DYNAMIC EARLY EXIT IN REASONING MODELS

Chenxu Yang^{$\diamond \heartsuit$}, Qingyi Si^{\diamond}, Yongjie Duan^{\diamond}, Zheliang Zhu^{$\diamond \heartsuit$}, Chenyu Zhu^{\diamond}, Zheng Lin ^{$\diamond \heartsuit$}, Li Cao^{\diamond}, Weiping Wang^{\diamond},

▲Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China
 [∞] School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China
 [◊] Huawei Technologies Co., Ltd.

{yangchenxu,linzheng}@iie.ac.cn, siqingyi@huawei.com



"He who knows when he has enough is rich." — Lao Tzu

Figure 1: (a) Motivation and (b) Performance of Dynamic Early Exit in Reasoning.

ABSTRACT

Recent advances in large reasoning language models (LRLMs) rely on test-time scaling, which extends long chain-of-thought (CoT) generation to solve complex tasks. However, overthinking in long CoT not only slows down the efficiency of problem solving, but also risks accuracy loss due to the extremely detailed or redundant reasoning steps. We propose a simple yet effective method that allows LLMs to self-truncate CoT sequences by early exit during generation. Instead of relying on fixed heuristics, the proposed method monitors model behavior at potential reasoning transition points (e.g., "Wait" tokens) and dynamically terminates the next reasoning chain's generation when the model exhibits high confidence in a trial answer. Our method requires no additional training and can be seamlessly integrated into existing o1-like reasoning LLMs. Experiments on multiple reasoning benchmarks MATH-500, AMC 2023, GPQA Diamond and AIME 2024 show that the proposed method is consistently effective on deepseek-series reasoning LLMs, reducing the length of CoT sequences by an average of 31% to 43% while improving accuracy by 1.7% to 5.7%.

^{*} Equal Contribution. Work done when Chenxu Yang was an intern at Huawei.

[†] Corresponding Author.

1 INTRODUCTION

The emergence of large reasoning language models (Xu et al., 2025a), such as DeepSeek-R1 (DeepSeek-AI et al., 2025) and GPT-O1 (OpenAI, 2025), has marked a significant breakthrough in natural language processing, particularly in solving complex and intricate tasks(WANG et al., 2025). These models leverage the test-time scaling (Snell et al., 2024) law by generating a longer chain-of-thought (CoT) (Wei et al., 2023) with rich and diverse reasoning paths, unleashing the potential of their reasoning ability.

However, the generation of overlong CoT sequences also significantly increases computational overload and reasoning latency, which hinders their deployment in computationally sensitive real-life applications. Moreover, recent research (Chen et al., 2025; Team et al., 2025a) reveals an intrinsic overthinking problem in LRLMs: These models persistently generate verbose reasoning sequences (Wu et al., 2025; Cuadron et al., 2025), introducing irrelevant information and unnecessary thought steps. Such redundant processing not only wastes computational resources but also leads to accuracy degradation by derailing from correct reasoning paths to erroneous ones (see Questions 19 and 26 in Figure 1(a)). This redundancy can be attributed to the design of the supervised fine-tuning (Achiam et al., 2023; Wei et al., 2021; Ouyang et al., 2022) or Reinforcement Learning (Bai et al., 2022; Ouyang et al., 2022; Schulman et al., 2017; Ramesh et al., 2024) stage, where the ability to dynamically adjust its reasoning length during generation is overlooked, leaving a gap in the inference efficiency of LRLMs.

Intuitively, as the number of reasoning paths increases, more information is referenced when generating conclusions. If we can identify the critical point where the reasoning information becomes just sufficient (termed **Pearl Reasoning**) and force the model to stop further thinking and directly output conclusions at this point, we can achieve both accuracy and efficiency. The aim of this paper can be viewed as *finding such pearls in long CoT sequences*. To validate our motivation, we forced the model to switches from thinking to directly generating answers, at different transition points in the thought process. If the answers obtained are correct, the existence of such pearl reasoning is verified. As shown in Figure 1(a), about 75% samples do contain such pearls (early exit yields correct answers), even 36.7% samples required only less than half of the original reasoning paths to reach correct conclusions. Therefore, how to find the pearl reasoning is a high potential and valuable topic to achieve efficient reasoning.

To this end, we propose a novel, training-free approach **DEER** that allows large reasoning language models to achieve **D**ynamic **E**arly **E**xit in **R**easoning. It regards the key moments when the model switches thought chains in reasoning as chances of early exit, and prompting LRLMs to stop thinking and generate trial answers at these moments. The confidence of each trial answer is the decision-making reference of early exit in reasoning. Specifically, the proposed method contains three actions: 1) **Reasoning Transition Monitoring**: Inspired by the budget forcing technique (Muennighoff et al., 2025), which extends the reasoning chain by introducing tokens like "*Wait*" to force deeper thinking, we identify words such as "*Wait*" as critical points for reasoning transitions and monitoring its appearance, but repurposed them inversely to terminate the reasoning chain. When the reasoning transition points are generated, the action of 2) **Trial Answer Inducing** is triggered: we replace it with "*final answer*" tokens to induce the model to immediately generate a trial answer, which will be used for 3) **Confidence Evaluating**. If the confidence is sufficiently high, set the model to stop further thinking and generate a conclusion based on the generated thoughts. Otherwise, revoke the action of Trial Answer Inducing and continue reasoning along the original path.

Our method is simple yet effective, and can be seamlessly extended to different sizes of reasoning models, achieving excellent results in the six most popular benchmarks, including AIME 2024 (MAA Committees), AMC 2023 (AI-MO, 2024), GPQA Diamond (Rein et al., 2023), MATH-500 (Hendrycks et al., 2021), HumanEval Chen et al. (2021) and BigCodeBench Zhuo et al. (2024). Specifically, our method, when integrated into a series of Deepseek-distillation models, can reduce the length of CoT sequences by an average of 31% to 43% while improving accuracy by 1.7% to 5.7% across different reasoning benchmarks. On programming benchmarks, it achieves an average reduction of 64.9% in generation length with a 2.1 points increase in pass@1. Our DEER offers a plug-and-play solution for improving both the efficiency and accuracy of large reasoning language models.

2 PRELIMINARIES

2.1 The Generation Pattern of Large reasoning language models

In contrast to traditional large language models (*System 1*), large reasoning language models (*System 2*) (Li et al., 2025) exhibit distinct generation patterns during the inference stage. (1) LRLMs use delimiters to divide the output into two processes: slow thinking and conclusion. LRLMs conduct systematic and thorough reasoning in the slow thinking, ultimately summarizing the thought process and providing the final answer in the conclusion. (2) During the slow thinking process, LRLMs engage in complex thinking actions, such as Problem Restatement & Comprehension, Approach Exploration, and Result Verification (Luo et al., 2025b). We refer to each thinking action as a thinking chunk, and the transitions between these chunks are often marked by **action transition points**, including "*Wait*", "*Alternatively*", and "*Hmm*".

