# 💡 INFTYTHINK: Breaking the Length Limits of Long-Context Reasoning in Large Language Models

Yuchen Yan[1,2]*        Yongliang Shen[1]†        Yang Liu[2]        Jin Jiang[2,3]

Mengdi Zhang[2]        Jian Shao[1]†        Yueting Zhuang[1]

[1]Zhejiang University [2]Meituan Group [3]Peking University
{yanyuchen, syl, jshao}@zju.edu.cn

## Abstract

Advanced reasoning in large language models has achieved remarkable performance on challenging tasks, but the prevailing long-context reasoning paradigm faces critical limitations: quadratic computational scaling with sequence length, reasoning constrained by maximum context boundaries, and performance degradation beyond pre-training context windows. Existing approaches primarily compress reasoning chains without addressing the fundamental scaling problem. To overcome these challenges, we introduce INFTYTHINK, a paradigm that transforms monolithic reasoning into an iterative process with intermediate summarization. By interleaving short reasoning segments with concise progress summaries, our approach enables unbounded reasoning depth while maintaining bounded computational costs. This creates a characteristic sawtooth memory pattern that significantly reduces computational complexity compared to traditional approaches. Furthermore, we develop a methodology for reconstructing long-context reasoning datasets into our iterative format, transforming OpenR1-Math into 333K training instances. Experiments across multiple model architectures demonstrate that our approach reduces computational costs while improving performance, with Qwen2.5-Math-7B showing 3-13% improvements across MATH500, AIME24, and GPQA_diamond benchmarks. Our work challenges the assumed trade-off between reasoning depth and computational efficiency, providing a more scalable approach to complex reasoning without architectural modifications.

📦 Project Page:  https://zju-real.github.io/InftyThink
🔘 Code:  https://github.com/ZJU-REAL/InftyThink

## 1 Introduction

Recent studies have demonstrated the remarkable reasoning capabilities of large language models (LLMs), with models like OpenAI o1 [1], DeepSeek-R1 [2], Gemini 2.0 Flash Thinking [3], QwQ [4], and Kimi-1.5 [5] surpassing human performance on high-difficulty tasks including mathematical competitions [6, 7]. These advanced reasoning models are typically developed through methodical

---

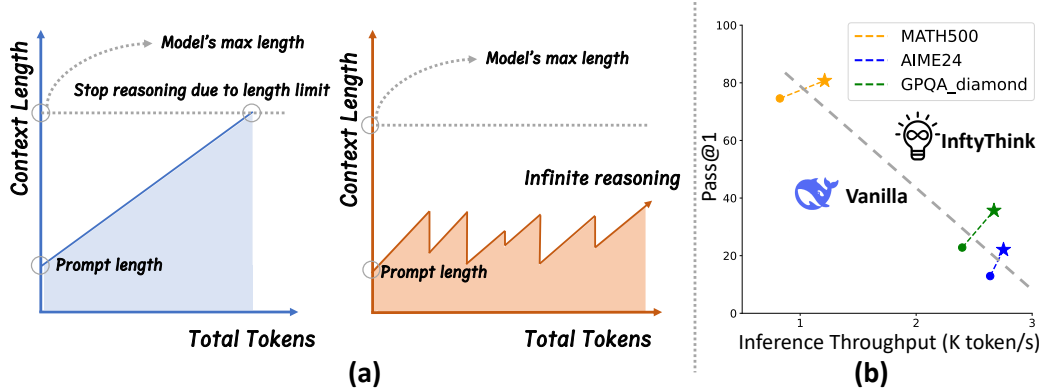*Contribution during internship at Meituan Group.
†Corresponding author.

Figure 1: **(a)** Computational complexity comparison between vanilla long-context reasoning (blue, left) and **INFTYTHINK** (orange, right). The sawtooth pattern of **INFTYTHINK** demonstrates how periodic summarization creates a bounded memory footprint, substantially reducing computational costs (smaller area under curve) while enabling deeper reasoning. **(b)** **INFTYTHINK** improves inference throughput while enhancing performance compared to the vanilla long-context reasoning.

techniques such as test-time scaling [8–11], post-training on long-thought trajectories [2, 12, 13], or large-scale reinforcement learning [2, 14] to generate effective reasoning paths that reach correct answers. A defining characteristic of these models is their ability to perform long-context reasoning, demonstrating advanced cognitive techniques including intent comprehension, multi-perspective analysis, self-reflection, and error correction [15, 16]. This evolution from simpler reasoning patterns to extensive deliberation has significantly improved problem-solving capabilities, particularly for complex challenges requiring multi-step inference.

However, this substantial improvement in reasoning quality comes with significant computational costs [17–19]. The computational complexity of decoder-based LLMs grows quadratically with sequence length, resulting in prohibitive resource requirements for long-form reasoning. This efficiency bottleneck manifests in three primary challenges: **First**, current reasoning models often generate thousands of tokens even for moderately complex problems [17, 20], creating substantial memory and processing overhead during inference. **Second**, reasoning processes are constrained by the model's maximum context length (aka. `max_length`) [21], frequently resulting in truncated reasoning that fails to reach conclusive answers. **Third**, most LLM architectures are pre-trained with relatively small context windows (4k-8k tokens), causing performance degradation when reasoning extends beyond these boundaries [22, 23].

Existing approaches [24–26] to address these limitations have explored various solutions with mixed success. Some methods attempt to compress reasoning chains post-generation [27–29], while others aim to train models to reason more concisely from the outset [30–33]. Chain-compression techniques like those employed in CoT-Valve [34] show promise but require predefined compression ratios during training, limiting their flexibility at inference time. TokenSkip [33] reduces redundant tokens by assessing each token's significance, though this impacts the model's reasoning performance. LightThinker [26] employs special tokens to dynamically compress the CoT process into a latent representation but lacks the ability to adaptively determine compression requirements for each step. Despite these advances, most approaches still operate within the traditional paradigm of generating a single, continuous reasoning chain, which merely attempting to make it more compact without addressing the fundamental computational scaling problem. This raises a critical question: ***Instead of optimizing within the constraints of monolithic reasoning, could we fundamentally re-imagine the reasoning process itself?***

In this paper, we propose a fundamentally different approach to long-context reasoning. Rather than viewing reasoning as a single extended process, we introduce **INFTYTHINK**, a novel paradigm that divides complex reasoning into multiple interrelated short reasoning segments. Each segment remains within a computationally efficient context length while maintaining the coherent flow of thought across iterations. This approach draws inspiration from human cognitive processes, where complex

problem-solving frequently involves breaking problems into manageable parts and summarizing intermediate progress.

The core mechanism of **INFTYTHINK** is an iterative process where the model generates a partial reasoning chain, summarizes its current thinking, and builds upon these summaries in subsequent iterations. As illustrated in Figure 1, traditional approaches (left, blue) face inevitable termination when context length reaches the model's maximum limit, often before completing the reasoning. In contrast, **INFTYTHINK** (right, orange) creates a sawtooth pattern through periodic summarization, enabling unbounded reasoning depth while maintaining a bounded memory footprint. This approach both reduces computational complexity (smaller area under the curve) and overcomes the fundamental ceiling on reasoning depth imposed by context length constraints. Beyond computational efficiency, **INFTYTHINK** offers a crucial advantage: it enables reasoning of arbitrary depth without architectural changes to the underlying model. By summarizing and building upon previous reasoning in manageable segments, models can effectively navigate complex problem spaces that would otherwise exceed context limitations.

