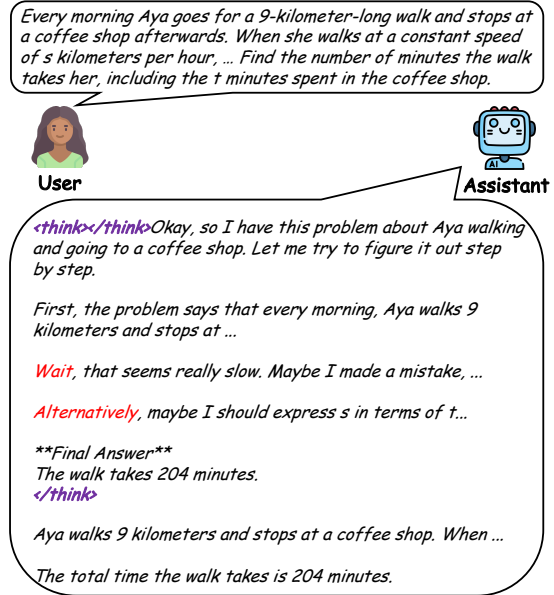


Figure 3: An example of the reward hacking problem occurs when Assistant uses the AutoThink model (Tu et al., 2025). The first generated token in the response is `</think>` (`<think>` is included in the input prompt x), indicating that the response is classified as the non-thinking mode. However, the response clearly demonstrates characteristics of the thinking mode, as evidenced by the usage of keywords like "Wait" and "Alternatively" along with the regeneration of the thinking's termination token `</think>`. The answer 204 is correct. This response is undeniably part of the thinking mode despite being incorrectly assigned a higher reward associated with the correct answer of the non-thinking mode. This discrepancy results in a clear instance of the reward hacking problem, where the reward allocation does not align with the true mode of the response.

strated in Figure 3. Specifically, a reward associated with the non-thinking mode for yielding the correct answer is allocated. However, it is the thinking mode that genuinely derives the correct answer, which is substantiated by the presence of keywords such as "Wait" and "Alternatively", as well as the regeneration of the thinking's termination token `</think>`. This is a classic reward hacking problem, where the reward is not allocated to the true mode.

Surprisingly, for models that do not consider the reward hacking problem, such as AutoThink (Tu et al., 2025), our evaluation of its Stage 1 model performance on AIME24 reveals that responses identified as non-thinking mode exhibit an extremely high average token usage. Specifically, the average token usage for non-thinking mode responses is 10845, while for thinking mode responses, it is 11976. This finding highlights that failing to address the reward hacking problem issue could lead to the collapse of the entire training process.

Algorithm 1 Thinking-Based Non-Thinking (TNT)

Input: policy model π_θ ; dataset \mathcal{D} ; hyperparameters $\omega, \mathcal{L}^\emptyset, \alpha, \beta$

- 1: **for** step = 1, . . . , M **do**
- 2: Sample a batch \mathcal{D}_b from \mathcal{D}
- 3: Sample K responses $\{y_x^k\}_{k=1}^K \sim \pi_\theta(\cdot|x)$ for each $x \in \mathcal{D}_b$
- 4: Split $\{y_x^k\}_{k=1}^K$ into $\mathcal{M}_x^T = \{y_x^j \mid p(y_x^j) = 1\}$ and $\mathcal{M}_x^N = \{y_x^j \mid p(y_x^j) = 0\}$ for each prompt x
- 5: Compute \mathcal{L}_x^N via \mathcal{M}_x^T (Equation (2))
- 6: Assign a reward for each response y_x^k of each prompt x via \mathcal{L}_x^N (Equation (5))
- 7: Update the policy model π_θ via GPRO and the assigned rewards
- 8: **end for**

Output: π_θ

4 Our Method

To mitigate the reward hacking problem that reduces the performance, there are two primary methods: adding SFT, which subsequently raises computational costs, or setting a smaller maximum token usage for the thinking mode than that of the non-thinking mode. However, existing research imposes a uniformly maximum token usage for the non-thinking mode across all queries without considering their varying levels of difficulty, resulting in only limited mitigation of the reward hacking problem (see details in Section A.2).

To set the maximum token usage of the non-thinking mode for different queries, we propose *Thinking-Based Non-Thinking* (TNT). The core insight of TNT is to utilize the solution component of the thinking mode responses to determine the maximum token usage for the non-thinking mode. Specifically, LLMs’ thinking mode is trained on large-scale datasets, ensuring that the solution component is devoid of any additional thinking. As a result, this solution component closely aligns with the output of the non-thinking mode, which is defined by generating solutions without thinking. Thus, we can appropriately utilize the token usage derived from this solution component to set the maximum token usage for the non-thinking mode.

4.1 Details of TNT

We now present the details of TNT. The pseudocode for TNT is provided in Algorithm 1.

Consider a reasoning model π_θ and a dataset \mathcal{D} . Let $(x, y_x^*) \in \mathcal{D}$ denote the input prompt and golden answer, respectively. We adopt the *ellipsis prompt* from Tu et al. (2025), defining x as $[x_1, \dots, x_n, \langle \text{think} \rangle, \backslash n, \dots, \backslash n]$, to enable sampling the non-thinking mode responses without the complex off-policy sampling. For each $(x, y_x^*) \in$

\mathcal{D} , we use $y_x^k \sim \pi_\theta(\cdot|x)$ to denote the k -th sampled response of the model π_θ given the prompt x . Also, we employ $p(y_x^k)$ to represent the discrimination function that determines whether the mode of the response y_x^k is thinking or non-thinking. This function outputs 1 if the mode is classified as thinking and 0 if it is classified as non-thinking. Following Tu et al. (2025); Zhang et al. (2025); Lou et al. (2025), y is classified as non-thinking mode if y starts with $\langle \text{think} \rangle$, otherwise as thinking mode.

The core of TNT is to determine the maximum token usage for the non-thinking mode by leveraging the solution component of responses in the thinking mode. Specifically, for each $(x, y_x^*) \in \mathcal{D}$, we sample K responses, denoted as $y_x^1, y_x^2, \dots, y_x^k, \dots, y_x^K$. Let $\mathcal{M}_x^T = \{y_x^j \mid p(y_x^j) = 1\}$ and $\mathcal{M}_x^N = \{y_x^j \mid p(y_x^j) = 0\}$, which correspond to sets of responses in the thinking mode and non-thinking mode given the prompt x , respectively. We denote \mathcal{L}_x^T and \mathcal{L}_x^N as the maximum token usages for the thinking and non-thinking modes given the prompt x . We focus on the computation of \mathcal{L}_x^N , as \mathcal{L}_x^T is equal to the training context length and remains unchanged across all prompt x .

For each $y_x^j \in \mathcal{M}_x^T$, we evaluate the number of tokens following $\langle \text{think} \rangle$ in the response $y_x^j = [y_{x,1}^j, \dots, y_{x,\tau}^j, \langle \text{think} \rangle, y_{x,\tau+2}^j, \dots, y_{x,m}^j]$, which is denoted as $h(y_x^j) = m - \tau - 1$. In other words, $h(y_x^j)$ is the length of the token sequence $[y_{x,\tau+2}^j, \dots, y_{x,m}^j]$. We then define \mathcal{L}_x^N as:

$$\mathcal{L}_x^N = \begin{cases} \omega \frac{\sum_{y_x^j \in \mathcal{M}_x^T} h(y_x^j)}{|\mathcal{M}_x^T|}, & \text{if } \mathcal{M}_x^T \neq \emptyset, \\ \mathcal{L}^\emptyset, & \text{if } \mathcal{M}_x^T = \emptyset, \end{cases} \quad (2)$$

where $\omega \geq 1$ and \mathcal{L}^\emptyset are positive constants. The role of ω in Equation (2) is crucial to reduce the risk of misidentifying reward hacking problem. Formally, we employ average value of $h(y_x^j)$, where $y_x^j \in \mathcal{M}_x^T$. If the token usage of a response from

\mathcal{M}_x^N exceeds the average value by precisely 1, it may be incorrectly identified as exhibiting reward hacking problem. To mitigate this risk, we incorporate a weight $\omega > 1$. This adjustment helps in distinguishing genuine anomalies from minor deviations, thereby addressing the potential misjudgment. In our experiments, we set $\omega = 2$. In addition, the parameter \mathcal{L}^θ in Equation (2) is introduced because we may not sample the thinking mode responses when we use the on-policy sampling. When this occurs, we are unable to calculate the expression $\omega \sum_{y_x^j \in \mathcal{M}^\tau} h(y_x^j) / |\mathcal{M}_x^\tau|$. Thus, the introduction of the parameter \mathcal{L}^θ is necessary. In our experiments, we typically set \mathcal{L}^θ to 1000.