System 1: [Prompt] + [Completion],

System 2: [Prompt] + $\langle \text{think} \rangle$ + [Slow Thinking] + $\langle \text{/think} \rangle$ + [Conclusion], where $\langle \text{think} \rangle$ and $\langle \text{/think} \rangle$ are begin-of-thinking and end-of-thinking delimiters respectively.

2.2 BUDGET FORCING

The budget forcing method proposed by Muennighoff et al. (2025) cleverly leverages the unique generation patterns in LRLMs to control test-time computation, thus achieving performance improvement without additional training or modules. They proposed a simple decoding-time intervention by forcing a maximum and minimum number of slow-thinking tokens at test time. Specifically, upon reaching the maximum token count, they appended the end-of-thinking token delimiter along with "*Final Answer*:" to early exit the thinking stage. To enforce a minimum, they suppressed the generation of the end-of-thinking token delimiter and appended action transition points to the LRLM's ongoing thought process, thus encouraging the model to double-check its answer or attempt new reasoning approaches. However, the budget forcing method they proposed is static, leaving significant room for improvement. We further explored the unique generation patterns in LRLMs and introduced a dynamic budget forcing method to enable early exits.

3 MOTIVATIONS AND OBSERVATIONS

In this section, we analyze the overthinking problem in LRLMs and explore the impact of static early exits on model performance.

We selected AIME2024 (MAA Committees) as the test set for exploratory experiments to perform a qualitative analysis and further conducted a quantitative analysis through experiments on MATH-500 (Hendrycks et al., 2021), GPQA-Diamond (Rein et al., 2023) and AIME2024. All experiments were conducted on DeepSeek-R1-Distill-Qwen-14B (DeepSeek-AI et al., 2025). In our experiments, we first enabled the LRLM to perform a complete inference on the test set (including both the slow thinking and conclusion contents). Then, we preserved the thinking content and divided it into thinking chunks based on the action transition points. Samples with more than five thinking chunks were retained. For these samples, we retained varying proportions (20%-90%) of their thinking chunks and appended an end-of-thinking token delimiter to each truncated reasoning sequence to forcibly terminate the slow-thinking process. The model then generated its final conclusion based on the partial reasoning contents. For the conclusions obtained with varying thinking contents, we investigated the number of samples that remained correct after early exiting when they were originally correct, as well as the number of samples that became correct after early exiting when they were originally incorrect, across three datasets in Figure 2.

We define the critical point where the reasoning information becomes just sufficient as the pearl reasoning. Figure 1(a) demonstrates that approximately 75% samples contain such pearls (early exit yields correct answers), and 36.7% of the samples exhibit pearl reasoning within the first half of the reasoning path. Additionally, we observed that there are some samples for which the correct answer can only be obtained through early exits (e.g., Questions 11, 19 and 26 in Figure 1(a)). These observations validate our motivation that early exiting at the pearl reasoning position and directly



(b) The number of originally incorrect samples that become correct with early exiting across different datasets.

Figure 2: Quantitative experiment results on three test sets.

outputting conclusion improves accuracy while reducing computational costs. Specifically, early exiting at the pearl reasoning position for all 30 samples can reduce inference costs by 40% while achieving an 11% improvement in accuracy. The quantitative results in Figure 2(a) show that for MATH-500, 60.8% of correctly answered samples remain correct with early exits using only 20% of the reasoning steps. For GPQA, 35.1% of correctly answered samples can still be correct with early exits after only 20% of the reasoning steps.

Figure 2(b) illustrates that exiting at different positions corrects varying proportions of wrong answers. For the MATH dataset, the highest correction rate is achieved when exiting at 40% of the reasoning steps, whereas for the GPQA dataset, the optimal correction occurs when exiting at 50%. It seems that the optimal early exit point varies for each problem and is closely related to the inherent difficulty of the problem itself. Therefore, it is intuitive that relying on a static early exit strategy based on fixed heuristics is suboptimal, underscoring the necessity of designing a dynamic early exit mechanism. This quantitative analysis further demonstrates that pearl reasoning can enhance accuracy while reducing the length of reasoning through error correction.

4 Method

4.1 DYNAMIC EARLY EXIT IN REASONING

In this section, we introduce the Dynamic Early Exit in Reasoning (DEER) method to determine optimal positions for early exits (pearl reasoning path), thereby alleviating the overthinking issue.

The core idea behind DEER is that a model's confidence in its trial answer dynamically indicates whether the thinking information required for LRLMs to generate the final answer is sufficient. We observe that when the model's reasoning process is incomplete or flawed, the trial answer tends to exhibit significantly lower confidence. Conversely, when the reasoning is comprehensive and logically sound, the model generates answers with higher confidence, as illustrated in Figure 10. This suggests that while the model's parameter space inherently encodes an implicit awareness of the occurrence of the Pearl Reasoning, the oversight of dynamically varying lengths of reasoning chains during training prevents it from explicitly terminating reasoning early. Our method addresses this limitation by explicitly activating and utilizing this implicit reasoning-awareness to enable dynamic early exits.

As shown in Figure 3, the specific design of our DEER involves three modules to determine whether to exit early: reasoning transition monitor, answer inducer, and confidence evaluator.

As mentioned in Section 2.1, LRLM incorporates action transition points between thinking chunks, which DEER regards as potential early exit points. The reasoning transition monitor is designed to



Figure 3: An overview of the Dynamic Early Exit in Reasoning (DEER) method.

track the emergence of these points, preparing for the identification of optimal early exit opportunities.

When the LRLM pauses at a potential early exit point, the trial answer inducer module prompts the model to generate an intermediate answer based on the reasoning content produced so far. We incorporated the answer delimiters ($boxed{}$) into the prompt to facilitate a more precise identification of the trial answers.

$$\boldsymbol{A} = \text{LRLM}(\boldsymbol{P}, \boldsymbol{T}, \boldsymbol{I}) \tag{1}$$

where P denotes the input prompt, T denotes the generated thoughts, I denotes the answer inducer prompt, and $A = [a_0, a_1, \dots, a_n]$ is the trial answer.

The confidence evaluator module computes the confidence of the induced trial answer. We take the maximum predicted probability of each token as its confidence. For multi-token trial answers, the overall confidence is computed as the mean confidence across all constituent tokens.

$$p(a_t) = \operatorname{softmax}(\mathcal{M}(\boldsymbol{P}, \boldsymbol{T}, \boldsymbol{I}, \boldsymbol{a_{< t}}))$$
(2)

$$C = \frac{1}{n} \sum_{i=1}^{n} \max_{a_t \in \mathcal{V}} p(a_t)$$
(3)

where \mathcal{M} is the LM Head of the LRLM and its preceding components, with logits as the output.