To validate our approach, we reconstructed the existing SFT dataset OpenR1-Math, which was distilled from DeepSeek-R1, adapting it to conform to our proposed **INFTYTHINK** paradigm. This reconstruction process transformed the original long-form reasoning examples into multiple interconnected reasoning segments with corresponding summaries. We then fine-tuned multiple base architectures on this reconstructed dataset and conducted comprehensive comparisons against traditional single-round long-context reasoning methods. Our experimental results demonstrate consistent improvements across various benchmarks, with Qwen2.5-Math-7B showing 3% improvement on MATH500, 13% improvement on AIME24, and 10% improvement on GPQA_diamond.

Our contributions are summarized as follows:

- We introduce **INFTYTHINK**, which transforms monolithic long-form reasoning into iterative reasoning with summarization, mimicking human working memory patterns and reducing the quadratic computational complexity of transformer-based models to a more manageable form.
- We develop a technique to reconstruct existing long-context reasoning datasets (demonstrated on OpenR1-Math) into our iterative format, preserving reasoning quality while enabling more efficient computation without specialized architectures.
- Across multiple model architectures, our approach achieves significant improvements while substantially reducing computational costs, challenging the assumed trade-off between reasoning depth and efficiency.

## 2    Methods

In this section, we present **INFTYTHINK**, a novel reasoning paradigm that addresses the computational inefficiency of conventional long-context reasoning in large reasoning models. First, we formalize our proposed iterative reasoning framework that enables unbounded reasoning depth while maintaining a bounded memory footprint. Then, we detail a principled approach for reconstructing existing long-context reasoning datasets to conform to our paradigm.

### 2.1    **INFTYTHINK** Reasoning Paradigm

To address the computational challenges inherent in long-context reasoning, we propose **INFTY-THINK**, a novel paradigm that fundamentally reimagines how language models approach complex reasoning tasks. This paradigm decomposes complex reasoning into a series of bounded-length segments with intermediate summarization steps, enabling theoretically unlimited reasoning depth without the quadratic computational scaling of traditional approaches. Figure 2 illustrates the key differences between our approach and conventional reasoning. Below, we first formalize the conventional reasoning approach before presenting our iterative framework.

### 2.1.1    Conventional Reasoning Paradigm

Contemporary reasoning models, particularly those in the class of OpenAI o1 and similar architectures, rely on extended single-round generation for complex reasoning tasks. These models generate content comprising two principal components: a comprehensive "thinking" phase that captures the exploratory

reasoning process, followed by a "conclusion" phase that distills key insights into a structured response. This conventional reasoning paradigm can be formalized as:

$$\texttt{<user>}\text{Question}\texttt{<assistant>}\texttt{<think>}RP\texttt{</think>}\text{Conclusion}$$

where `<user>` and `<assistant>` demarcate the dialogue structure, `<think>` and `</think>` encapsulate the model's reasoning process $RP$, and the final Conclusion synthesizes the reasoning into a coherent answer.

This established approach, while effective for many problems, faces a fundamental limitation: as reasoning complexity increases, the token length of $RP$ grows substantially, often exceeding context window constraints and incurring quadratic computational costs. To address this limitation, we introduce **INFTYTHINK**, a paradigm that transforms monolithic reasoning into an iterative process with intermediate summarization steps.
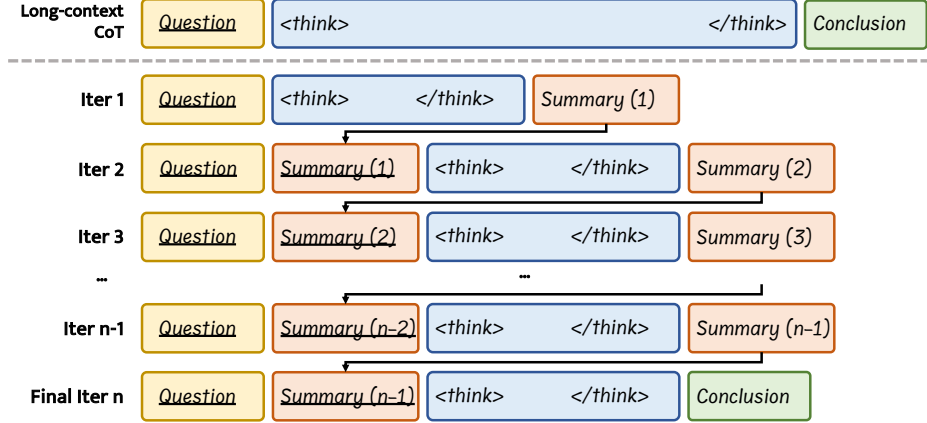


Figure 2: Illustration of **INFTYTHINK** versus vanilla long-context reasoning. **Upper panel:** Vanilla long-context reasoning generates continuous tokens until reaching maximum context length. **Lower panel:** Our **INFTYTHINK** approach divides reasoning into multiple iterations. The underlined segments represent content included in the prompt as model input, while non-underlined segments show model-generated output. Each iteration in **INFTYTHINK** consists of: (1) summarizing previous reasoning progress, (2) generating a focused reasoning segment within an efficient token budget, and (3) producing a concise progress summary. This iterative process enables arbitrarily deep reasoning chains without architectural modifications to the underlying model, while maintaining significantly lower computational complexity compared to traditional approaches.

### 2.1.2 Iterative Reasoning with Summarization: INFTYTHINK

In the **INFTYTHINK** framework, reasoning proceeds through multiple connected segments, each maintaining computational efficiency while preserving the coherent progression of thought. The initial reasoning iteration is formalized as:

$$\texttt{<user>}\text{Question}\texttt{<assistant>}\texttt{<think>}RP_1\texttt{</think>}\texttt{<summ>}S_1\texttt{</summ>}$$

where $RP_1$ represents the first segment of reasoning constrained to an efficient length, and $S_1$ denotes a concise summary of this segment encapsulated by the special tokens `<summ>` and `</summ>`. This summary serves as a compressed representation of the reasoning state, capturing essential information while discarding unnecessary details.

For subsequent iterations ($i > 1$), the model builds upon previous reasoning by incorporating the prior summary:

$$\texttt{<user>}\text{Question}\texttt{<assistant>}\texttt{<hist>}S_{i-1}\texttt{</hist>}\texttt{<think>}RP_i\texttt{</think>}\texttt{<summ>}S_i\texttt{</summ>}$$

where `<hist>` and `</hist>` delimit the previous summary $S_{i-1}$, which provides critical context for the current reasoning segment $RP_i$. Each iteration maintains a bounded token length while building upon accumulated knowledge through the summary mechanism.
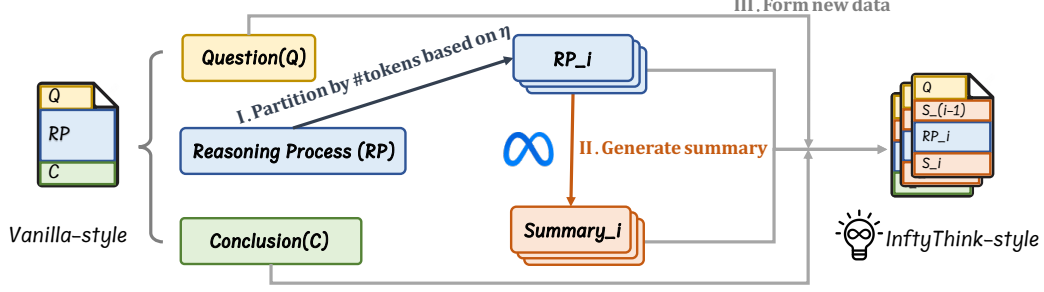
Figure 3: Systematic pipeline for reconstructing vanilla-style long-context reasoning data into the **INFTYTHINK**-style format. **I.** Original reasoning processes are partitioned into optimally sized fragments based on parameter ($\eta$), preserving semantic coherence. **II.** Meta-Llama-3.3-Instruct generates concise yet comprehensive summaries for each reasoning fragment. **III.** The original fragments and their generated summaries are systematically recombined to create **INFTYTHINK**-style training instances that teach the model to reason iteratively.