After computing \mathcal{L}_x^N , we proceed to describe its application in constructing the reward function $R(x, y_x^k, y_x^*, p(y_x^k), \mathcal{L}_x^N)$. Let $r(y_x^k, y_x^*)$ denote the answer extraction function, which returns 0 if y_x^k is incorrect and 1 if y_x^k is correct. Firstly, if the response y_x^k is classified as the thinking mode, e.g., $p(y_x^k) = 1$, we use the following reward function:

$$R^T(x, y_x^k, y_x^*, \mathcal{L}_x^N) = \begin{cases} 1, & \text{if } r(y_x^k, y_x^*) = 1, \\ 0, & \text{if } r(y_x^k, y_x^*) = 0. \end{cases} \quad (3)$$

This function indicates that when the response y_x^k is determined to be in thinking mode, a correct answer returns a reward of 1, while an incorrect response yields 0. Secondly, if the response y_x^k is classified as the non-thinking mode, e.g., $p(y_x^k) = 0$, We use the following reward function that is modified from the one proposed by Tu et al. (2025):

$$R^N(x, y_x^k, y_x^*, \mathcal{L}_x^N) = \begin{cases} 2, & \text{if } r(y_x^k, y_x^*) = 1 \wedge |y_x^k| \leq \mathcal{L}_x^N, \\ -1, & \text{if } r(y_x^k, y_x^*) = 0 \wedge |y_x^k| \leq \mathcal{L}_x^N, \\ -2, & \text{if } |y_x^k| > \mathcal{L}_x^N, \end{cases} \quad (4)$$

where $|y_x^k|$ represents the length of y_x^k . This function implies that if $|y_x^k|$ is less than or equal to \mathcal{L}_x^N , indicating no reward hacking problem, we directly employ the naive reward function as presented in Tu et al. (2025). This incentivizes using the non-thinking mode when both non-thinking and thinking modes can yield correct responses, thereby allowing for automatic mode selection based on query complexity. Conversely, when $|y_x^k|$ exceeds \mathcal{L}_x^N , reward hacking problem is indicated. By setting the reward as -2 for this scenario, we ensure that both incorrect and correct answers are assigned lower rewards compared to scenarios without reward hacking problem, thereby minimizing the reward hacking problem behavior.

By integrating Equations (3) and (4), which correspond to the reward functions for thinking and non-thinking mode responses, respectively, the ultimate reward function can be expressed as follows:

$$R(x, y_x^k, y_x^*, p(y_x^k), \mathcal{L}_x^N) = \begin{cases} R^T(x, y_x^k, y_x^*, \mathcal{L}_x^N), & \text{if } p(y_x^k) = 1, \\ R^N(x, y_x^k, y_x^*, \mathcal{L}_x^N), & \text{if } p(y_x^k) = 0. \end{cases} \quad (5)$$

Finally, according to the dataset \mathcal{D} and the reward function $R(x, y_x^k, y_x^*, p(y_x^k), \mathcal{L}_x^N)$ defined in Equation (5), we employ the GRPO algorithm with a token-level policy gradient loss for training (Shao et al., 2024), as defined by:

$$\mathcal{J}(\theta) = \mathbb{E}_{(x, y_x^*) \sim \mathcal{D}, \{y_x^k\}_{k=1}^K \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\frac{1}{\sum_{k=1}^K |y_x^k|} \sum_{k=1}^K \sum_{t=1}^{|y_x^k|} \min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1-\varepsilon, 1+\varepsilon) \hat{A}_{i,t} \right) \right],$$

where $r_{i,t}(\theta)$ is the token-level importance weight defined as the ratio between the new and old token probabilities, and $\hat{A}_{i,t}$ represents the token-level advantage estimated through the reward function $R(x, y_x^k, y_x^*, p(y_x^k), \mathcal{L}_x^N)$ defined in Equation (5). The overall loss is normalized by the total number of tokens across all sampled trajectories.

4.2 Discussions of TNT

The core of TNT focuses on determining the maximum token usage for the non-thinking mode, which enables TNT to be compatible with RL algorithms. This implies that, while GRPO is used in the design of TNT, it can be replaced with other advanced GRPO variants, such as Dr. GRPO (Liu et al., 2025), DAPO (Yu et al., 2025), and GSPO (Zheng et al., 2025), or the classic PPO (Schulman et al., 2017). Also, as discussed in Section 4.1, TNT can also be effectively combined with the off-policy sampling, as did in Zhang et al. (2025) and Jiang et al. (2025). Additionally, techniques from other research on hybrid reasoning models can be applied. For instance, we can assign different weights to the control and response tokens in the loss function (Fang et al., 2025; Lou et al., 2025; Zhang et al., 2025). Furthermore, the two methods employed by Tu et al. (2025)—Batch-Level Reward Balancing and Length-Aware Reward—can also be utilized to respectively alleviate model collapse and further reduce token usage. Moreover, the reference model reward introduced in Zhang et al. (2025) can also be incorporated to provide additional signals during the training.

5 Experiments

5.1 Experimental Setups

Datasets and Models. Following previous hybrid reasoning model training methods (Fang et al., 2025; Zhang et al., 2025; Tu et al., 2025), we utilize DeepScaleR (Luo et al., 2025) as the training dataset, which consists of 40,000 mathematical problems with varying levels of difficulty. Following previous hybrid reasoning model training methods, we select DeepSeek-R1-Distill-Qwen-1.5B (Guo et al., 2025) as the base model in this section, given that it serves as the only base model shared across all previous hybrid reasoning model training methods. Besides, we show more details for the scenario when the base models are DeepScaleR-1.5B (Luo et al., 2025) or DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025) in Finding 4 and Section B.1. Our evaluation is performed on five standard mathematical benchmarks: AIME24, AIME25, Minerva, AMC23, and Olympiad. We do not conduct tests on the GSM8K (Cobbe et al., 2021) and MATH (Lightman et al., 2023), as both are too simplistic to adequately reflect the model’s performance. During evaluation, all models use a 32K context window. For evaluation metrics, we consider accuracy and response token usage.

Training Details. All experiments are conducted on a single node with 8 H800 GPUs. All experiments are implemented using the verl framework (Sheng et al., 2025). The batch size and training context length are set to 64 and 16K, respectively. We selected the checkpoint at step 600 when the base model is DeepSeek-R1-Distill-Qwen-1.5B, as it observed convergence and yielded the optimal performance. We set $\omega = 2$ and $\mathcal{L}^\theta = 1000$, as mentioned in Section 4.1.

Baselines. We mainly evaluate TNT against representative methods for hybrid reasoning model training method, such as Thinkless (Fang et al., 2025), AdaptThink (Zhang et al., 2025), and AutoThink (Tu et al., 2025). They are selected based on the public availability of their trained models, as well as their base models consistent with our experiments. Furthermore, regarding specific configurations: for Thinkless, we exclusively analyze its RL model due to the suboptimal performance of its SFT model; for AdaptThink, we evaluate configurations with $\delta \in \{0.0, 0.02, 0.05\}$, noting that $\delta = 0.05$ serves as the default; and for AutoThink, we assess performance across its three distinct training stages (Stage 1 to Stage 3).