Finally, the comparison between the obtained confidence and the empirical threshold λ determines whether to exit early. If $C > \lambda$, we consider the reasoning information currently generated by the LRLM to be sufficient, indicating that the model has reached the Pearl Reasoning. At this point, DEER stops further reasoning actions and proceeds to deliver the conclusion. Otherwise, the model reverts to the previous transition point to generate the next thinking chunk.

4.2 BRANCH-PARALLEL DECODING ACCELERATION

Intuitively, the computation of **Answer Inducer** and **Confidence Evaluator** in DEER introduces additional latency during inference, especially for the code generation task where the trial answers are still long. This further reduces the efficiency gains we've earned by greatly shortening the generated CoT sequences. To solve this problem, we integrate DEER with a branch-parallel acceleration strategy (Figure 4) to further address these efficiency limitations: (1) Multiple branches are linearized into a single sequence and generated in parallel using a special causal attention mask; (2) Dynamic KV cache management through confidence-based pruning. This strategy enables temporal overlap between the trial answer evaluation and the ongoing reasoning-chain generation, thus optimizing overall inference efficiency. More discussion about end-to-end latency will be added in the forthcoming version.



Figure 4: Branch-parallel decoding and dynamic KV cache management.

Algorithm 1 Dynamic Early Exit in Reasoning (DEER)

1: Initialization: Large Reasoning Language Model LRLM(·), zero-shot-CoT zs_cot, question, answer inducer prompt I, set of action transition points \mathbb{P} , end-of-thinking delimiter $\langle / think \rangle$, maximum length max_len , and confidence threshold λ . 2: $x \leftarrow zs_cot + question, r \leftarrow []$ 3: while $len(\mathbf{x}) < max_len$ do 4: $y \leftarrow \text{LRLM}(\boldsymbol{x})$ 5: if $y \in \mathbb{P}$ then ▷ Generate thoughts until meets action transition points 6: $A \leftarrow \text{LRLM}(x + I)$ > Prompt LRLM to generate trial answer tokens 7: Get C according to Equation 2 & 3 > Calculate the confidence of the trial answer 8: if $C > \lambda$ then 9: $x \leftarrow x + \langle / \text{think} \rangle, r \leftarrow r + \langle / \text{think} \rangle$ ▷ Exit when thinking is sufficient 10: end if 11: else 12: $oldsymbol{x} \leftarrow oldsymbol{x} + y, oldsymbol{r} \leftarrow oldsymbol{r} + y$ 13: end if 14: end while

```
15: return r
```

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Benchmarks. To thoroughly evaluate the models' performance across various reasoning capabilities, we have chosen four representative reasoning benchmarks widely used in the field: MATH-500 Hendrycks et al. (2021) is a challenging benchmark comprising competition-level problems drawn from diverse high school mathematics domains, including Prealgebra, Algebra, and Number Theory. For consistency with previous research, we adopt the same 500-problem subset originally curated by OpenAI for evaluation (Lightman et al., 2023). AMC 2023 (AI-MO, 2024) contains 40 mathematical problems, covering algebra, geometry, number theory, and combinatorics. The American Mathematics Competitions (AMC), organized by the Mathematical Association of America (MAA), are prestigious contests designed to develop problem-solving skills and identify mathematical talent. AIME 2024 (MAA Committees) comprises 30 challenge problems selected from the 2024 American Invitational Mathematics Examination (AIME). This prestigious contest evaluates participants' mathematical reasoning abilities across diverse domains, including arithmetic, algebra, counting, geometry, number theory, probability, and other secondary school math topics. Beyond these three math problems, we also conducted evaluations on scientific questions. GPQA (Rein et al., 2023) is a PhD-level benchmark consisting of high-quality questions spanning physics, chemistry, and biology subdomains. Notably, domain experts with PhDs in these fields achieved only 69.7% accuracy on this dataset OpenAI (2024). For our experiments, we specifically select the highest quality subset, known as GPQA Diamond (composed of 198 questions). We also used two programming benchmarks to validate our method: HumanEval Chen et al. (2021) is proposed by OpenAI, containing 164 hand-crafted (to avoid data leakage) Python programming tasks focusing on basic algorithms, each with function signatures, docstrings, canonical solutions, and unit tests. BigCodeBench Zhuo

Size	Method	MATH-500		AMC 2023		GPQA Diamond		AIME 2024		Overall	
		$\text{ACC} \uparrow$	$\text{LEN}\downarrow$	ACC \uparrow	$\text{LEN}\downarrow$						
32B	Vanilla	89.2	2036	87.5	2992	48.5	4293	56.7	5818	70.5	3785
	DEER	88.6 _{-0.6}	942 _{-54%}	87.5 _{+0.0}	2256 _{-25%}	59.1 _{+10.6}	2553 _{-41%}	63.3 _{+6.6}	3839 _{-34%}	74.6 _{+4.1}	2398 _{-37%}
14B	Vanilla	86.0	1747	82.5	3505	52.0	3842	50.0	6307	67.6	3850
	DEER	87.0 _{+1.0}	1001 _{-43%}	77.5 _{-5.0}	1799 _{-49%}	52.5 _{+0.5}	2281 _{-41%}	60.0 _{+10.0}	3639 _{-42%}	69.3 _{+1.7}	2180 _{-43%}
7B	Vanilla	86.0	1731	70.0	2513	24.8	3453	26.7	3787	51.9	2871
	DEER	87.4 _{+1.4}	975 _{-44%}	82.5 _{+12.5}	1622 _{-35%}	27.3 _{+2.5}	2265 _{-34%}	33.3 _{+6.6}	3235 _{-15%}	57.6 _{+5.7}	2024 _{-30%}
1.5B	Vanilla	68.4	2188	42.5	3312	5.6	4128	20.0	4200	34.1	3457
	DEER	73.2 _{+4.8}	1118 _{-49%}	47.5 _{+5.0}	2384 _{-28%}	6.6 _{+1.0}	2752 _{-33%}	20.0 _{+0.0}	3302 _{-21%}	36.9 _{+2.8}	2389 _{-31%}

Table 1: Experimental results on the open-source DeepSeek-R1-Distill series of models (DeepSeek-R1-Distill-Qwen models with different sizes).

et al. (2024) is designed as a real-world-oriented benchmark, which includes 1,140 tasks requiring interactions with 139 libraries and diverse function calls.

Reasoning Models. We conducted experiments on the open-source DeepSeek-R1-Distill-Qwen series of models (DeepSeek-AI et al., 2025), including models with sizes of 1.5B, 7B, 14B, and 32B parameters. All the distilled models are supervised fine-tuned with reasoning data generated by DeepSeek-R1 model. In addition, we also use the popular and competitive reasoning QwQ-32B (Team, 2025; Qwen et al., 2025), to further verify the robustness of our method.