The final iteration ($n$) culminates in a conclusion rather than another summary:

<user>Question<assistant><hist>$S_{n-1}$</hist><think>$RP_n$</think>Conclusion

Throughout these expressions, blue denotes reasoning segments, orange represents intermediate summaries, and green indicates the final conclusion. This formulation elegantly handles edge cases: when problems are simple enough to be solved in a single iteration, the model bypasses summary generation, defaulting to the conventional paradigm.

During inference, the model iteratively generates reasoning segments and corresponding summaries, with each summary becoming the context for the subsequent iteration. This process continues until the model produces a conclusion instead of a summary, signaling completion of the reasoning task. To prevent potential infinite loops, we impose a hyperparameter `max_epochs` that terminates iteration if exceeded, though our empirical results indicate that well-trained models naturally converge within a reasonable number of iterations.

## 2.2 Data Reconstruction

While our **INFTYTHINK** paradigm offers a theoretically compelling approach to unbounded reasoning, it requires appropriate training data to enable models to learn this iterative reasoning process. Prior work has established that models can acquire sophisticated reasoning capabilities through supervised fine-tuning on data generated by highly capable reasoners. Building on this insight, we develop a principled methodology for transforming existing long-context reasoning datasets into our iterative format. We select OpenR1-Math [35][3] as our source dataset, which is a collection of mathematical reasoning generated by DeepSeek-R1 in response to questions from NuminaMath-1.5 [36]. This dataset spans a diverse spectrum of mathematical domains and difficulty levels, from elementary mathematics to competition-level problems, making it an ideal testbed for our approach. Our reconstruction pipeline comprises three key stages:

**Reasoning Segmentation** For each instance in the dataset, we partition the original reasoning process ($RP$) into segments based on a hyperparameter $\eta$ that determines the maximum token length of each segment. Rather than applying arbitrary truncation, we implement a semantically-aware segmentation algorithm: we first decompose the reasoning process into semantic units by identifying natural breakpoints at sentence or paragraph boundaries. These units are then tokenized and sequentially aggregated into segments, optimizing for coherence while ensuring each segment remains below the $\eta$ threshold. This process yields a sequence of reasoning segments $RP_1, RP_2, \ldots, RP_n$, formally expressed as:

$$\text{Partition}(RP, \eta) \Rightarrow RP_1, RP_2, \ldots, RP_n \tag{1}$$

---

[3]All data usage in this paper is in full compliance with the terms and conditions of the Apache License 2.0.

**Summary Generation**    For each reasoning segment, we generate a concise summary that captures its essential insights and progress toward the solution. To ensure information continuity across iterations, we employ a sophisticated foundation model $M$ (specifically Meta-Llama-3.3-70B-Instruct [37]) with carefully crafted prompting (provided in Appendix C):

$$S_i = \text{summarize}(M, RP_i, RP_1, \ldots, RP_{i-1}) \tag{2}$$

The summarization model receives not only the current reasoning segment $RP_i$ but also all preceding segments and their summaries, enabling it to create summaries that maintain reasoning continuity.

**Training Instance Construction**    From the segmented reasoning and generated summaries, we construct training instances that teach the model to perform iterative reasoning with summarization. These instances follow the structure of our **INFTYTHINK** paradigm:

$$D_i = \begin{cases} (\text{Question}, RP_1, S_1) & \text{for } i = 1, \\ (\text{Question}, S_{i-1}, RP_i, S_i) & \text{for } 1 < i < n, \\ (\text{Question}, S_{n-1}, RP_n, \text{Conclusion}) & \text{for } i = n. \end{cases} \tag{3}$$

For the initial reasoning step ($i = 1$), the model learns to generate the first reasoning segment followed by its summary. For intermediate steps ($1 < i < n$), it learns to continue reasoning based on previous summaries and generate new summaries. For the final step ($i = n$), it learns to produce a conclusive answer. This reconstruction process transforms each original example into $n$ training instances, where $n$ is the number of reasoning segments. The complete pipeline is illustrated in Figure 3. Applying this methodology to the OpenR1-Math dataset with $\eta$=4k, we expand the original 220K examples into 333K training instances, forming our OpenR1-Math-Inf dataset. This dataset enables models to learn the **INFTYTHINK** approach through supervised fine-tuning.

## 3    Experiments

### 3.1    Settings

We employ instruction fine-tuning to validate the proposed reasoning paradigm and associated dataset. Specifically, akin to the distilled model discussed in DeepSeek-R1 [2], training is conducted on five base models of varying sizes: Qwen2.5-Math-1.5B, Qwen2.5-Math-7B [38], Qwen2.5-14B, Qwen2.5-32B [39], and Meta-Llama-3.1-8B [37]. Instruction-based fine-tuning is applied using both OpenR1-Math and the newly introduced OpenR1-Math-Inf. The trained models are evaluated across multiple benchmarks, including MATH500 [40, 41], AIME24, and GPQA_diamond [42]. The detailed experimental setup is provided in the Appendix B.

### 3.2    Main Results

Table 1 presents our comprehensive evaluation of **INFTYTHINK** across five model architectures of varying scales and specializations. Several important patterns emerge from these results that provide insight into how our proposed reasoning paradigm affects model performance.

**Consistent Improvements Across Model Families and Scales.** Our **INFTYTHINK** consistently outperforms the vanilla reasoning approach across all model sizes and architectures. Notably, the improvements generalize beyond the Qwen architecture family to Meta-Llama-3.1-8B, demonstrating that the benefits of our iterative reasoning paradigm are not architecture-specific but rather represent a fundamental improvement in how models approach complex reasoning.

**Extended Reasoning Depth and Increased Output Throughput.** **INFTYTHINK** mitigates the computational overhead associated with the $O(n)$ complexity of LLMs at extended inference lengths by decomposing a single long generation into multiple shorter generation steps. This approach consistently enhances throughput across LLMs of varying parameter scales. Furthermore, the iterative mechanism inherent in **INFTYTHINK** allows models to efficiently handle extended reasoning tasks, maintaining high inference speed even as the generation length significantly increases. We further discuss this aspect in Appendix F.

**Scaling Trends with Model Size.** We observe an interesting relationship between model scale and the magnitude of improvement from **INFTYTHINK**. The relative gains are most pronounced in

| Base Model | Method | MATH500 | | | AIME24 | | | GPQA_diamond | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc | Tok | TPS | Acc | Tok | TPS | Acc | Tok | TPS |
| Qwen2.5-Math-1.5B | Vanilla | 74.58 | 5.20 | 0.83 | 12.92 | 13.25 | 2.64 | 22.82 | 9.68 | 2.40 |
| | INFTYTHINK | 80.75 | 7.09 | 1.21 | 22.08 | 25.96 | 2.75 | 35.67 | 10.18 | 2.67 |
| Qwen2.5-Math-7B | Vanilla | 88.88 | 4.02 | 0.66 | 26.46 | 12.93 | 2.36 | 41.76 | 8.43 | 1.80 |
| | INFTYTHINK | 91.65 | 4.90 | 0.69 | 40.00 | 22.25 | 2.67 | 51.86 | 9.93 | 1.89 |
| Qwen2.5-14B | Vanilla | 92.79 | 3.67 | 0.44 | 48.96 | 11.05 | 1.55 | 56.98 | 7.61 | 1.22 |
| | INFTYTHINK | 93.05 | 4.58 | 0.46 | 51.46 | 22.70 | 1.55 | 59.31 | 8.95 | 1.45 |
| Qwen2.5-32B | Vanilla | 95.91 | 3.38 | 0.30 | 55.63 | 10.22 | 0.90 | 64.14 | 6.80 | 0.71 |
| | INFTYTHINK | 95.99 | 3.85 | 0.37 | 62.50 | 17.39 | 0.94 | 65.62 | 7.55 | 0.82 |
| Meta-Llama-3.1-8B | Vanilla | 80.90 | 5.10 | 0.44 | 20.83 | 13.60 | 1.23 | 38.48 | 9.76 | 1.14 |
| | INFTYTHINK | 82.56 | 6.16 | 0.53 | 32.50 | 24.90 | 1.34 | 49.27 | 10.10 | 1.62 |