5.2 Experimental Results

Finding 1. TNT enhances accuracy and simultaneously significantly reduces token usage.

Table 1 presents the accuracy and token usage of various models when the base model is DeepSeek-R1-Distill-Qwen-1.5B across the tested datasets. To evaluate the trade-off between accuracy and conciseness, we adopt the Token Efficiency (TE) metric from Hong et al. (2025), defined as A/\sqrt{L} , where A and L represent accuracy and token usage, respectively. The average accuracy and token usage across all datasets are used to compute this metric. Additionally, the token usage on each benchmark for both the thinking mode and non-thinking mode responses are shown in Table 2.

Firstly, compared to the base model, TNT reduces the average response token usage by 46.2% and enhances the average accuracy by 4.1%. Secondly, TNT surpasses Thinkless across all datasets in terms of both accuracy and token usage. On the AIME24 dataset, TNT utilizes approximately 80% of the tokens used by Thinkless while exceeding their accuracy by over 10%. Notably, TNT and Thinkless employed the same RL dataset—DeepScaleR. Thirdly, compared to AdaptThink, while TNT primarily focuses on enhancing accuracy, it surpasses three AdaptThink models trained with different δ by evaluating the trade-off between accuracy and conciseness through the TE. Lastly, TNT consistently outperforms the final output model of AutoThink (AutoThink-Stage 3) across all benchmarks in both accuracy and token usage. Notably, TNT achieves this without the CoT Compression method used in AutoThink. While incorporating CoT Compression methods may further reduce TNT’s token usage, we opt not to pursue this as our focus remains on the hybrid reasoning models rather than CoT Compression methods.

In fact, we also examine the variations in accuracy and token usage on both the training and testing datasets during the training process. See details in Appendix B.1.

Finding 2. TNT learns to autonomously select between the non-thinking mode and thinking mode based on the difficulty of the query.

We also investigate the relationship between thinking behavior and the inherent difficulty of tasks. In Table 3, we report the ratio of responses classified as belonging to the non-thinking mode across different models. By considering the re-

Models	Accuracy (%) \uparrow						Token Usage \downarrow						TE
	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	
Base Model	28.6	24.6	26.2	62.1	43.6	37.0	16865	16464	7490	11050	11808	12736	0.33
Thinkless	28.4	24.2	26.6	58.7	42.5	36.1	11207	10664	5463	7153	7379	8373	0.39
Adapt-0	31.9	25.0	27.3	68.2	45.4	39.6	10140	10039	4145	5661	6796	7356	0.46
Adapt-2	32.5	24.2	23.3	<u>66.4</u>	43.9	38.1	9087	8367	<u>2283</u>	4254	5252	<u>5849</u>	<u>0.50</u>
Adapt-5	29.5	22.9	25.9	60.8	41.2	36.1	8105	8038	1661	3475	4240	5104	<u>0.50</u>
Auto-S1	27.1	20.8	22.3	60.7	41.6	34.5	11701	11317	6553	<u>3822</u>	7308	8140	0.38
Auto-S2	<u>35.0</u>	26.3	<u>28.0</u>	65.1	47.0	<u>40.3</u>	11858	11186	5861	7178	7768	8770	0.43
Auto-S3	34.1	25.0	28.7	65.8	46.2	40.0	9742	8625	3172	4875	5345	6352	<u>0.50</u>
TNT (Ours)	37.7	<u>26.1</u>	28.7	65.9	<u>46.4</u>	41.0	<u>8537</u>	<u>8252</u>	3150	4508	<u>5020</u>	5893	0.53

Table 1: Comparison of accuracy, token usage, and TE on mathematical benchmarks across hybrid reasoning models when the base model is DeepSeek-R1-Distill-Qwen-1.5B. The best and second results are bolded and underlined, respectively. Gray represents the base model, cyan represents baseline hybrid reasoning models training methods, and orange represents TNT model. We will use the same color settings below as well.

Models	Non-Thinking Mode Tokens						Thinking Mode Tokens					
	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	AIME24	AIME25	Minerva	AMC23	Olym.	AVG
Base Model	–	–	–	–	–	–	16865	16464	7490	11050	11808	12736
Thinkless	–	–	372	367	467	402	11439	10670	5535	7510	8649	8761
Adapt-0	5468	2283	1059	3507	3912	3246	10536	11216	4448	6513	8759	8294
Adapt-2	6549	4724	1586	3034	4624	4103	11050	10258	3807	5957	8206	7856
Adapt-5	4076	2556	961	1499	2314	2281	11137	11041	5558	8514	9268	9104
Auto-S1	10845	8858	1893	3980	3806	5876	11976	12223	4212	10114	10033	9712
Auto-S2	7870	6507	4995	5024	5521	5983	12777	12066	6815	11060	11081	10760
Auto-S3	4537	1890	1559	2107	2759	2570	10655	8601	3729	6950	7130	7413
TNT (Ours)	995	795	601	859	937	837	8633	8325	3475	7159	6086	6736

Table 2: Comparison of non-thinking and thinking mode token usage on mathematical benchmarks across hybrid reasoning models when the base model is DeepSeek-R1-Distill-Qwen-1.5B.

Models	No-Thinking Mode Ratio (%)				
	AIME24	AIME25	Minerva	AMC23	Olym.
Base Model	0.0	0.0	0.0	0.0	0.0
Thinkless	0.0	0.0	1.4	5.0	2.8
Adapt-0	11.3	9.6	9.0	28.3	40.5
Adapt-2	42.9	35.4	68.6	58.3	71.3
Adapt-5	44.2	36.7	84.7	71.8	72.3
Auto-S1	22.8	25.4	16.8	58.1	37.3
Auto-S2	19.0	17.9	52.4	64.3	59.6
Auto-S3	13.3	10.2	29.0	47.6	46.8
TNT (Ours)	1.7	0.8	11.3	29.4	20.7

Table 3: Non-thinking mode ratio on mathematical benchmarks across hybrid reasoning models when the base model is DeepSeek-R1-Distill-Qwen-1.5B.

sults in Table 1, we observe a negative correlation between the non-thinking rate and task difficulty. Specifically, accuracy serves as a proxy for dataset difficulty. We find that on more challenging datasets, TNT tends to invoke explicit thinking more frequently. In addition, the results in Appendix B.1 (Figure 5) illustrate the variation of the non-thinking mode ratio on both the training and

testing datasets during the training process. We observe that while such ratio on the training set exhibits a U-shaped pattern similar to Thinkless, such ratio on the testing dataset consistently remains at a low level. This observation indicates that TNT does not exhibit overfitting.

Finding 3. TNT effectively mitigates the reward hacking problem.

We analyze the probability of thinking-related verbs such as "Wait," "Alternatively," and "Double-Check" appearing in responses classified under the non-thinking mode across different models, as shown in Figure 4. This probability reflects the likelihood of reward hacking problem in the model responses, as the presence of these words indicates actual thinking by the model. Our findings show that Thinkless, which utilizes SFT, exhibits the lowest probability, followed by our TNT model. Additionally, we find that the uniform maximum token usage for all queries in the non-thinking mode used in AdaptThink significantly increases the likelihood of reward hacking problem compared to TNT.

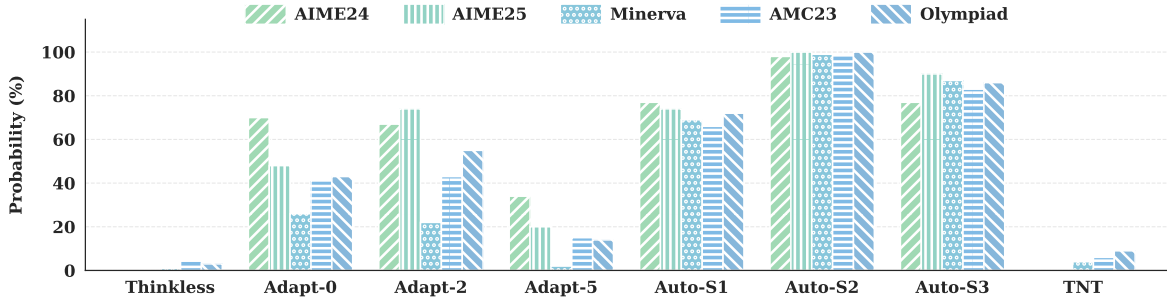


Figure 4: Probability of thinking-related verbs appearing in non-thinking mode responses across models with base model DeepSeek-R1-Distill-Qwen-1.5B on mathematical benchmarks.