Metrics. The goal of DEER is to maintain the correctness performance of LRLMs while avoiding the redundant token overhead caused by overthinking. To this end, we selected *Accuracy* (ACC) and *Generation Length* (LEN) as the evaluation metrics. *Accuracy* (ACC) is calculated as follows: $Accuracy = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}\{\mathcal{M}(\mathcal{LRLM}(x_i)) = y_i\}$, where x_i is the question and y_i is the ground-truth answer from the dataset. $\mathcal{M}(\cdot)$ extracts the answer from the LRLM's response. $\mathbb{I}\{\cdot\}$ is an indicator function that determines whether the inside given condition is valid. The accuracy evaluation is based on the evaluation framework publicly released by Ye et al. (2025) (LIMO). Intuitively, the longer the generated text, the greater the inference cost for LRLMs. Therefore, we calculate the average generation length per sample to evaluate the cost as follows: *Generation Length*(LEN) = $\frac{1}{N} \sum_{i=1}^{N} |\mathcal{LRLM}(x_i)|$, where $|\cdot|$ measures the number of generated words. For the two programming benchmarks, we use the Pass@1 metric to measure generated code correctness.

Implementation details. All evaluations are conducted in a Zero-shot Chain-of-Thought (CoT) setting with the following prompt: "*Please reason step by step, and put your final answer within* boxed{}." For the decoding strategy, we employ greedy decoding with a single sample for the correctness evaluation. The ground-truth answers to the evaluation problems in our experiments are all well-structured numerical values or options. Therefore, we apply rule-based evaluations directly to verify mathematical equivalence. We set the maximum generation length at 16,384 to ensure that the evaluation captures complete problem-solving attempts.

5.2 EXPERIMENTAL RESULTS

Main Results. DEER consistently and significantly reduces the reasoning length while enhancing accuracy across reasoning models of different scales and four challenging reasoning benchmarks. As shown in Table 1, DEER reduces the length of CoT sequences by an average of 31% to 43% while improving accuracy by an average of 1.7% to 5.7% across different reasoning benchmarks. These consistent and significant improvements highlight the robustness and effectiveness of DEER.

Notably, when integrated into DeepSeek-R1-Distill-Qwen-1.5B, DEER outperforms vanilla CoT by 4.8 points in accuracy while utilizing only 51% of the inference tokens on MATH-500. On the larger DeepSeek-R1-Distill-Qwen-32B model, DEER achieves a 10.6 point improvement on GPQA-Diamond and a 6.6 point improvement on AIME 2024. Overall, for the reasoning models



Figure 5: More detailed experimental results of DEER compared to Vanilla CoT.

DeepSeek-R1-Distill-Qwen-1.5b and -7b, DEER is more effective than vanilla CoT on the simpler MATH-500 and AMC 2023. Conversely, the larger DeepSeek-R1-Distill-Qwen-14B and -32B saw the greatest benefits on the harder GPQA-Diamond and AIME24. This indicates that LRLMs also exhibit overthinking phenomenon when tackling challenging problem, while such phenomenon is more obvious when the model reasoning capability matches the benchmark difficulty.

Fine-Grained Analysis. As shown in Figure 5, we have divided the samples into four categories, each of which is represented by four colors. Green bars indicate samples that can only be answered correctly by early exit, red bars indicate samples that are only answered correctly without early exit, blue bars indicate samples that can always be answered correctly by early exit or not, and gray bars indicate samples that would be answer incorrectly anyway. It can be seen that DEER corrects more answers (green bars) than it alters incorrectly (red bars) through early exits.

The difficulty of the two datasets on the left side (MATH-500 and AMC23) of the figure is lower than that of the two datasets on the right side (GPQA-Diamond, AIME24). The gray bars in Figure 5 are considered as samples whose difficulty exceeds the capability of the model itself, so the sum of the green, blue, and red samples can be considered as the upper limit of the capability of DEER. The higher the proportion of green bars in these three colors (namely green-in-correct ratio), the more severe the overthinking phenomenon. By comparing horizontally from left to right, we observe that green-in-correct ratio decreases, indicating that LRLMs are more likely to experience overthinking on difficult problems than on simple ones. By comparing vertically from top to bottom, green-in-correct ratio increases, indicating that smaller models are more prone to overthinking compared to larger models.

Performance Ceiling of DEER. In addition, we found that there is still room for further improvement in DEER. The red and gray numbers in the figure represent the ceiling of DEER for smaller models and the accuracy of vanilla CoT for larger models, respectively. If DEER could avoid early exits on the samples marked by the red bars, the performance of the 7B model on MATH-500 and AMC23 could surpass that of the 32B model. Similarly, on GPQA-Diamond and AIME24, the 14B model's performance could exceed that of the 32B model. This suggests that more precisely identifying the pearl reasoning, for better decisions on whether and when to early exit, is a meaningful direction for future research.

Performance on Programming Tasks. Figure 6(b) presents the experimental results of DEER using DeepSeek-R1-Distill-Qwen-14B as the backbone on two programming datasets. When the threshold is set to 0.97, compared to vanilla CoT, DEER achieves a 62.7% reduction in reasoning length while improving pass@1 by 4.3 points on HumanEval. On BigCodeBench, it reduces the reasoning length by 67.1% with only a 0.2-point decrease in pass@1. These results demonstrate that our method effectively addresses overthinking in programming tasks.

5.3 DISCUSSION

Impact of threshold λ . Figure 6(a) shows the performance of DEER (DeepSeek-R1-Distill-Qwen-7B) on the MATH-500 dataset with different threshold λ . In the left subplot of Figure 6(a), we present the experimental results for λ values set to 0.9 (too early), 0.95, and 1.0 (too late). The results indicate that when the threshold is set too low, a minor additional reduction in reasoning



Figure 6: More experimental results.

Dataset	Method	ACC	LEN	Chunk Size	Chunk Num	Early Exit Ratio
MATH-500	Vanilla	86.0	1747	1747	1.0	0.0%
	DEER(W)	87.0 _{+1.0}	1001 _{-43%}	118.8	14.7	97.4%
	DEER(A)	88.6 _{+2.6}	1354_22%	329.6	5.3	95.6%
AMC 2023	Vanilla	82.5	3505	3505	1.0	0.0%
	DEER(W)	77.5 _{-5.0}	1799 _{-49%}	193.6	18.1	95.0%
	DEER(A)	82.5 _{+0.0}	2299 _{-34%}	687.3	5.1	87.5%
GPQA Diamond	Vanilla	52.0	3842	3842	1.0	0.0%
	DEER(W)	52.5 _{+0.5}	2281 _{-41%}	77.9	49.3	73.2%
	DEER(A)	52.5 _{+0.5}	2749 _{-28%}	205.5	18.7	46.0%
AIME 2024	Vanilla	50.0	6307	6307	1.0	0.0%
	DEER(W)	$60.0_{+10.0}$	3639 _{-42%}	113.4	55.6	90.0%
	DEER(A)	66.7 _{+16.7}	4209-33%	548.4	11.5	70.0%
Overall	Vanilla	67.6	3850	3850	1.0	0.0%
	DEER(W)	$69.3_{\pm 1.7}$	2180 _{-43%}	125.9	34.4	88.9%
	DEER(A)	72.6 _{+5.0}	2653 _{-31%}	442.7	10.2	74.8%