Table 1: Our main experimental results. The results are obtained by sampling the model 16 times with a temperature of 0.7. **Acc** stands for average accuracy(%), **Tok** stands for average number of generated tokens (K), and **TPS** stands for average number of tokens generated per second (K/s).

smaller models (e.g., 6.17%, 9.16%, and 12.85% improvements for Qwen2.5-Math-1.5B on the three benchmarks) and gradually diminish as model size increases, particularly on the MATH500 benchmark. This suggests that **INFTYTHINK** provides a form of algorithmic enhancement that partially compensates for limited model capacity, effectively allowing smaller models to perform more complex reasoning than their size would typically permit.

The iterative summarization mechanism in **INFTYTHINK** appears to effectively mitigate the limitations of traditional long-context reasoning by enabling more structured exploration of the solution space. The pattern of improvements suggests that our approach particularly benefits complex problems requiring multi-step reasoning, which are precisely the scenarios where long-context reasoning is most challenged by computational constraints. Our findings also suggest important implications for model scaling: **INFTYTHINK** may offer a more computationally efficient path to improved reasoning capabilities than simply scaling model size, particularly for smaller models where the relative improvements are most pronounced. This could have significant practical implications for deploying advanced reasoning capabilities in resource-constrained environments.

## 4 Analysis

### 4.1 Endowing Short-context Models with Long-context Reasoning Ability

Many foundational LLMs are pretrained with limited context windows (4k or 8k tokens), yet a significant portion of reasoning datasets exceeds these boundaries. Analysis of OpenR1-Math shows only 54% of samples contain fewer than 4k tokens, and 83% are within 8k tokens (Appendix D), revealing a critical mismatch between model architecture and reasoning requirements.

Table 1 shows consistent performance improvements when applying **INFTYTHINK** across all model configurations. The gains are particularly notable on complex benchmarks like AIME24 and GPQA_diamond, where problems typically require longer reasoning chains that would exceed standard context windows. For instance, Qwen2.5-Math-7B achieves a 13.54% improvement on AIME24 and a 10.1% improvement on GPQA_diamond using our approach. These improvements suggest that **INFTYTHINK** effectively addresses context length limitations by restructuring long reasoning into manageable segments with summarization.

To validate our approach against alternative context extension methods, we implemented RoPE positional encoding interpolation [43–45], a common technique for extending context windows beyond pretraining lengths. While this approach yielded modest improvements, it consistently underperformed compared to **INFTYTHINK** across all benchmarks (detailed results in Appendix E). This comparison is particularly revealing: rather than attempting to stretch architectural limitations

through embedding manipulation, **INFTYTHINK** fundamentally restructures the reasoning process itself to work within existing constraints.

These findings suggest that **INFTYTHINK** offers a more effective solution to the long-context reasoning challenge than traditional context window extension techniques. By allowing models to periodically summarize and build upon previous reasoning, our approach enables more flexible and adaptable reasoning capabilities that aren't bound by fixed architectural constraints. This has important implications for deploying reasoning systems in environments where context length would otherwise be a limiting factor.

## 4.2 Influence of Context Window Size Parameter $\eta$

Parameter $\eta$ plays a crucial role in **INFTYTHINK**, controlling the maximum token length for each reasoning iteration. This parameter creates a fundamental tradeoff: larger values reduce the number of iterations but increase per-iteration computational cost, while smaller values distribute computation more evenly but potentially fragment reasoning.

Table 2 presents performance across three $\eta$ values (2k, 4k, and 6k tokens) on our benchmarks. Surprisingly, all configurations consistently outperform the baseline with no clear optimal value across all datasets. On MATH500, performance increases marginally with $\eta$, suggesting that longer uninterrupted reasoning benefits simpler problems. In contrast, for AIME24, smaller $\eta$ values yield better results (17.91% improvement at $\eta$=2k versus 12.91% at $\eta$=6k), indicating that more frequent summarization helps with complex reasoning.

| Method | $\eta$ | MATH500 | AIME24 | GPQA_D |
|--------|--------|---------|--------|--------|
| Vanilla | / | 88.88 | 26.46 | 41.76 |
| InftyThink | 2k | 91.60+2.72 | 44.37+17.91 | 51.89+10.13 |
| | 4k | 91.65+2.77 | 40.00+13.54 | 51.86+10.1 |
| | 6k | 92.26+3.38 | 39.37+12.91 | 49.84+8.08 |

Table 2: Evaluation results across different $\eta$. GPQA_D refers to GPQA_diamond. Experiments are conducted on Qwen2.5-Math-7B.

These results challenge the intuition that fragmented reasoning necessarily harms performance. Even with $\eta$=2k, where reasoning is interrupted every 2,000 tokens, **INFTYTHINK** maintains or improves performance across all benchmarks. This suggests that well-designed summarization mechanisms effectively preserve critical information while discarding redundant computation.

The robustness to different $\eta$ values demonstrates that **INFTYTHINK**'s benefits derive primarily from its iterative summarization approach rather than specific segmentation boundaries. This flexibility allows practitioners to select $\eta$ values based on hardware constraints or specific application requirements without significant performance penalties.

## 4.3 Performance across Reasoning Iteration Rounds

A defining characteristic of **INFTYTHINK** is its ability to transcend the constraints of maximum context length through iterative reasoning with summarization. Figure 4 quantifies this capability by tracking performance across reasoning iterations and comparing it against traditional reasoning methods with fixed token limits. The results reveal three key insights:

**Progressive Performance Improvement** Unlike traditional reasoning approaches that hit a performance ceiling determined by their maximum token limit, **INFTYTHINK** enables continuous improvement through successive iterations. On AIME24, a challenging benchmark, performance steadily increases from iterations 1 through 10 for all ($\eta$) settings, demonstrating that complex problems benefit substantially from extended reasoning beyond conventional context limits.