Also, AutoThink, which does not consider the reward hacking problem, shows the highest probability. Furthermore, results in Table 2 corroborate that TNT effectively mitigates the reward hacking problem. Specifically, in terms of the token usage for responses in non-thinking mode, Thinkless has the lowest value, with TNT as the second lowest.

Although Thinkless addresses the reward hacking problem more effectively, it relies on SFT, which results in significant computational costs. Specifically, the SFT dataset used by Thinkless is approximately 2000K in size, making it about 50 times larger than the RL dataset, DeepScaleR, which is employed by Thinkless, AdaptThink, AutoThink, and ours. Moreover, constructing the SFT dataset incurs substantial computational overhead. Thinkless utilizes DeepSeek-R1-671B (Guo et al., 2025) to generate the SFT dataset, a process that is computationally intensive. In addition, SFT leads to significant performance degradation, which is why we do not compare against Thinkless’s SFT model. Precisely, the accuracy of Thinkless’s SFT model on AIME24 is merely around 10%!

Finding 4. TNT becomes more pronounced on better base models.

To verify the scalability of TNT, we conducted experiments on DeepScaleR-1.5B and DeepSeek-R1-Distill-Qwen-7B. Interestingly, TNT attains TE scores of 0.70 and 0.79, respectively, significantly surpassing the runners-up (0.54 and 0.67). Our results demonstrate that the advantage of TNT becomes more pronounced as the capability of the base model increases. See details in Appendix B.1.

Finding 5. TNT exhibits better performance compared with CoT compression methods.

Besides, we compared our accuracy and token usage against CoT compression methods (Hou et al., 2025; Arora and Zanette, 2025). Also, we exclusively evaluate open-source checkpoints.

While CoT compression methods typically show a drop in accuracy as compression intensifies, TNT maintains or even improves the model’s accuracy, achieving the highest TE scores. This establishes TNT not merely as a compression tool, but as a reasoning enhancement framework that fundamentally reorganizes how models allocate computational resources. See details in Appendix B.2.

Finding 6. TNT exhibits excellent performance in Out of Distribution settings.

We also evaluate the performance of TNT in out-of-distribution (OOD) tasks. The experimental results show that TNT exhibits the best TE across all tested models. See details in Appendix B.3.

Finding 7. Reward hacking problem appears without our reward function.

Finally, we conduct the ablation study. Specifically, we test under the scenario where the components of our reward function specifically related to reward hacking problem are removed. Consequently, reward hacking problem occurs. See more details in Appendix B.4.

6 Conclusion

We explore the reward hacking problem in training hybrid reasoning models via RL, where correct answers from the no-thinking mode are rewarded more highly than those from the thinking mode. To mitigate this problem, we propose TNT, which leverages the solution component of the thinking mode to adaptively determine the maximum token usage of the no-thinking mode for different queries. This improves the poor detection of reward hacking in previous approaches that set a uniform maximum token usage across all queries, thereby effectively mitigates the reward hacking problem and achieves superior performance. Experiments show that TNT significantly improves accuracy and token efficiency while minimizing reward hacking.

Limitations

While TNT shows promising adaptive thinking capabilities, it encounters several limitations. First, similar to most prior open-source studies, we confine our model training to mathematical datasets. These datasets are preferred because they are easily accessible and can offer accurate, verifiable rewards. Furthermore, as discussed in Section 4.2, while TNT has the potential to be integrated with various techniques, our current study lacks sufficient testing due to computational resource constraints. As such, we have not explored whether combining TNT with these techniques can enhance performance.

Acknowledgments

This work is funded by Nanjing University China Mobile Communications Group Co.,Ltd. Joint Institute, National Natural Science Foundation of China (62192783, 62276128, 62406111), Jiangsu Natural Science Foundation (BK20243051), Jiangsu Science and Technology Major Project (BG2024031), the Fundamental Research Funds for the Central Universities(14380128, KG202514), the Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of Education of China (No. JYB2025XDXM118), the Collaborative Innovation Center of Novel Software Technology and Industrialization, and Shanghai Artificial Intelligence Laboratory.

References

- Daman Arora and Andrea Zanette. 2025. Training language models to reason efficiently. *arXiv preprint arXiv:2502.04463*.
- Hanting Chen, Yasheng Wang, Kai Han, Dong Li, Lin Li, Zhenni Bi, Jinpeng Li, Haoyu Wang, Fei Mi, Mingjian Zhu, and 1 others. 2025. Pangu embedded: An efficient dual-system llm reasoner with metacognition. *arXiv preprint arXiv:2505.22375*.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, and 1 others. 2024. Do not think that much for $2+3=?$ on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Gongfan Fang, Xinyin Ma, and Xinchao Wang. 2025. Thinkless: Llm learns when to think. *arXiv preprint arXiv:2505.13379*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Jialiang Hong, Taihang Zhen, Kai Chen, Jiaheng Liu, Wenpeng Zhu, Jing Huo, Yang Gao, Depeng Wang, Haitao Wan, Xi Yang, and 1 others. 2025. Reconsidering overthinking: Penalizing internal and external redundancy in cot reasoning. *arXiv preprint arXiv:2508.02178*.
- 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*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Lingjie Jiang, Xun Wu, Shaohan Huang, Qingxiu Dong, Zewen Chi, Li Dong, Xingxing Zhang, Tengchao Lv, Lei Cui, and Furu Wei. 2025. Think only when you need with large hybrid-reasoning models. *arXiv preprint arXiv:2505.14631*.
- Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, and 1 others. 2025. From system 1 to system 2: A survey of reasoning large language models. *arXiv preprint arXiv:2502.17419*.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.
- Chenwei Lou, Zewei Sun, Xinnian Liang, Meng Qu, Wei Shen, Wenqi Wang, Yuntao Li, Qingping Yang, and Shuangzhi Wu. 2025. Adacot: Pareto-optimal adaptive chain-of-thought triggering via reinforcement learning. *arXiv preprint arXiv:2505.11896*.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Li Erran Li, and 1 others.

2025. Deepscaler: Surpassing o1-preview with a 1.5 b model by scaling rl. *Notion Blog*.
- Alexander Pan, Erik Jones, Meena Jagadeesan, and Jacob Steinhardt. 2024. Feedback loops with language models drive in-context reward hacking, 2024. URL <https://arxiv.org/abs/2402.06627>.
- Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian Liang, Junxian He, and 1 others. 2025. A survey of efficient reasoning for large reasoning models: Language, multimodality, and beyond. *arXiv preprint arXiv:2503.21614*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems*, pages 1279–1297.
- Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, and 1 others. 2025. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*.
- Songjun Tu, Jiahao Lin, Qichao Zhang, Xiangyu Tian, Linjing Li, Xiangyuan Lan, and Dongbin Zhao. 2025. Learning when to think: Shaping adaptive reasoning in rl-style models via multi-stage rl. *arXiv preprint arXiv:2505.10832*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35.
- Fengli Xu, Qianyu Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, and 1 others. 2025. Towards large reasoning models: A survey of reinforced reasoning with large language models. *arXiv preprint arXiv:2501.09686*.
- Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.
- Zizheng Zhan, Ken Deng, Huaixi Tang, Wen Xiang, Kun Wu, Weihao Li, Wenqiang Zhu, Jingxuan Xu, Lecheng Huang, Zongxian Feng, and 1 others. 2025. Kat-v1: Kwai-autothink technical report. *arXiv preprint arXiv:2507.08297*.
- Jiajie Zhang, Nianyi Lin, Lei Hou, Ling Feng, and Juanzi Li. 2025. Adaptthink: Reasoning models can learn when to think. *arXiv preprint arXiv:2505.13417*.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*.