Table 2: Experimental results on the DeepSeek-R1-Distill-Qwen-14B model with different action transition points are presented. In DEER(W), the action transition point is denoted by "*Wait*", whereas in DEER(A), it is denoted by "*Alternatively*".

length leads to a significant drop in accuracy, reflecting an overcorrection of overthinking. Conversely, when the threshold is set too high, the model exits reasoning too late, resulting in prolonged reasoning lengths with a decline in accuracy. In the right subplot of Figure 6(a), we present the experimental results for λ values set to 0.94, 0.95, 0.96, and 0.97. It indicates that our method is robust to hyperparameters, as minor fluctuations in the parameters have little impact on the results.

The results on the programming task in Figure 6(b) demonstrate that as λ increases, DEER's pass@1 improves while the generation length also becomes longer. We attribute this to the extended length of trial answers in programming tasks, wherein the probabilities of many subsequent consecutive tokens tend to approach 1 due to the presence of fixed pattern segments (e.g., the last line often being "return variables"). The entropy differences among critical tokens are diluted when averaged. To address this, we plan to explore a design that compares the probabilities of key tokens against a threshold for optimization.

Impact of Thought Switch Signals. In the main experiments, we achieved promising results by using "*Wait*" as the signal for early exit, denoted as DEER(W). Considering that LRLMs may generate different action transition points, we added a set of comparative experiments using "*Alternatively*" as the checkpoint token, denoted as DEER(A). Table 2 shows that DEER(A) achieves a higher accuracy, surpassing DEER(W) by an average of 3.7 points. However, its reduction in inference length is 12% less. We collect statistics on the length and number (namely Chunk Size and Chunk Num) of reasoning chunks obtained by dividing the original thought chain with the two thought switch signals, respectively. As shown in Table 2, the average length of each reasoning chunk of DEER(A) is 3.5 larger than that of DEER(W), which makes the reasoning chains referenced when DEER(A) generates a trial answer more complete, thereby improving accuracy. Correspondingly, the number of reasoning chunks of DEER(A) is significantly smaller than that of DEER(W), which makes DEER(A) less likely to attempt early exit. Therefore, DEER(A) has 14.1% fewer samples than DEER(W) to implement early exit. This also explains why the average infer-



Figure 8: Experimental results on the QwQ-32B model.

ence length of DEER(A) is higher than that of DEER(W). This phenomenon inspires us that a better reason chunk segmentation can further improve early exit in reasoning, which we leave for future discussion.

Efficiency Improvement. To accurately verify the gains brought by DEER in end-to-end inference efficiency, we measured the average latency on MATH and AMC datasets based on huggingface Wolf et al. (2020) without any inference acceleration frameworks such as vllm Kwon et al. (2023)). As shown in Table 7, our method still reduces the latency by 43.4% to 47.3% when the proposed branch-parallel decoding acceleration are not used. Considering that multiple actions of trial answer inducing and the confidence evaluation will introduce a large amount of calculation delay, once the DEER uses the parallel acceleration strategy, all actions' extra delays except the last one can be overlapped, and the efficiency is expected to be further improved. This result and the efficiency gain analysis combined with the inference framework vllm will be presented in the next version.



Figure 7: Generation latency

Performance on QWQ-32B. Figure 8 presents the results on the QwQ-32B model. Compared to vanilla CoT, DEER reduces the generation length by 24% but makes three additional errors. This result is slightly weaker compared to the DeepSeek-R1-Distill series of models. Through observation, we found that appending the end-of-thinking delimiter does not reliably control QwQ-32B to exit reasoning, which makes QwQ-32B unable to fully absorb the gains of DEER. To explore the causes of this phenomenon, we compare the CoTs generated by QwQ-32B and the DeepSeek-R1-Distill models in Figure 11. It can be seen that $\langle /think \rangle$ token is not the only one in the QwQ-32B's CoT, which indicates that QwQ-32B itself cannot strictly follow the effect of terminating thinking brought by $\langle /think \rangle$. This further leads to a slightly less effective mitigation of the overthinking problem compared to other models, but overall, DEER remains effective. We conjecture that this behavior may be related to the format reward setting during the RL training phase of the QwQ-32B model.

5.4 CASE STUDY

In Figure 9, we provide examples of results on MATH-500 to visually demonstrate the effectiveness of DEER. The design of DEER ensures that it follows the same reasoning process as the vanilla CoT method before early exiting. Both methods arrive at the correct answer during the first reasoning step, as shown in the green box. The difference lies in the fact that our method exits early after evaluating the confidence of the trial answer as sufficiently high, thus producing the correct result. In contrast, the vanilla CoT method proceeds to the next reasoning action. After double-checking and switching reasoning approaches, the model becomes trapped in an endless cycle of verification due to inconsistent answers from the two approaches, ultimately failing to provide a final answer.

Table 10 shows the detailed process of DEER applied on a mathematical example. It can be observed that, at each reasoning switch point ("*Wait*" token), DEER generates a trial answer and evaluates its confidence. The change in confidence is consistent with the reliability of the current reasoning chunks and trial answers. This shows that LRLMs implicitly know when to leave early, and our method is simple and effective to realize such potential of the model itself.

6 RELATED WORK

The advent of Open-AI o1 (OpenAI, 2025) established test-time scaling (Snell et al., 2024) as a pivotal research direction in the LLM community. This approach enhances LLMs' slow thinking capabilities, enabling breakthroughs in complex problem solving. The recent open-sourcing of DeepSeek-R1 (DeepSeek-AI et al., 2025) has further intensified interest in locally deployed reasoning models. However, two critical challenges have emerged: 1) excessively long CoT generated significantly degrades inference efficiency, and 2) growing empirical evidence (Chen et al., 2025; Team et al., 2025a) reveals their susceptibility to overthinking – a phenomenon where models continue reasoning beyond the point of optimal output. Following the taxonomy of efficient reasoning established in (Sui et al., 2025; Wang et al., 2025), we categorize related work into three classes: post-traning based, prompt-based, and output-based efficient reasoning methods.