**Efficiency of Iterative Summarization** Even with the smaller context setting of $\eta$=2k, **INFTYTHINK** eventually reaches comparable performance to larger $\eta$ values across all benchmarks. This is particularly evident on AIME24, where by iteration 10, the $\eta$=2k configuration approaches the performance of $\eta$=6k despite using reasoning segments one-third the size. This demonstrates that effective summarization can preserve critical reasoning information even with frequent compression.
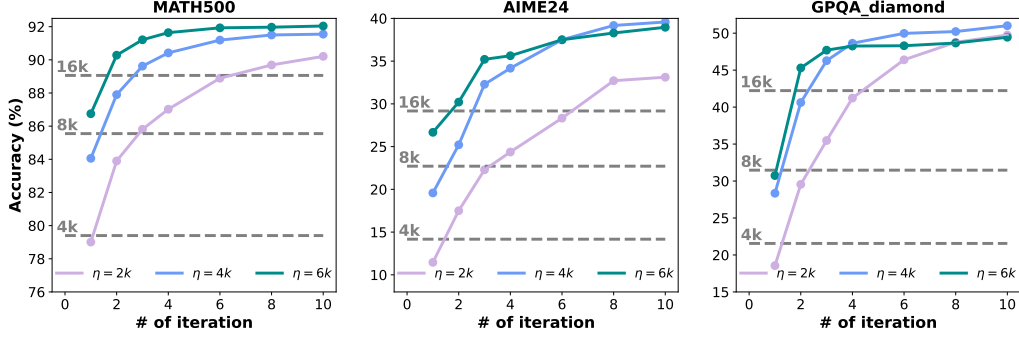
Figure 4: Model performance across iterations compared to traditional reasoning. Horizontal lines represent conventional long-context reasoning with different `max_new_tokens` settings (4k, 8k, 16k). Curves show **INFTYTHINK**'s accuracy evolution across iterations under different $\eta$ settings (2k, 4k, 6k). **INFTYTHINK** rapidly surpasses fixed-length reasoning constraints, with performance continuing to improve beyond traditional token limits. Experiments conducted on Qwen2.5-Math-7B.

**Early-Stage Performance Tradeoffs**   Models with larger $\eta$ values (6k) consistently outperform those with smaller segments (2k) in early iterations across all benchmarks. However, this advantage diminishes and sometimes reverses in later iterations, particularly on GPQA_diamond where $\eta$=4k eventually surpasses $\eta$=6k. This suggests that while larger segments provide initial advantages, they may commit the model to reasoning paths that become difficult to revise, whereas smaller segments allow more flexible exploration over multiple iterations.

## 5 Discussion

**Alignment with Human Reasoning**   **INFTYTHINK**'s iterative reasoning approach shares interesting parallels with human problem-solving strategies. Humans rarely solve complex problems through a single, exhaustive thought process but instead work through incremental steps, summarizing intermediate progress, and building on previous insights. The strong performance of our approach, particularly on complex problems, suggests that structuring AI reasoning to better align with these natural problem-solving patterns may yield both efficiency and effectiveness benefits. This connection between iterative reasoning with summarization and improved performance offers potential insights for developing more effective AI reasoning systems.

**Adaptive Reasoning Depth**   Unlike conventional reasoning approaches with fixed computational budgets, **INFTYTHINK** can adaptively allocate computational resources based on problem difficulty. Our analysis in Section 4.4 shows that simpler problems (e.g., in MATH500) reach ceiling performance with fewer iterations, while more complex problems (in AIME24 and GPQA_diamond) benefit from extended reasoning. This adaptive depth capability has important implications for practical deployment, as it enables efficient resource allocation across heterogeneous problem sets without requiring predetermined computation limits.

**Limitations**   Despite its advantages, **INFTYTHINK** faces several limitations. First, the quality of reasoning depends heavily on the model's summarization capabilities—poor summarization can lead to information loss that hinders subsequent reasoning. Second, breaking reasoning into segments might disrupt the coherent flow of thought for certain problem types that benefit from maintaining a complete chain of reasoning. Finally, while our approach reduces computational complexity, it introduces additional inference steps that may increase latency in time-sensitive applications.

**Future Directions**   Several promising directions could extend this work. First, integrating reinforcement learning techniques such as GRPO could help models better learn when and what to summarize, potentially improving information retention across iterations. Second, exploring variable-length reasoning segments that adapt based on problem complexity could further optimize the tradeoff between computational efficiency and reasoning coherence. Third, applying **INFTYTHINK** to multi-modal

reasoning tasks could expand its applicability to domains requiring integration of visual, textual, and numerical reasoning. Finally, investigating how to parallelize different reasoning paths within the INFTYTHINK framework could further accelerate complex problem-solving.

## 6 Conclusion

In this paper, we introduced INFTYTHINK, a novel reasoning paradigm that transforms monolithic long-context reasoning into an iterative process with periodic summarization. By generating partial reasoning, summarizing current understanding, and building upon these summaries in subsequent iterations, our approach effectively addresses the quadratic computational complexity and context length limitations of conventional approaches. Experiments across multiple model architectures demonstrate consistent performance and throughput improvements on challenging reasoning benchmarks . Our analysis confirms that INFTYTHINK not only reduces computational costs but also enables models to transcend their native context window limitations without architectural modifications. This paradigm represents a step toward more cognitively plausible AI reasoning through iterative refinement rather than exhaustive single-pass analysis, opening promising avenues for more efficient, flexible reasoning in language models that decouples reasoning depth from computational complexity.

## References

[1] OpenAI. Introducing openai o1. https://openai.com/index/introducing-o1-preview/, 2024.

[2] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, January 2025.

[3] Google DeepMind. Gemini flash thinking. https://deepmind.google/technologies/gemini/flash-thinking/, 2025.

[4] Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown. https://qwenlm.github.io/blog/qwq-32b-preview/, November 2024.

[5] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, et al. Kimi k1.5: Scaling reinforcement learning with llms, January 2025.

[6] Ehsan Latif, Yifan Zhou, Shuchen Guo, Lehong Shi, Yizhu Gao, et al. Can openai o1 outperform humans in higher-order cognitive thinking?, December 2024.

[7] Ehsan Latif, Yifan Zhou, Shuchen Guo, Yizhu Gao, Lehong Shi, et al. A systematic assessment of openai o1-preview for higher order thinking in education, October 2024.

[8] Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Alphamath almost zero: Process supervision without process, September 2024.

[9] Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training, February 2024.

[10] Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, et al. rstar-math: Small llms can master math reasoning with self-evolved deep thinking, January 2025.

[11] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, et al. S1: Simple test-time scaling, March 2025.

[12] Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning, February 2025.

[13] Open Thoughts Team. Open thoughts. open-thoughts, January 2025.

[14] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, February 2024.

[15] Zhiyuan Zeng, Qinyuan Cheng, Zhangyue Yin, Bo Wang, Shimin Li, et al. Scaling of search and learning: A roadmap to reproduce o1 from reinforcement learning perspective, December 2024.

[16] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, et al. Evaluation of openai o1: Opportunities and challenges of agi, September 2024.

[17] OpenAI. Learning to reason with llms. https://openai.com/index/learning-to-reason-with-llms/, 2024.

[18] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling, July 2024.

[19] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, August 2024.

[20] Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, et al. Thoughts are all over the place: On the underthinking of o1-like llms, February 2025.

[21] Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack, November 2024.

[22] Mo Li, Songyang Zhang, Yunxin Liu, and Kai Chen. Needlebench: Can llms do retrieval and reasoning in 1 million context window?, July 2024.

[23] Tao Yuan, Xuefei Ning, Dong Zhou, Zhijie Yang, Shiyao Li, et al. Lv-eval: A balanced long-context benchmark with 5 length levels up to 256k, October 2024.

[24] Jianhui Pang, Fanghua Ye, Derek Fai Wong, Xin He, Wanshun Chen, and Longyue Wang. Anchor-based large language models, June 2024.

[25] Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. Token-budget-aware llm reasoning, February 2025.

[26] Jintian Zhang, Yuqi Zhu, Mengshu Sun, Yujie Luo, Shuofei Qiao, et al. Lightthinker: Thinking step-by-step compression, February 2025.

[27] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks, April 2024.

[28] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, et al. H$_2$o: Heavy-hitter oracle for efficient generative inference of large language models, December 2023.

[29] Guoxuan Chen, Han Shi, Jiawei Li, Yihang Gao, Xiaozhe Ren, et al. Sepllm: Accelerate large language models by compressing one segment into one separator, December 2024.

[30] Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. C3ot: Generating shorter chain-of-thought without compromising effectiveness, December 2024.