A Related Work

To train hybrid reasoning models, current methods fall into two categories: SFT-based and RL-only.

A.1 SFT-Based Training Method

The first method uses SFT to ensure a model possesses the ability to automatically select between thinking and non-thinking modes according to the difficulty of the query (Chen et al., 2025; Jiang et al., 2025; Lou et al., 2025; Zhan et al., 2025). Specifically, this method fine-tunes the model on a specially constructed dataset in which simple queries are paired with the non-thinking mode while complex queries are paired with the thinking mode. In some cases, RL is applied following SFT to further increase this capability (Jiang et al., 2025; Lou et al., 2025; Zhan et al., 2025). However, a primary drawback of this method is that it heavily relies on large-scale, high-quality labeled data. The process of curating the SFT dataset—namely classifying queries by difficulty and generating corresponding high-quality responses for both modes—remains highly computationally intensive. In addition, performing SFT with the obtained SFT dataset is likewise highly computationally intensive.

A.2 RL-Only Training Method

To avoid the substantial computational overhead associated with SFT, other works rely solely on RL to teach the model to automatically select between thinking and non-thinking modes, such as Thinkless, AutoThink, and AdaptThink (Fang et al., 2025; Tu et al., 2025; Zhang et al., 2025). Specifically, the model is incentivized with a higher reward for correct non-thinking answers than for correct thinking answers. Unfortunately, this training method is highly susceptible to reward hacking problem (Skalse et al., 2022; Pan et al., 2024). Precisely, this training method often rely on superficial signals—such as special control tokens—to assign rewards, without visibility into the model’s internal computational process. Consequently, the model might learn to alternative the tags to receive higher reward for generating the non-thinking control tokens while still performing thinking internally.

To mitigate this reward-hacking problem, Fang et al. (2025) augment an SFT process. Unlike the SFT-based training method mentioned previously, this SFT is used solely to fix the model’s outputs in thinking and non-thinking modes and does not teach the model to select between these

modes based on task difficulty. Although the data-collection process for this SFT is simpler, the computational overhead of this SFT training remains substantial.

To avoid such high computational overhead, Zhang et al. (2025) restrict the maximum token usage for the non-thinking mode to ensure the model does not masquerade as non-thinking mode while using the thinking mode to obtain higher returns. Unfortunately, existing work imposes the same maximum token usage for the non-thinking mode across all queries, whereas different queries require different maximum token usage. Specifically, consider a simple query, such as what is $1 + 1$ equal to, and a complex query, such as an AIME query; even when employing a long CoT to solve the simple query and continuously doing exploration, reflection, and self-verification for the response, token usage may still fail to beat solving the difficult AIME query via a CoT. In this situation, setting the maximum token usage for the non-thinking mode at a low level ensures the identification of reward hacking problem in the non-thinking mode responses to the simple query. However, this method results in almost all non-thinking mode responses to the complex query mode being classified as experiencing reward hacking problem. Conversely, if the maximum token usage is set high enough to distinguish reward hacking problem in the non-thinking mode responses to the complex query effectively, then almost all the non-thinking mode responses to the simple query are classified as not experiencing reward hacking problem. Therefore, applying a uniform maximum token usage for the non-thinking mode across all queries results in poor detection of reward hacking problem, leading only limited mitigation of the reward hacking problem.

Our TNT leverages the solution component of the thinking mode responses to determine the maximum token usage for the non-thinking mode across various queries. Consequently, TNT enhances its ability to identify whether reward hacking problem occurs, thereby effectively mitigating this problem and achieving superior performance.

B Addition Experiments

B.1 Experiment on DeepScaleR-1.5B and DeepSeek-R1-Distill-Qwen-7B

To further verify the universality and scalability of TNT, we extended our experiments to DeepScaleR-1.5B and DeepSeek-R1-Distill-Qwen-

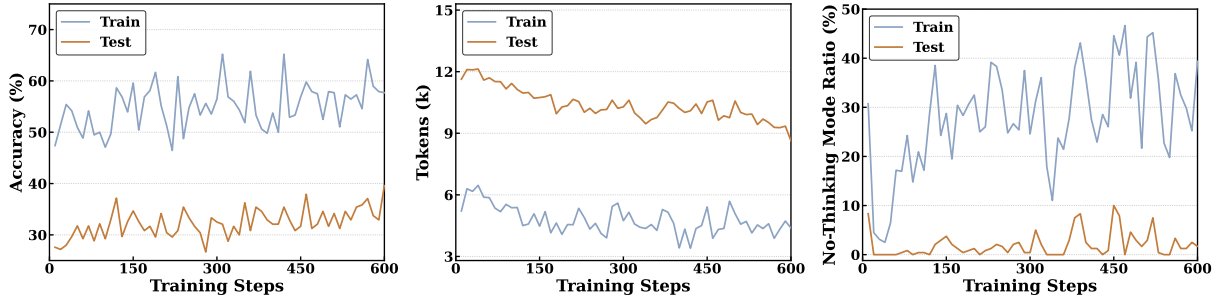


Figure 5: Accuracy (left), token usage (middle), and non-thinking ratio (right) in RL training on DeepSeek-R1-Distill-Qwen-1.5B. Due to constraints on computational resources, we only use AIME24 as the test dataset.

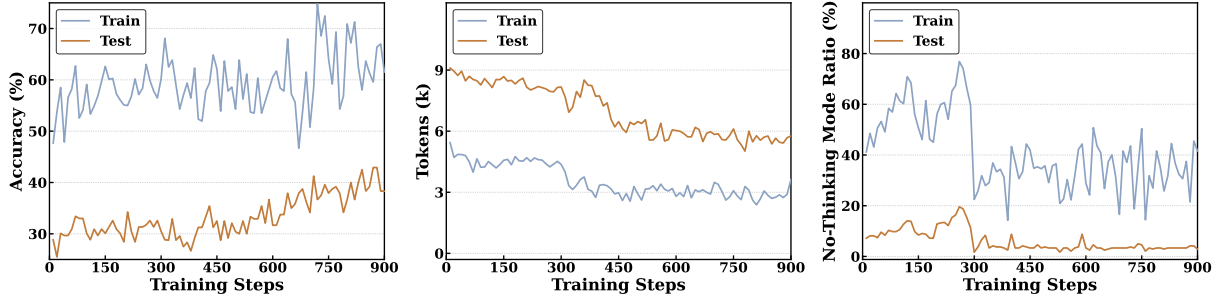


Figure 6: Accuracy (left), token usage (middle), and non-thinking ratio (right) in RL training on DeepScaleR-1.5B. Due to constraints on computational resources, we only use AIME24 as the test dataset.

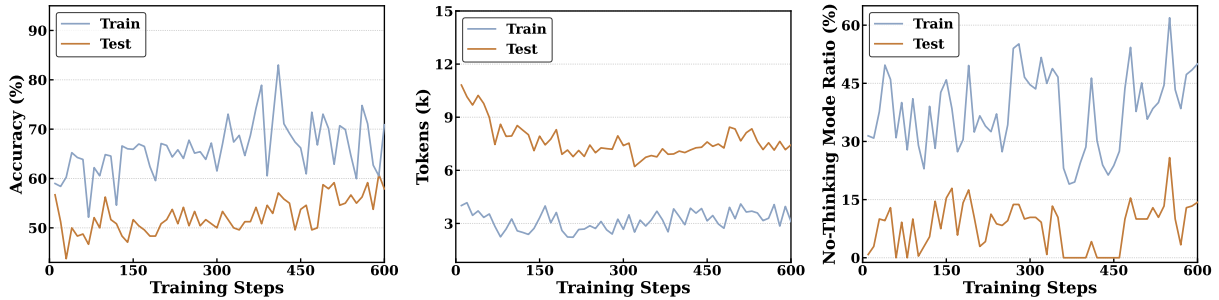


Figure 7: Accuracy (left), token usage (middle), and non-thinking ratio (right) in RL training on DeepSeek-R1-Distill-Qwen-7B. Due to constraints on computational resources, we only use AIME24 as the test dataset.