Post-training based efficient reasoning methods use supervised fine-tuning (Yu et al., 2024; Kang et al., 2025; Xia et al., 2025; Ma et al., 2025b; Munkhbat et al., 2025; Liu et al., 2024; Han et al., 2024) with variable-length CoT data or incorporate length rewards (Team et al., 2025b; Luo et al., 2025a; Aggarwal & Welleck, 2025; Arora & Zanette, 2025; Yeo et al., 2025; Shen et al., 2025b; Qu et al., 2025; Cui et al., 2025) in reinforcement learning to enable the model to adaptively generate chains of thought of different lengths. However, these methods often require a large amount of computational resources and face challenges in dataset construction. Recently, some work (Hao et al., 2024; Shen et al., 2025c; Cheng & Van Durme, 2024; Shen et al., 2025a; Su et al., 2025; Saunshi et al., 2025) has shown that using latent representations to replace explicit textual reasoning steps allows reasoning models to be more efficient. However, such methods often require extensive-epoch SFT on carefully curated datasets(Hao et al., 2024; Xu et al., 2025c), leading to overfitting on the output format and consequently compromising the model's inherent expressiveness and generalization ability.

Prompt-based efficient reasoning methods (Han et al., 2024; Xu et al., 2025b; Lee et al., 2025; Renze & Guven, 2024; Chen et al., 2024) use varying prompts to enforce reasoning models to generate concise CoT with less unnecessary reasoning steps. Especially, (Aytes et al., 2025; Chuang et al., 2024; 2025; Ong et al.) assign different prompts to queries based on their difficulty, thereby adjusting the length of the CoT generated by reasoning models. We also explored the performance of our method combined with prompt design in Table X, demonstrating further reductions in the length of reasoning chains while maintaining considerable accuracy.

Most of the **Output-based efficient reasoning** methods focus on optimizing the best-of-N sampling for LLMs, such as pruning low-quality samples (Xie et al., 2023; Liao et al., 2025) and implementing early stopping (Li et al., 2024; Manvi et al., 2024; Aggarwal et al., 2023) when multiple samples achieve self-consistency. However, following the introduction of advanced reasoning models like R1, there is less reliance on best-of-N sampling methods, as these models exhibit strong reasoning capabilities independently. Very recently, two concurrent works share similar motivations with ours. Zhang et al. (2025) also proposes to terminate early based on trial answers, but requires an additional probe model to determine the correctness. They focus on enhancing the verification capabilities of the probe model, whereas our method explore how to enable the model to self-determine when to exit early and integrate seamlessly into existing reasoning logic. Ma et al. (2025a) prompts reasoning models to directly output final answers during decoding, but only achieves better performance in the low-budget regime or being adapted to best-of-N methods compared to baselines, which limits the applicability and generalization.

7 CONCLUSION

This paper empirically verifies the rationale behind the early exit motivation in reasoning models during chain-of-thought generation through pilot experiments. Based on this, we propose a training-free dynamic early exit algorithm, which makes the reasoning model withdraw from subsequent thinking when the thinking amount is just enough. Our method comprehensively evaluated across reasoning models of varying model sizes and demonstrates superior performance with fewer tokens on six classical reasoning benchmarks (including math, science, and programming tasks), which offers a win-win solution to the trade-off between accuracy and efficiency commonly encountered in test-time scaling.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. arXiv preprint arXiv:2503.04697, 2025.
- Pranjal Aggarwal, Aman Madaan, Yiming Yang, et al. Let's sample step by step: Adaptiveconsistency for efficient reasoning and coding with llms. arXiv preprint arXiv:2305.11860, 2023.
- AI-MO. Amc 2023, 2024. URL https://huggingface.co/datasets/AI-MO/ aimo-validation-amc.
- Daman Arora and Andrea Zanette. Training language models to reason efficiently. *arXiv preprint* arXiv:2502.04463, 2025.
- Simon A Aytes, Jinheon Baek, and Sung Ju Hwang. Sketch-of-thought: Efficient llm reasoning with adaptive cognitive-inspired sketching. arXiv preprint arXiv:2503.05179, 2025.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.
- Qiguang Chen, Libo Qin, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. Unlocking the capabilities of thought: A reasoning boundary framework to quantify and optimize chain-of-thought. *Advances in Neural Information Processing Systems*, 37:54872–54904, 2024.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Do not think that much for 2+3=? on the overthinking of o1-like llms, 2025. URL https://arxiv.org/abs/2412.21187.
- Jeffrey Cheng and Benjamin Van Durme. Compressed chain of thought: Efficient reasoning through dense representations. arXiv preprint arXiv:2412.13171, 2024.
- Yu-Neng Chuang, Helen Zhou, Prathusha Sarma, Parikshit Gopalan, John Boccio, Sara Bolouki, and Xia Hu. Learning to route llms with confidence tokens. *arXiv preprint arXiv*, 2410, 2024.
- Yu-Neng Chuang, Leisheng Yu, Guanchu Wang, Lizhe Zhang, Zirui Liu, Xuanting Cai, Yang Sui, Vladimir Braverman, and Xia Hu. Confident or seek stronger: Exploring uncertainty-based ondevice llm routing from benchmarking to generalization. arXiv preprint arXiv:2502.04428, 2025.
- Alejandro Cuadron, Dacheng Li, Wenjie Ma, Xingyao Wang, Yichuan Wang, Siyuan Zhuang, Shu Liu, Luis Gaspar Schroeder, Tian Xia, Huanzhi Mao, et al. The danger of overthinking: Examining the reasoning-action dilemma in agentic tasks. arXiv preprint arXiv:2502.08235, 2025.
- Yingqian Cui, Pengfei He, Jingying Zeng, Hui Liu, Xianfeng Tang, Zhenwei Dai, Yan Han, Chen Luo, Jing Huang, Zhen Li, et al. Stepwise perplexity-guided refinement for efficient chain-ofthought reasoning in large language models. arXiv preprint arXiv:2502.13260, 2025.

- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. Tokenbudget-aware llm reasoning. *arXiv preprint arXiv:2412.18547*, 2024.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL https://arxiv.org/abs/2103.03874.
- 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, pp. 24312–24320, 2025.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL https://arxiv.org/abs/2309.06180.
- Ayeong Lee, Ethan Che, and Tianyi Peng. How well do llms compress their own chain-of-thought? a token complexity approach. *arXiv preprint arXiv:2503.01141*, 2025.
- Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Xinglin Wang, Bin Sun, Heda Wang, and Kan Li. Escape sky-high cost: Early-stopping self-consistency for multi-step reasoning. *arXiv* preprint arXiv:2401.10480, 2024.
- Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong, Zhijiang Guo, Le Song, and Cheng-Lin Liu. From system 1 to system 2: A survey of reasoning large language models, 2025. URL https://arxiv.org/abs/2502.17419.