[31] Daman Arora and Andrea Zanette. Training language models to reason efficiently, February 2025.

[32] Tengxiao Liu, Qipeng Guo, Xiangkun Hu, Cheng Jiayang, Yue Zhang, Xipeng Qiu, and Zheng Zhang. Can language models learn to skip steps?, November 2024.

[33] Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. Tokenskip: Controllable chain-of-thought compression in llms, February 2025.

[34] Xinyin Ma, Guangnian Wan, Runpeng Yu, Gongfan Fang, and Xinchao Wang. Cot-valve: Length-compressible chain-of-thought tuning, February 2025.

[35] open-r1. Openr1-math-220k. https://huggingface.co/datasets/open-r1/OpenR1-Math-220k, February 2025.

[36] Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, et al. Numinamath-1.5. https://huggingface.co/datasets/AI-MO/NuminaMath-1.5, February 2025.

[37] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, et al. The llama 3 herd of models, November 2024.

[38] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, et al. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, September 2024.

[39] Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, et al. Qwen2.5 technical report, January 2025.

[40] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, et al. Measuring mathematical problem solving with the math dataset, November 2021.

[41] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, et al. Let's verify step by step, May 2023.

[42] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, et al. Gpqa: A graduate-level google-proof q&a benchmark, November 2023.

[43] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, November 2023.

[44] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of large language models via positional interpolation, June 2023.

[45] kaiokendev. Things i'm learning while training superhot. https://kaiokendev.github.io/til, 2023.

[46] Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, Findings of the Association for Computational Linguistics: ACL 2023, pages 1049–1065, Toronto, Canada, July 2023. Association for Computational Linguistics.

[47] Jiankai Sun, Chuanyang Zheng, Enze Xie, Zhengying Liu, Ruihang Chu, et al. A survey of reasoning with foundation models, January 2024.

[48] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, et al. Deepseek-coder: When the large language model meets programming – the rise of code intelligence, January 2024.

[49] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, et al. Code llama: Open foundation models for code, January 2024.

[50] Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions, July 2024.

[51] Yuchen Yan, Jin Jiang, Yang Liu, Yixin Cao, Xin Xu, et al. S^3cmath: Spontaneous step-level self-correction makes large language models better mathematical reasoners, February 2025.

[52] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, et al. Metamath: Bootstrap your own mathematical questions for large language models, May 2024.

[53] Jin Jiang, Yuchen Yan, Yang Liu, Yonggang Jin, Shuai Peng, et al. Logicpro: Improving complex logical reasoning via program-guided learning, February 2025.

[54] Jeffrey Cheng and Benjamin Van Durme. Compressed chain of thought: Efficient reasoning through dense representations, December 2024.

[55] Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, et al. Scaling up test-time compute with latent reasoning: A recurrent depth approach, February 2025.

[56] 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, December 2024.

[57] Xuan Shen, Yizhou Wang, Xiangxi Shi, Yanzhi Wang, Pu Zhao, and Jiuxiang Gu. Efficient reasoning with hidden thinking, January 2025.

[58] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, et al. Efficient memory management for large language model serving with pagedattention, September 2023.

[59] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, et al. Sglang: Efficient execution of structured language model programs, June 2024.

# Contents

# A   Related Works

## A.1   Reasoning of Large Language Models

Reasoning ability is one of the fundamental competencies of large language models(LLMs), reflecting their capacity to tackle complex challenges in the human domain. Currently, LLMs demonstrate impressive performance across various reasoning tasks, including commonsense reasoning, mathematical reasoning, code reasoning, logical reasoning and etc. The reasoning capabilities of these models can be enhanced at several stages during training [46, 47]. For instance, during pre-training, the inclusion of extensive reasoning-related knowledge and examples helps the model to learn reasoning patterns from the data [14, 48, 49]. Similarly, in the supervised fine-tuning phase, incorporating high-quality reasoning question-answer pairs can further refine the model's reasoning patterns and enhance its capabilities [50–53]. In the reinforcement learning phase, the model's reasoning is monitored and guided through feedback on outcome or processes, providing additional improvements [2, 14]. The release of OpenAI's o1 [1] marked a significant breakthrough in the reasoning abilities of LLMs. OpenAI o1 demonstrated long-context reasoning capabilities, where the model utilized extended chains of thought to integrate planning, self-correction, and other cognitive functions, significantly boosting its reasoning performance [15]. More recently, DeepSeek-R1, an open-source o1-like reasoning model, has exhibited comparable reasoning abilities. Furthermore, distilled data from DeepSeek-R1 enables smaller LLMs to also acquire long-context reasoning skills [2, 12, 13].

## A.2   Compression of LLM's Reasoning Process

Current research on compressing the Chain-of-Thought (CoT) process to accelerate large language model (LLM) inference is primarily categorized into two approaches: CoT token compression and KV cache compression. CoT token compression enhances inference efficiency by reducing the number of tokens generated by the model. This approach can be further subdivided into discrete text token compression and continuous latent token compression methods. Discrete text token compression employs straightforward strategies such as prompt engineering [25], instruction fine-tuning [30], and reinforcement learning [31] to train models to produce more concise reasoning processes. Within this category, the skip-tokens method [32, 33] enables the model to intelligently skip unimportant tokens during inference, thereby achieving acceleration. In contrast, continuous latent token compression [54–57] explores a more innovative approach by attempting to compress reasoning steps into continuous latent representations. This allows LLMs to perform effective inference without explicitly generating discrete word tokens. On the other hand, KV cache compression optimizes inference performance by reducing the storage requirements and computational load of the KV cache. This approach mainly includes two types of methods: training-free and training-based. Training-free KV cache management strategies enhance efficiency by selectively retaining key tokens. The criteria for selection include prioritizing initial and most recent tokens for their temporal relevance [27], identifying tokens with significant historical attention [28], or selecting tokens based on structural cues such as punctuation marks [29]. Training-based KV cache management [24, 26] involves introducing special tokens and training LLMs to compress important historical information into these tokens, thereby achieving KV cache merging. This method instructs the model on when to perform compression during the training phase and applies corresponding interventions during the inference phase.

# B   Experiment Settings

For the training process, we utilize the Megatron-LM framework. The supervised fine-tuning is performed for 3 epochs, with a maximum sequence length of 16,384 tokens. The batch size is set to 32, the initial learning rate is 1e-5, and the warmup ratio is set at 0.03. The learning rate follows a cosine decay schedule to reach zero. To accelerate training, we pack all SFT samples to the maximum sequence length. Each packed sample retains its original positional embeddings, and attention values are computed independently for each instance. All experiments are conducted on 256 Ascend H910B-64G NPUs.

For models trained on OpenR1-Math, we conduct standard single-round inference with a maximum output length of 16,384 tokens. For models trained on OpenR1-Math-Inf, we apply the proposed **INFTYTHINK** reasoning paradigm, performing multi-round iterative inference with a maximum of 50 epochs and a single-round maximum reasoning length of 8,192 tokens. To mitigate potential

fluctuations in the evaluation results, each evaluation case is sampled 16 times with a temperature setting of 0.7, and the average accuracy is computed. All inferences are executed using *vLLM* [58] v1-engine on NVIDIA A100-80G GPUs. For models with 1.5B parameters, inference is performed on 1 GPU, for 7B/8B models, inference is performed on 2 GPUs, while for models with 14B and 32B parameters, inference is performed on 4 GPUs.