7B. As shown in Table 4, TNT demonstrates consistent superiority in TE across both base models, significantly reducing computational costs while maintaining competitive performance. Specifically, TNT reduces average token usage by approximately 41% for DeepScaleR-1.5B and over 51% for DeepSeek-R1-Distill-Qwen-7B compared to their respective base models. Importantly, this substantial reduction in token usage does not compromise model performance; TNT consistently matches or surpasses the accuracy of the base models and robust baselines. The training dynamics illustrated in Figure 6 and Figure 7 align perfectly with our observations on the 1.5B model, where accuracy steadily improves on the test set while token consumption naturally decreases throughout

the RL process. This confirms that the efficacy of our TNT is robust across different model scales and initial distributions without requiring explicit length penalties.

While the overarching trend of enhanced efficiency is consistent, the performance gains manifest with slight nuances between the two models. In the case of DeepScaleR-1.5B, which is already highly optimized for mathematical reasoning, the primary contribution of TNT lies in aggressive efficiency optimization. Although its average accuracy of 43.5% is comparable to the 43.7% achieved by AutoThink-Stage2, TNT achieves this result using only 3908 average tokens compared to the 6535 tokens required by the latter. This results in a superior TE of 0.70 compared to the 0.54 score for

Models	Accuracy (%) \uparrow						Token Usage \downarrow						TE
	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	
Base Model: DeepScaleR-1.5B													
Base Model	<u>41.9</u>	<u>30.6</u>	28.4	71.4	42.9	43.0	9123	9113	4175	5275	5375	6612	0.53
Auto-S1	32.3	24.8	26.4	65.2	41.5	38.0	12819	11779	5203	5467	7068	8467	0.41
Auto-S2	39.8	31.0	<u>31.6</u>	<u>72.2</u>	44.0	43.7	<u>8651</u>	<u>8772</u>	4467	5351	5434	6535	<u>0.54</u>
Auto-S3	42.0	30.0	30.7	71.4	42.9	43.4	9030	9204	<u>3819</u>	<u>4682</u>	<u>5073</u>	<u>6362</u>	<u>0.54</u>
TNT (Ours)	39.6	30.2	31.8	72.7	<u>43.1</u>	<u>43.5</u>	5479	5216	2402	3208	3236	3908	0.70
Base Model: DeepSeek-R1-Distill-Qwen-7B													
Base Model	<u>57.6</u>	38.4	38.0	81.7	<u>49.5</u>	<u>53.0</u>	12804	14786	5028	7733	8812	9833	0.53
Adapt-5	56.1	37.1	34.0	81.0	48.3	51.3	9966	11553	2528	5284	6381	7142	0.61
Auto-S1	51.7	36.7	31.8	77.7	46.9	49.0	8625	9575	<u>1763</u>	<u>4517</u>	<u>5117</u>	<u>5919</u>	0.64
Auto-S2	55.0	<u>40.1</u>	<u>37.6</u>	<u>82.4</u>	49.4	52.9	9078	9707	3019	4710	5723	6447	0.66
Auto-S3	53.7	<u>40.1</u>	37.4	81.3	49.2	52.3	<u>8210</u>	<u>9446</u>	2759	4861	5549	6165	<u>0.67</u>
TNT (Ours)	57.9	41.7	36.5	83.9	51.1	54.2	6855	7595	1628	3619	4103	4760	0.79

Table 4: Comparison of accuracy, token usage, and TE on mathematical benchmarks across hybrid reasoning models with base models DeepScaleR-1.5B and DeepSeek-R1-Distill-Qwen-7B. The best and second results are bolded and underlined, respectively. We did not compare ThinkLess on DeepScaleR-1.5B and DeepSeek-R1-Distill-Qwen-7B, nor did we compare AdaptThink on DeepSeek-R1-Distill-Qwen-7B, due to the lack of open-source checkpoints for these methods on the corresponding models. Especially, we exclusively present adapt-5 for DeepSeek-R1-Distill-Qwen-7B, since it serves as the only open-source checkpoint provided by AdaptThink based on this base model.

Models	Non-Thinking Mode Tokens						Thinking Mode Tokens					
	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	AIME24	AIME25	Minerva	AMC23	Olym.	AVG
TNT-DeepSeek-1.5B	995	795	601	859	937	837	8633	8325	3475	7159	6086	6736
TNT-DeepScaleR-1.5B	2779	-	812	967	1206	1153	5663	5216	2475	4020	3775	4230
TNT-DeepSeek-7B	1552	1359	702	1214	1185	1202	7738	8435	2924	6069	6327	6299

Table 5: No-thinking and thinking mode token usage of TNT on mathematical benchmarks. TNT-DeepSeek-1.5B/7B refer to TNT with base mode DeepSeek-R1-Distill-Qwen-1.5B/7B and TNT-DeepScaleR-1.5B refer to TNT with base mode DeepScaleR-1.5B. We also follow these abbreviation conventions in the table below.

Models	No-Thinking Mode Ratio (%)				
	AIME24	AIME25	Minerva	AMC23	Olym.
TNT-DeepSeek-1.5B	1.7	0.4	11.3	30.7	20.7
TNT-DeepScaleR-1.5B	2.9	0	4.4	26.6	21.0
TNT-DeepSeek-7B	14.3	11.9	58.3	50.5	43.2

Table 6: Non-thinking mode ratio of TNT on mathematical benchmarks across base models.

AutoThink-Stage2 and AutoThink-Stage3. And for DeepSeek-R1-Distill-Qwen-7B, TNT achieves a comprehensive improvement by securing both the highest average accuracy of 54.2%, which surpasses the base model by 1.2%, and the lowest token usage among all compared methods. Consequently, it attains a remarkable TE score of 0.79, significantly outperforming the best AutoThink-Stage3 baseline score of 0.67, demonstrating that TNT outperforms baselines on strong reasoners, with the advantage becoming more pronounced as model capability increases.

We analyze the token usage across different base models to investigate the impact of model capacity on reasoning length. As shown in Table 5, the token consumption in the non-thinking mode is not statically bounded; instead, it exhibits significant behavioral differences depending on the base model’s strength. We analyze the impact of base model scale and initialization on the mode selection mechanism as shown in Table 6. A clear positive correlation is observed between model parameter size and the utilization of the no-thinking mode. When scaling from 1.5B to 7B parameters, the TNT framework exhibits a significantly higher propensity to bypass the reasoning process; for instance, the no-thinking ratio on the Minerva benchmark surges from 11.3% with TNT-DeepSeek-1.5B to 58.3% with TNT-DeepSeek-7B. This suggests that larger models, possessing richer internalized knowledge, can more frequently solve problems directly without invoking the extensive reasoning chain. Con-