- Baohao Liao, Yuhui Xu, Hanze Dong, Junnan Li, Christof Monz, Silvio Savarese, Doyen Sahoo, and Caiming Xiong. Reward-guided speculative decoding for efficient llm reasoning. arXiv preprint arXiv:2501.19324, 2025.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step, 2023. URL https://arxiv.org/abs/2305.20050.
- Tengxiao Liu, Qipeng Guo, Xiangkun Hu, Cheng Jiayang, Yue Zhang, Xipeng Qiu, and Zheng Zhang. Can language models learn to skip steps? *arXiv preprint arXiv:2411.01855*, 2024.
- Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. arXiv preprint arXiv:2501.12570, 2025a.
- Yijia Luo, Yulin Song, Xingyao Zhang, Jiaheng Liu, Weixun Wang, GengRu Chen, Wenbo Su, and Bo Zheng. Deconstructing long chain-of-thought: A structured reasoning optimization framework for long cot distillation, 2025b. URL https://arxiv.org/abs/2503.16385.
- Wenjie Ma, Jingxuan He, Charlie Snell, Tyler Griggs, Sewon Min, and Matei Zaharia. Reasoning models can be effective without thinking. *arXiv preprint arXiv:2504.09858*, 2025a.
- Xinyin Ma, Guangnian Wan, Runpeng Yu, Gongfan Fang, and Xinchao Wang. Cot-valve: Lengthcompressible chain-of-thought tuning. *arXiv preprint arXiv:2502.09601*, 2025b.
- Rohin Manvi, Anikait Singh, and Stefano Ermon. Adaptive inference-time compute: Llms can predict if they can do better, even mid-generation. *arXiv preprint arXiv:2410.02725*, 2024.
- MAA Committees. Aime problems and solutions. https://artofproblemsolving.com/ wiki/index.php/AIME_Problems_and_Solutions.
- 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, 2025. URL https://arxiv.org/abs/2501.19393.
- Tergel Munkhbat, Namgyu Ho, Seo Hyun Kim, Yongjin Yang, Yujin Kim, and Se-Young Yun. Selftraining elicits concise reasoning in large language models. arXiv preprint arXiv:2502.20122, 2025.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms with preference data, 2024. URL https://arxiv. org/abs/2406.18665.
- OpenAI. Learning to reason with llms, September 2024. URL https://openai.com/index/ learning-to-reason-with-llms/.
- OpenAI. Learning to reason with llms. https://openai.com/research/ learning-to-reason-with-llms, 2025. Accessed: 15 March 2025.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Yuxiao Qu, Matthew YR Yang, Amrith Setlur, Lewis Tunstall, Edward Emanuel Beeching, Ruslan Salakhutdinov, and Aviral Kumar. Optimizing test-time compute via meta reinforcement finetuning. arXiv preprint arXiv:2503.07572, 2025.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.

- Shyam Sundhar Ramesh, Yifan Hu, Iason Chaimalas, Viraj Mehta, Pier Giuseppe Sessa, Haitham Bou Ammar, and Ilija Bogunovic. Group robust preference optimization in reward-free RLHF. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=PRAsjrmXXK.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
- 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), pp. 476–483. IEEE, 2024.
- Nikunj Saunshi, Nishanth Dikkala, Zhiyuan Li, Sanjiv Kumar, and Sashank J Reddi. Reasoning with latent thoughts: On the power of looped transformers. *arXiv preprint arXiv:2502.17416*, 2025.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Xuan Shen, Yizhou Wang, Xiangxi Shi, Yanzhi Wang, Pu Zhao, and Jiuxiang Gu. Efficient reasoning with hidden thinking. *arXiv preprint arXiv:2501.19201*, 2025a.
- 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*, 2025b.
- Zhenyi Shen, Hanqi Yan, Linhai Zhang, Zhanghao Hu, Yali Du, and Yulan He. Codi: Compressing chain-of-thought into continuous space via self-distillation. *arXiv preprint arXiv:2502.21074*, 2025c.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL https://arxiv.org/abs/2408.03314.
- DiJia Su, Hanlin Zhu, Yingchen Xu, Jiantao Jiao, Yuandong Tian, and Qinqing Zheng. Token assorted: Mixing latent and text tokens for improved language model reasoning. *arXiv preprint arXiv:2502.03275*, 2025.
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, and Xia Hu. Stop overthinking: A survey on efficient reasoning for large language models, 2025. URL https://arxiv.org/abs/2503. 16419.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. Kimi k1.5: Scaling reinforcement learning with llms, 2025a. URL https://arxiv.org/abs/2501.12599.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025b.

- Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL https://qwenlm.github.io/blog/qwq-32b/.
- Hongru WANG, Deng Cai, Wanjun Zhong, Shijue Huang, Jeff Z. Pan, Zeming Liu, and Kam-Fai Wong. Self-reasoning language models: Unfold hidden reasoning chains with few reasoning catalyst. In Workshop on Reasoning and Planning for Large Language Models, 2025. URL https://openreview.net/forum?id=p4wXiD8FX1.
- Rui Wang, Hongru Wang, Boyang Xue, Jianhui Pang, Shudong Liu, Yi Chen, Jiahao Qiu, Derek Fai Wong, Heng Ji, and Kam-Fai Wong. Harnessing the reasoning economy: A survey of efficient reasoning for large language models, 2025. URL https://arxiv.org/abs/2503.24377.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface's transformers: State-ofthe-art natural language processing, 2020. URL https://arxiv.org/abs/1910.03771.
- Yuyang Wu, Yifei Wang, Tianqi Du, Stefanie Jegelka, and Yisen Wang. When more is less: Understanding chain-of-thought length in llms. arXiv preprint arXiv:2502.07266, 2025.
- Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. Tokenskip: Controllable chain-of-thought compression in llms. arXiv preprint arXiv:2502.12067, 2025.
- Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu Zhao, Min-Yen Kan, Junxian He, and Michael Xie. Self-evaluation guided beam search for reasoning. Advances in Neural Information Processing Systems, 36:41618–41650, 2023.
- Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao, and Yong Li. Towards large reasoning models: A survey of reinforced reasoning with large language models, 2025a. URL https://arxiv.org/abs/2501.09686.
- Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. Chain of draft: Thinking faster by writing less. *arXiv preprint arXiv:2502.18600*, 2025b.
- Yige Xu, Xu Guo, Zhiwei Zeng, and Chunyan Miao. Softcot: Soft chain-of-thought for efficient reasoning with llms. *arXiv preprint arXiv:2502.12134*, 2025c.
- Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning, 2025. URL https://arxiv.org/abs/2502.03387.
- Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chainof-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.
- Ping Yu, Jing Xu, Jason Weston, and Ilia Kulikov. Distilling system 2 into system 1. arXiv preprint arXiv:2407.06023, 2024.
- Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. Reasoning models know when they're right: Probing hidden states for self-verification. *arXiv preprint arXiv:2504.05419*, 2025.
- Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. *arXiv preprint arXiv:2406.15877*, 2024.

A APPENDIX



Figure 9: Comparison of generated content samples on MATH-500. Blue represents the conclusion, while other colors represent different reasoning actions.



Figure 10: An example demonstrating LRLM's implicit ability to sense the emergence of pearl reasoning.