## C  Prompts for Summary Generation

> **Prompts for Summary Generation**
>
> <|begin_of_text|><|start_header_id|>system<|end_header_id|>
>
> Cutting Knowledge Date: December 2023
> Today Date: 26 Jul 2024
>
> <|eot_id|><|start_header_id|>user<|end_header_id|>
>
> **{question}**<|eot_id|><|start_header_id|>assistant<|end_header_id|>
>
> **{reasoning_process}**<|eot_id|><|start_header_id|>user<|end_header_id|>
>
> **Please list what you have achieved in your last response. Note that you should only output the summarization. You should list all the key steps and important intermediate conclusion. Please list them with '*'.** <|eot_id|><|start_header_id|>assistant<|end_header_id|>
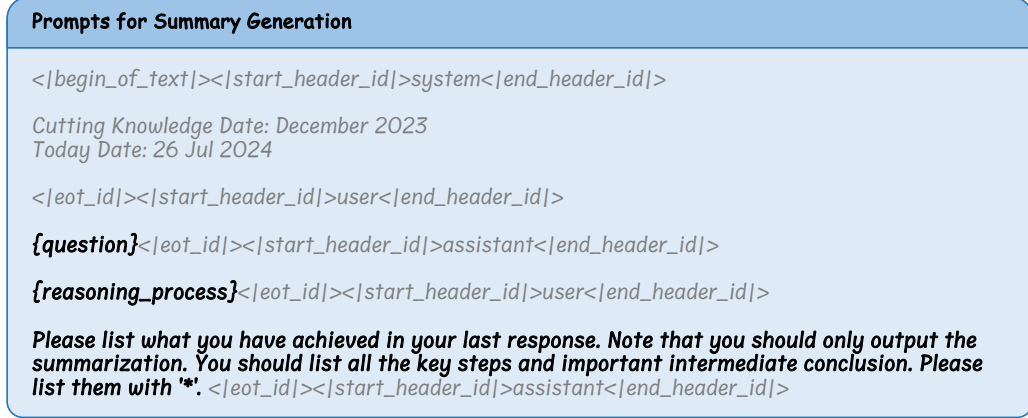
Figure 5: Prompt for generating a summary of a reasoning process fragment. A multi-turn dialogue approach is employed to generate the summary. The light-colored section in the figure represents the chat template, while the dark-colored section corresponds to the input we designed.
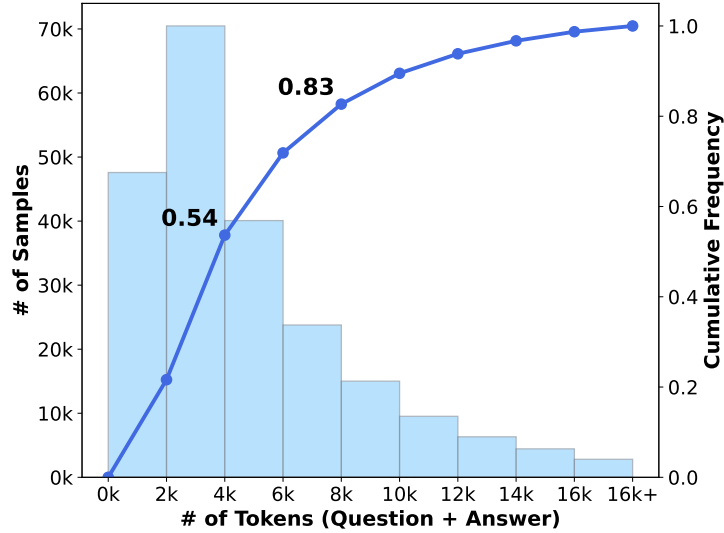
## D  Token Length Distribution of OpenR1-Math



Figure 6: Token distribution of OpenR1-Math. The statistics are obtained using the tokenizer of Qwen2.5-Math-7B.

# E  Experiment Details for RoPE-Scaled Models

We apply linear interpolation to RoPE with a scale factor of 8, extending the context window of Qwen2.5-Math-7B from 4k to 32k. The model is then trained using the same methodology as in the main experiment. The results are presented in Table 3.

| Data | Positional Embedding | MATH500 | AIME24 | GPQA_diamond |
|------|----------------------|---------|--------|--------------|
| OpenR1-Math | raw | 88.88 | 26.46 | 41.76 |
| | linear scale to 32k | 90.91+2.03 | 30.63+4.17 | 48.26+6.5 |
| OpenR1-Math-Inf | raw | 91.65+2.77 | 40.00+13.54 | 51.86+10.1 |

Table 3: Comparison results of RoPE linear interpolation experiments (%).

# F  Discussion about Inference-time Computational Cost

Contemporary LLMs face a fundamental efficiency bottleneck due to the quadratic ($O(n^2)$) computational scaling of attention with sequence length. **INFTYTHINK** addresses this by decomposing reasoning into shorter segments with periodic summarization.

Figure 7 demonstrates that **INFTYTHINK** (red line) consistently achieves higher accuracy than traditional reasoning (gray line) under equivalent computational budgets. This advantage increases with problem complexity, becoming most prominent on AIME24 and GPQA_diamond benchmarks. Simultaneously, the blue line shows that **INFTYTHINK** makes more efficient use of each token, particularly for complex problems where traditional approaches struggle to maintain performance scaling. By avoiding the quadratic growth pattern of traditional reasoning, **INFTYTHINK** fundamentally improves the relationship between computation and reasoning performance, offering a promising direction for deploying advanced reasoning in resource-constrained environments.
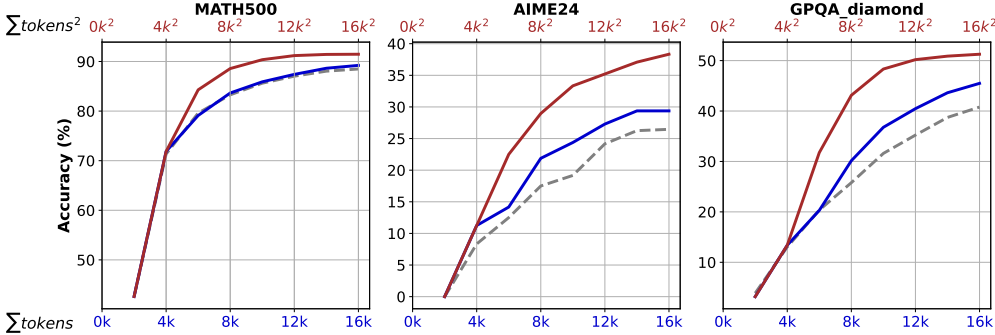


Figure 7: Accuracy across various benchmarks under different computational scales. The gray line represents traditional long-context reasoning trained on OpenR1-Math. The two colored lines correspond to **INFTYTHINK**, with the blue line indicating the total number of tokens computed by the model and the red line representing the squared sum of tokens computed across multiple inference iterations. The gray line can simultaneously represent the effects of traditional long-context reasoning in both dimensions. Comparing the gray line with the blue line illustrates the accuracy trend as the model reasons over a certain number of tokens, while the comparison between the gray line and the red line reflects the relationship between computational cost (with O(n²) complexity) and accuracy. Experiments are conducted on Qwen2.5-Math-7B.