Models	Accuracy (%) \uparrow						Token Usage \downarrow						TE
	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	AIME24	AIME25	Minerva	AMC23	Olym.	AVG	
<i>Base Model: DeepSeek-R1-Distill-Qwen-1.5B</i>													
Base Model	28.6	24.6	26.2	62.1	<u>43.6</u>	37.0	16865	16464	7490	11050	11808	12736	0.33
RE-0	<u>31.9</u>	<u>24.8</u>	<u>28.3</u>	<u>65.6</u>	40.6	<u>38.2</u>	13404	12807	4666	7912	8493	9456	0.39
RE-1	30.3	24.0	25.8	64.3	38.2	36.5	12660	12474	3633	6827	7665	8652	0.39
RE-4	28.9	21.4	19.7	62.4	33.5	33.2	9467	8960	1444	4902	5424	6039	0.43
TP-4k	29.2	21.2	<u>28.3</u>	65.3	39.4	36.7	8154	8261	3353	4843	5586	6039	0.47
TP-2k	27.3	22.2	26.7	64.4	38.2	35.8	<u>7541</u>	<u>7140</u>	<u>2430</u>	<u>4462</u>	<u>4815</u>	<u>5278</u>	0.49
TP-iter2k	27.3	20.6	26.2	65.1	39.0	35.6	6990	5947	2805	3602	3780	4625	<u>0.52</u>
<i>TNT (Ours)</i>	37.7	26.1	28.7	65.9	46.4	41.0	8537	8252	3150	4508	5020	5893	0.53
<i>Base Model: DeepScaleR-1.5B</i>													
Base Model	41.9	30.6	28.4	71.4	<u>42.9</u>	<u>43.0</u>	9123	9113	4175	5275	5375	6612	0.53
TP-4k	37.9	28.8	32.3	<u>72.6</u>	42.3	42.8	6592	6095	3181	3907	4070	4769	0.62
TP-2k	35.8	26.9	29.8	69.9	<u>42.9</u>	41.1	5687	5315	<u>2482</u>	3458	3575	4103	0.64
TP-iter2k	36.6	28.5	31.3	71.4	42.3	42.0	<u>5495</u>	<u>5074</u>	2595	<u>3355</u>	<u>3498</u>	<u>4003</u>	<u>0.66</u>
<i>TNT (Ours)</i>	<u>39.6</u>	<u>30.2</u>	<u>31.8</u>	72.7	43.1	43.5	5479	5216	2402	3208	3236	3908	0.70
<i>Base Model: DeepSeek-R1-Distill-Qwen-7B</i>													
Base Model	57.6	38.4	38.0	81.7	49.5	53.0	12804	14786	5028	7733	8812	9833	0.53
RE-0	59.5	<u>40.6</u>	<u>37.5</u>	<u>82.4</u>	<u>49.7</u>	<u>53.9</u>	11171	13217	4074	6703	7853	8604	0.58
RE-1	52.9	40.5	37.2	81.3	48.9	52.2	10696	11420	2972	5827	6761	7535	0.60
RE-4	55.2	38.8	32.9	81.0	48.3	51.2	<u>9775</u>	<u>10988</u>	<u>2611</u>	<u>5086</u>	<u>6088</u>	<u>6910</u>	<u>0.62</u>
<i>TNT (Ours)</i>	<u>57.9</u>	41.7	36.5	83.9	51.1	54.2	6855	7595	1628	3619	4103	4760	0.79

Table 7: Comparison of accuracy, token usage, and TE on mathematical benchmarks between TNT and CoT compression models with base models DeepSeek-R1-Distill-Qwen-1.5B, DeepScaleR-1.5B, and DeepSeek-R1-Distill-Qwen-7B. Gray represents the base models, green denotes efficient reasoning models, and orange represents TNT. RE- x refers to Efficient-Reasoning with $\alpha = x * 0.01$. TP- x refers to ThinkPrune- x . The best and second results are bolded and underlined, respectively. We did not compare RE on DeepScaler-1.5B, nor did we compare TP on DeepSeek-R1-Distill-Qwen-7B, due to the lack of open-source checkpoints for these methods on the corresponding models.

versely, the performance of TNT-DeepScaleR-1.5B closely mirrors that of TNT-DeepSeek-1.5B across all datasets. Although DeepScaleR is a derivative of DeepSeek-1.5B optimized for mathematical reasoning, its training process does not explicitly target or enhance direct answering (no-thinking) capabilities. Consequently, the intrinsic non-reasoning proficiency remains unchanged, resulting in a static distribution of mode selection ratios comparable to the original base model.

B.2 Comparison with CoT Compression Methods

To assess the competitiveness of TNT against recent CoT compression strategies, we compare it with two prominent methods: Efficient-Reasoning (RE) (Arora and Zanette, 2025), and ThinkPrune (TP) (Hou et al., 2025). As detailed in Table 7, TNT consistently outperforms both methods in

terms of TE across all three base models, demonstrating that our hybrid reasoning model training method offers a superior trade-off between accuracy and computational cost. Unlike standard CoT compression methods that often degrade performance in exchange for reduced token usage, TNT frequently enhances accuracy while simultaneously achieving deeper reductions in token usage. For instance, while RE and TP models typically show a drop in accuracy as CoT compression intensifies, TNT maintains or even improves upon the base model’s performance, achieving the highest TE scores. This establishes TNT not merely as a compression tool, but as a reasoning enhancement framework that fundamentally reorganizes how models allocate computational resources.

Analyzing the specific behaviors on different base models reveals distinct advantages of TNT. On the DeepSeek-R1-Distill-Qwen-1.5B model,

Models	Accuracy (%)	Token Usage	TE
Base Model	36.9	17343	0.28
Thinkless	28.8	5958	<u>0.37</u>
Adapt-0	33.3	19388	0.23
Adapt-2	35.9	18463	0.26
Adapt-5	37.4	17627	0.28
Auto-S1	33.3	9794	0.33
Auto-S2	34.3	12374	0.30
Auto-S3	31.3	8816	0.33
TNT (Ours)	35.9	<u>8803</u>	0.38

Table 8: Accuracy, token usage and TE comparison across models on GPQA Diamond benchmark.

	Step 100	Step 300	Step 500
RF in Eq. (6)	3.9	17.9	94.2
TNT	0	7.2	5.1

Table 9: Average probability of thinking-related verbs appearing in non-thinking mode responses across Minerva, AMC23, and Olympiad when using the reward function in Equation (6). RF refers to reward function.

TNT achieves a remarkable average accuracy of 41.0%, which is substantially higher than the 36.7% achieved by the best performing TP-4k and the 38.2% of RE-0. While TP-iter2k manages to lower token usage slightly more than TNT, its accuracy suffers significantly, dropping to 35.6%. In contrast, TNT balances a moderate token usage with superior accuracy, resulting in a leading TE of 0.53. Similarly, on the DeepScaleR-1.5B base, TNT surpasses all TP baselines in both efficiency and effectiveness. It achieves an average token usage of 3908, which is lower than the 4003 tokens used by TP-iter2k, while maintaining a higher average accuracy of 43.5% compared to 42.0% for the latter. The advantage is even more pronounced on the larger DeepSeek-R1-Distill-Qwen-7B model, where TNT attains an exceptional TE of 0.79. Here, TNT drastically cuts average token usage to 4760 compared to 6910 for the most aggressive baseline RE-4, all while securing the highest average accuracy of 54.2%. This confirms that TNT is uniquely capable of scaling efficiency gains without the accuracy penalties commonly associated with traditional length-penalty or pruning methods.

B.3 Experiments in OOD

To assess the out-of-distribution (OOD) generalization capabilities of TNT, we conducted experiments on the GPQA Diamond, which features challenging graduate-level queries in biology, physics, and

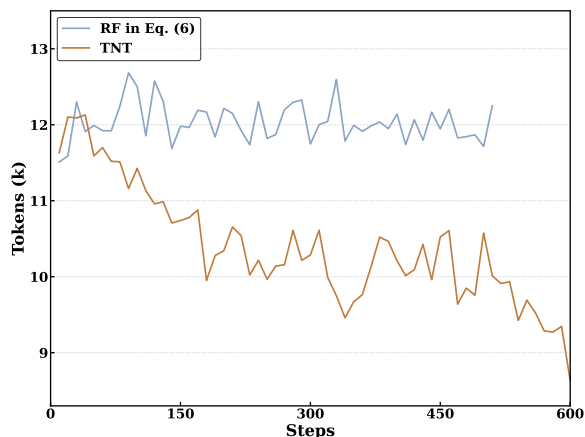


Figure 8: Token usage varies on AIME24 when using the reward function in Equation (6).