$\label{eq:Question: Define \[p = \sum_{k = 1}^\infty \frac{1}{k^2} \q a = \sum_{k = 1}^\infty \frac{1}{k^3}.\] Find a way to write \[\sum_{j = 1}^\infty \frac{1}{k^3}.\] Find a way to write \[\sum_{j = 1}^\infty \frac{1}{(j + k)^3} \] find a way to write \[\sum_{k = 1}^\infty \frac{1}{(j + k)^3} \] find a way to write \[\sum_{k = 1}^\] find \[\sum_{k = 1}^\] find$

Length Answer: DeepSeek-R1-Distill-Qwen-32B $Okay, so I have this problem where I need to express the double sum \((\sum_{j = 1}^\\infty \sum_{k = 1}^\\infty \)) in terms of other in the second secon$ function at 3, which is known as Apery's constant, approximately 1.202... But I don't think I need the exact values; instead, I need to relate the double can change variables to simplify this. Let me set ||(n = j + k||). Then, for each fixed ||(n||), ||(j||) can range from 1 to ||(n - 1||), and ||(k = n - j||). So, perhaps I can rewrite the double sum as a single sum over \\(n\), where each term is the sum over \\(j\) from 1 to \\(n - 1\\) of \\(\\frac{1}{n^3}\\).\n\nLet me write that down:\n\\\[\n\\sum_{j = 1}^\\infty \\sum_{k = 1}^\\infty \\frac{1}{(j + k)^3} = \\sum_{n = 2}^\\infty \\sum_{j = 1}^{n - 1} \\frac{1}{n^3}\\\]\n\Yes that makes sense because the smallest \\(n\) can be is 2 (when \\(j = 1\\)) and \\(k = 1\\)), and it goes 408 to infinity. Now, for each (n, N), the inner sum is just adding $((n+1) n^3)$ a total of (n - 1) times. So, the inner sum simplifies to $((n+1) n^3)$ $\begin{array}{l} \left\{n^{3}\right\} \left(\left|n^{3}\right\} \right) \left(n^{1} \left(n^{2}\right) \right) \left(n^{1} \left(n^{2} \left(n^{2}\right) \right) \left(n^{1} \left(n^{2}\right) \right) \left(n^{1} \left(n^{2}\right) \right) \left(n^{1} \left(n^{2} \left(n^{2}\right) \right) \left(n^{2} \left(n$ \\frac {1} {n^3}\\). So, if I subtract the \\(n = 1\\) terms from both series, I can express the sums starting from \\(n = 2\\).\n\nLet's compute that:\n\nFor $\label{eq:linear} the first sum: ln/l[n/lsum_{n = 2}^linfty lfrac{1}{n^2} = p - l/rac{1}{1^2} = p - 1/n/l]n/nSimilarly, for the second sum: ln/l/[n/lsum_{n = 2}^linfty lfrac{1}{n^3} = q - l/rac{1}{1^3} = q - 1/n/l]n/nSo, substituting back into our expression: ln/l/[n/lsuff(p - 1 \/right) - \/left(q - 1 \/right) - l/right) = l/right) - l/$ $= \mathbf{p} - \mathbf{1} - \mathbf{q} + \mathbf{1} = \mathbf{p} - \mathbf{q} \cdot \mathbf{n}$ </think> $\n\nSo$, the double sum simplifies to $\(p - q\)$. Therefore, the answer is: $\n\n\(p - q)\n\)$ 15 QwQ-32B $Okay, so \ I \ need to \ find \ a \ way to \ express the \ double \ sum \ ((\ sum_{j = 1} \ (\ sum_{k = 1}$ (q_{1}) , where (p_{1}) is the sum of $(1/k^{2})$ and (q_{1}) is the sum of $(1/k^{3})$. Hmm, let me think about how to approach this. In First, I is the sum of $(1/k^{3})$. remember that double sums can sometimes be simplified by changing variables or switching the order of summation. Maybe I can let \\(n = j + k\\), so that for each (n(k)), the terms where (j + k = n(k)) contribute to the sum. That might help me rewrite the double sum as a single sum over (n(k)), with each term multiplied by the number of pairs (j, k)) that add up to (n(k)). In Let me try that. Let's set (n = j + k(k)). Then for each (n (n + j + k(k))). $2 \parallel (since \mid i \mid j) such that \mid i \mid j = k = n \parallel j = k = n \parallel j = k \parallel j = n \parallel j = k \parallel j = n \parallel j \mid j = n \parallel j \mid j = n \parallel$ only one pair (1,1); when l(n = 3l), there are two pairs (1,2) and (2,1), and so on. So the double sum can be rewritten as: $l(n)l[ln]sum_{n=2}(l(n+3)l(n+$ k)() would be grouped by (n = i + k), so each (n) contributes (n = 1) terms each of (n + 1). Therefore, the total contribution $\label{eq:linearized_linearized$ 387 Now I need to express this in terms of (p) and (q). Let me see. Let's first write out the sum: $n/n/[n]/m_{n=2}^{(n-1)}(n+1) \{n^3\} = 0$ $\frac{1}{n^2} - \frac{1}{n^2} - \frac{1$ $(|sum_{k=1}^{(1)}| + 2)), then ((|sum_{n=2}^{(1)}| + 2)), then ((|sum_{n=2}^{(1)}| + 2)))$ $\label{eq:lastic_last} \label{eq:lastic_last} \label{eq:lastic_last} \label{eq:lastic_last} \label{eq:lastic_last} \label{eq:lastic_lastic_last} \label{eq:lastic_lastic_last} \label{eq:lastic_last$ $\frac{|q_{1}|}{|q_{1}|} = \frac{|q_{1}|}{|q_{1}|} = \frac{|q_{1}|}{|q_{1}|}$ Let's summarize the above thought process and provide the result. The double sum can be rewritten by considering the sum over (n = j + k). For $= ach \left(\frac{1}{n} \right)^{n} = are \left(\frac{1}{n} \right)^{n} = are \left(\frac{1}{n} \right)^{n} = \frac{1}{n} =$ 260 $\label{eq:linear} $$ \eqref{linear} $$ \eqref{$ Answer** $\n\boxed{p - q}\n\n$ </think> $\ln \pi = 1^{\infty} \left(\frac{1}{(j + k)^3} \right)$ in terms of $(0, k) \in \mathbb{R}^{-1}$. changing variables. Let (n = j + k). For each (n (geq 2)), the number of pairs ((j, k)) such that (j + k = n) is (n - 1). This allows us to 128 l(n) = p - q(n) + q - q(n) + p - q(n) + p - q(n) + p - q(n) + p - q(n) + p

Figure 11: An example demonstrating the presence of more than one $\langle /\text{think} \rangle$ tokens in QwQ-32B. For both models, the green thinking chunk is sufficient to draw the correct conclusion. However, QwQ-32B proceeds with an additional summary (red chunk) and generates its own $\langle /\text{think} \rangle$ token. Based on all the above content, it arrives at the conclusion.