To quantitatively analyze the computational cost of our proposed **INFTYTHINK** during inference, we introduce two key metrics: the total token count and the sum of squared token counts. Specifically, for an iterative generation process with $n$ iterations, we define the token count at each step as:

$$\text{Tokens}_i = \big|\text{tokenize}(\text{Question})\big| + \big|\text{tokenize}(S_{i-1})\big| + \big|\text{tokenize}(RP_i)\big| + \text{tokenize}(S_i)$$

where $\left|\text{tokenize}(x)\right|$ indicates the number of tokens after tokenization of string $x$. In particular, the number of tokens generated during the first inference step is:

$$\text{Tokens}_1 = \left|\text{tokenize}(\text{Question})\right| + \left|\text{tokenize}(RP_i)\right| + \left|\text{tokenize}(S_1)\right|$$

As the final inference step generate a conclusion instead of a summary, the token count during the final inference step is defined as:

$$\text{Tokens}_n = \left|\text{tokenize}(\text{Question})\right| + \left|\text{tokenize}(S_{i-1})\right| + \left|\text{tokenize}(RP_i)\right| + \left|\text{tokenize}(\text{Conclusion})\right|$$

The first metric, the total sum of tokens, is defined as:

$$\sum \text{Tokens} = \sum \text{Tokens}_i, i \in [1, n]$$

The second metric, the sum of squared token counts, is defined as:

$$\sum \text{Tokens}^2 = \sum \text{Tokens}_i^2, i \in [1, n]$$

For a standard long-context reasoning task with a single generation, where $n = 1$, the relationship $\left(\sum \text{Tokens}\right)^2 = \sum \text{Tokens}^2$ holds.

In Figure 7, we illustrate the relationship between these metrics by employing a dual-axis design. The lower axis (colored in blue) tracks the first metric, $\sum \text{Tokens}$, while the upper axis (colored in red) represents the second metric, $\sum \text{Tokens}^2$, with its scale being the square of the lower axis.

For traditional long-context reasoning, the theoretical relationship $\left(\sum \text{Tokens}\right)^2 = \sum \text{Tokens}^2$ holds, allowing us to depict both metrics using a single curve, shown as a gray line. In contrast, for **INFTYTHINK**, we differentiate the two metrics by employing distinct lines, each colored to correspond with its respective axis.

To plot the curve shown in the figure, we calculate the number of correct instances at specific token thresholds. Specifically, we set eight token thresholds: 2k, 4k, 6k, 8k, 10k, 12k, 14k, and 16k. We select all the correct completions from the evaluation, tokenize them, and then count how many samples fall under each of these thresholds. The accuracy is computed by dividing the number of samples for each threshold by the total number of completions.

There are two ways to analyze this figure. The first approach is to compare the accuracy for the same computational cost, by fixing the value of x and comparing the corresponding y values. The second approach is to compare the computational cost for the same accuracy, by fixing the value of y and comparing the corresponding x values.

We would like to emphasize that the computational complexity calculations provided above were based on the most rigorous methodology. However, in practical applications, inference frameworks like vLLM [58] and sglang [59] already support prefix-caching, which eliminates the need to recompute the attention values of the question during each inference. Despite this, under the strictest computational model, **INFTYTHINK** demonstrates superior efficiency, underscoring the effectiveness of the proposed approach.

### F.1 Computational Cost across Different $\eta$

In order to compare the trade-off between computational cost and performance at different values of $\eta$, we also plotted the variations of these two metrics with accuracy for different $\eta$ values, as shown in Figure 8 and 9. Specifically, the choice of $\eta$ demonstrates a clear trade-off with performance. Smaller $\eta$ values lead to higher reasoning efficiency in the early stages, whereas larger $\eta$ values result in improved reasoning performance. Based on our observations, among the comparisons of $\eta = 2k$, $\eta = 4k$, and $\eta = 6k$, $\eta = 2k$ strikes the optimal balance between these factors.
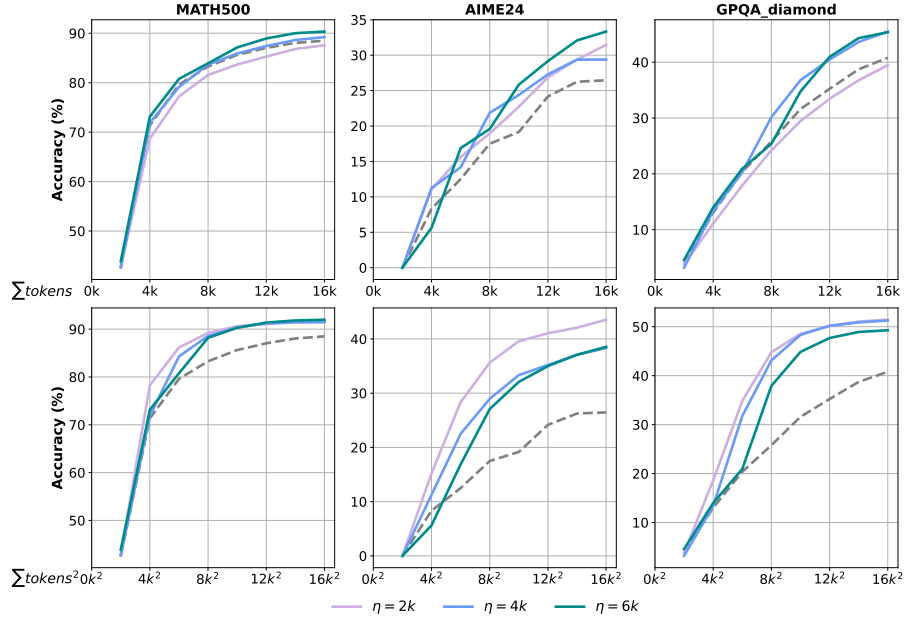
Figure 8: Accuracy(%) across various benchmarks under different computational scales on different $\eta$ settings. The three subplots above illustrate the relationship between $\sum$ Tokens and accuracy, while the three subplots below depict the relationship between $\sum$ Tokens$^2$ and accuracy. Experiments are conducted on Qwen2.5-Math-7B.



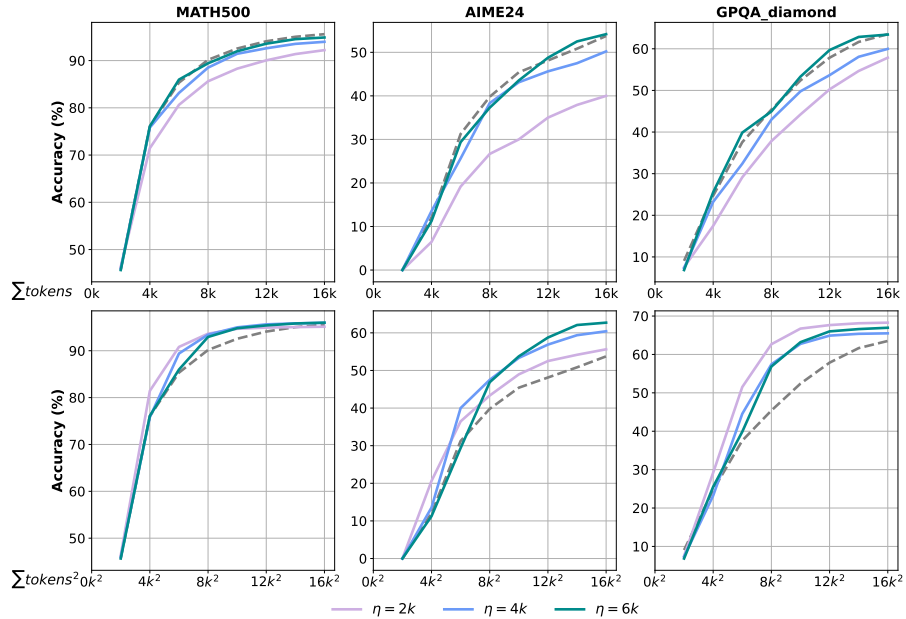Figure 9: Accuracy(%) across various benchmarks under different computational scales on different $\eta$ settings. The three subplots above illustrate the relationship between $\sum$ Tokens and accuracy, while the three subplots below depict the relationship between $\sum$ Tokens$^2$ and accuracy. Experiments are conducted on Qwen2.5-32B.