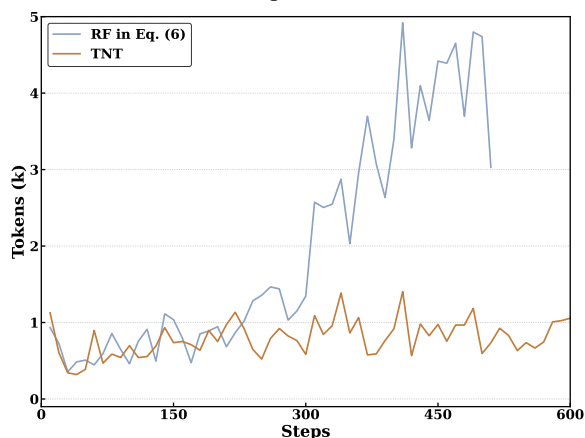


Figure 9: Token usage of the non-thinking mode responses varies on the training dataset when using the reward function in Equation (6).

chemistry. As shown in Table 8, while TNT does not achieve the single best score in accuracy or token usage alone, it strikes a superior balance between performance and computational cost, making it the best overall model. Specifically, TNT achieves the highest TE of 0.38, surpassing all other models. This indicates that it provides high accuracy while remaining highly efficient, a crucial advantage over models that are either slightly more accurate but far more resource-intensive or less accurate. This exceptional balance underscores TNT’s robust generalization and practical effectiveness on complex, scientific OOD tasks.

B.4 Ablation Study

Finally, we conduct tests under the scenario where the components of our reward function related to reward hacking problem are removed. Specifically, the reward function reduces to the naive reward function in Tu et al. (2025), as shown in the follow-

ing:

$$R(x, y_x^k, y_x^*, p(y_x^k), \mathcal{L}_x^N) = \begin{cases} 1, & \text{if } p(y_x^k) = 1 \wedge r(y_x^k, y_x^*) = 1, \\ 0, & \text{if } p(y_x^k) = 1 \wedge r(y_x^k, y_x^*) = 0, \\ 2, & \text{if } p(y_x^k) = 0 \wedge r(y_x^k, y_x^*) = 1, \\ -1, & \text{if } p(y_x^k) = 0 \wedge r(y_x^k, y_x^*) = 0. \end{cases} \quad (6)$$

We evaluate three key metrics: (i) average probability of thinking-related verbs appearing in responses classified as the non-thinking mode on the testing dataset, (ii) the average token usage across all responses on the testing dataset, and (iii) the average token usage for responses categorized under the non-thinking mode on the training dataset. The average probability of thinking-related verbs appearing in non-thinking mode responses on the testing dataset are presented in Table 9. Due to the high variance in calculating such probabilities for AIME24 and AIME25—originating from the limited sampling of non-thinking mode responses in these datasets—we focus solely on the average probabilities across Minerva, AMC23, and Olympiad. In addition, due to the high computational cost of evaluations on these three datasets, we are able to assess only a limited number of checkpoints. This constraint prevents us from performing evaluations every 10 steps, as shown in the subsequent experimental results. The average token usage across all responses on the testing dataset are presented in Figure 8. Note that the evaluation is conducted solely on AIME24 due to computational resource constraints. Specifically, we cannot evaluate checkpoints on all five mathematical benchmarks every 10 steps. The token usage for responses classified as the non-thinking mode on the training dataset are shown in Figure 9.

We observe that as training progresses, utilizing the reward function shown in Equation (6), token usage for responses classified as non-thinking mode on the training dataset significantly increases after a certain step. This phenomenon suggests a high probability of reward hacking problem since such token usage is too large. This finding aligns with results presented in Table 9, where TNT’s responses seldom exhibit reward hacking problem, whereas models using the reward function in Equation (6) frequently demonstrate such behavior. Consequently, the elevated probability of reward hacking problem leads to a persistent lack of reduction in token usage for responses on AIME24.

C A Case Study

To provide a more intuitive and qualitative understanding of TNT, we present a series of case studies comparing TNT with the DeepSeek-R1-Distill-Qwen-1.5B baseline across mathematical queries of varying difficulty. These cases, illustrated in Figures 10 and 11, demonstrate TNT’s ability to dynamically engage in the thinking or non-thinking mode. Our analysis highlights that for the simple query, TNT opts for the non-thinking mode that is resource-efficient; for the complex query that cause the baseline to fail, TNT engages in the thinking mode to arrive at the correct solution. This adaptive behavior underscores TNT’s superior flexibility and effectiveness in solving a wide spectrum of queries.

Question

$\frac{m}{n}$ is the irreducible fraction value of

$$3 + \frac{1}{3 + \frac{1}{3 + \frac{1}{3}}}$$

What is the value of $m + n$?

Ground Truth: 142

DeepSeek-R1-Distill-Qwen-1.5B Response

Thinking (1342tokens)

Okay, so I have this problem here: I need to find the irreducible fraction value of the continued fraction $3 + \frac{1}{3 + \frac{1}{3 + \frac{1}{3}}}$ and then determine the sum of the numerator and denominator of that fraction.

...

Adding 3 to that gives $\frac{33}{10}$, reciprocal is $\frac{10}{33}$, which adds correctly to 3, giving $\frac{109}{33}$. And since 109 is prime, it's irreducible. So, adding numerator and denominator gives 142. I think that's solid.

Final Answer: The value of $m + n$ is 142.

Solution (479tokens)

To find the irreducible fraction value of the continued fraction $3 + \frac{1}{3 + \frac{1}{3 + \frac{1}{3}}}$ and determine the sum of the numerator and denominator, we start by simplifying from the innermost part.

...

The fraction $\frac{109}{33}$ is irreducible because 109 is a prime number and does not share any common factors with 33 other than 1. Therefore, the sum of the numerator and the denominator is: $109 + 33 = 142$

The value of $m + n$ is 142.

TNT Response

Solution (579tokens)

To find the value of $\frac{m}{n}$ where $\frac{m}{n}$ is the irreducible fraction value of the continued fraction: $3 + \frac{1}{3 + \frac{1}{3 + \frac{1}{3}}}$

...

Conclusion:

The irreducible fraction value is $\frac{109}{33}$. Therefore, $m = 109$ and $n = 33$.

The value of $m + n$ is: $m + n = 109 + 33 = 142$

Final Answer: 142

Figure 10: For a easy query, TNT can solve it by switching to the non-thinking mode to use fewer tokens compared to DeepSeek-R1-Distill-Qwen-1.5B.

Question

Find the number of triples of nonnegative integers (a, b, c) satisfying $a + b + c = 300$ and $a^2b + a^2c + b^2a + b^2c + c^2a + c^2b = 6,000,000$.

Ground Truth: 601

DeepSeek-R1-Distill-Qwen-1.5B Response

Thinking (11156tokens)

Alright, so I've got this problem here where I need to find the number of triples of nonnegative integers (a, b, c) that satisfy two conditions. The first condition is that $a + b + c$ equals 300.

...

So, solving for b and c : the roots of $x^2 - 200x + 10,000 = 0$, which are $x = 100$ and $x = 100$. So, $b = c = 100$.

Thus, the only solution is $a = b = c = 100$.

Therefore, the number of ordered triples is $\boxed{1}$.

Solution (671tokens)

Given the problem of finding the number of triples of nonnegative integers $((a, b, c))$ that satisfy the conditions $a + b + c = 300$ and $a^2b + a^2c + b^2a + b^2c + c^2a + c^2b = 6,000,000$, we start by analyzing the given equations.

...

Thus, the only solution is the triple $(100, 100, 100)$, giving us one ordered triple. $\boxed{1}$

TNT Response

Thinking (10616tokens)

Okay, so I have this problem here where I need to find the number

...

Wait, maybe that's correct.

Therefore, the number of solutions is 601.

Hence, the final answer is $\boxed{601}$.

Solution (734tokens)

We need to find the number of triples of nonnegative integers (a, b, c) that satisfy the following conditions:

...

Summing these, we get $3 \times 201 = 603$. However, the case where all three variables are 100 is counted three times, so we subtract 2 to avoid overcounting, resulting in $603 - 2 = 601$.

Thus, the number of triples (a, b, c) is: $\boxed{601}$

Figure 11: For a hard query, TNT is able to solve it correctly through careful thinking, while DeepSeek-R1-Distill-Qwen-1.5B was not.