

VTC-R1: Vision-Text Compression for Efficient Long-Context Reasoning

Yibo Wang¹ Yongcheng Jing¹ Shunyu Liu¹ Hao Guan¹ Rong-Cheng Tu¹
Chengyu Wang² Jun Huang² Dacheng Tao¹

Abstract

Long-context reasoning has significantly empowered large language models (LLMs) to tackle complex tasks, yet it introduces severe efficiency bottlenecks due to the computational complexity. Existing efficient approaches often rely on complex additional training or external models for compression, which limits scalability and discards critical fine-grained information. In this paper, we propose VTC-R1, a new efficient reasoning paradigm that integrates vision-text compression into the reasoning process. Instead of processing lengthy textual traces, VTC-R1 renders intermediate reasoning segments into compact images, which are iteratively fed back into vision-language models as "optical memory." We construct a training dataset based on OpenR1-Math-220K achieving $3.4\times$ token compression and fine-tune representative VLMs—Glyph and Qwen3-VL. Extensive experiments on benchmarks such as MATH500, AIME25, AMC23 and GPQA-D demonstrate that VTC-R1 consistently outperforms standard long-context reasoning. Furthermore, our approach significantly improves inference efficiency, achieving $2.7\times$ speedup in end-to-end latency, highlighting its potential as a scalable solution for reasoning-intensive applications. Our code is available at <https://github.com/w-yibo/VTC-R1>.

1. Introduction

Reasoning capability (Li et al., 2025b; Lightman et al., 2023; Yao et al., 2023; Huang & Chang, 2023; Yao et al., 2025) has emerged as a powerful technique of large language models (LLMs), enabling them to tackle complex tasks such as mathematical problem solving (Hendrycks et al., 2021; Luo et al., 2025a; Hu et al., 2025) and code generation (Chen et al., 2021; Jiang et al., 2024). Recent advancements, exemplified by OpenAI o1 (OpenAI, 2024) and DeepSeek-

¹Nanyang Technical University ²Alibaba Cloud Computing.

Preprint.

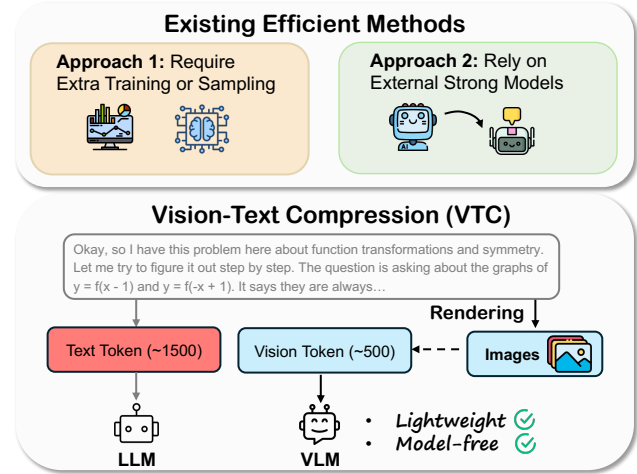


Figure 1. Comparison between existing efficient reasoning approaches and vision-text compression (VTC). Existing methods either require additional training or sampling procedures, or rely on external strong models. In contrast, VTC leverages lightweight rendering to transform long textual reasoning traces into compact visual representations, enabling VLMs to encode information with significantly fewer vision tokens ($3\text{-}4\times$ compression). This approach is both lightweight and model-free.

R1 (Guo et al., 2025), leverage reinforcement learning to further scale this capability to long-context reasoning, substantially improving performance on challenging real-world tasks (Wang et al., 2026b). Despite recent progress, long-context reasoning introduces severe efficiency bottlenecks. The computational complexity of the transformer architecture (Zaheer et al., 2020; Beltagy et al., 2020; Kitaev et al., 2020) grows quadratically with sequence length, causing both computation and memory costs to increase rapidly as the context expands. This leads to degraded inference speed, reduced training efficiency, and limited scalability, which significantly hinders real-world deployment.

To mitigate these issues, several efficient approaches are proposed (Chen et al., 2025; Munkhbat et al., 2025; Lee et al., 2025; Liu et al., 2024). Existing methods can be broadly categorized into two groups. i) Extra training or sampling stages beyond standard training. For example, CoT-Valve (Ma et al., 2025b) adopts a multi-stage training procedure to obtain models specialized for different reasoning lengths and O1-Pruner (Luo et al., 2025b) applies offline

reinforcement learning with multiple sampled trajectories (16 responses per problem). These approaches increase training and inference cost. ii) External strong models to guide reasoning compression. TokenSkip (Xia et al., 2025) requires an additional model to estimate token importance, while R1-Compress (Wang et al., 2025) and InfyThink (Yan et al., 2025) depend on powerful external summarization models (e.g., Llama-3.3-70B-Instruct) to condense long reasoning traces. Although both categories of methods are effective, they often restrict exploration space and discard fine-grained information that is critical for reasoning.

Without additional training or external models, how can we achieve efficient reasoning while preserving fine-grained information?

Motivated by this, a promising yet underexplored direction is vision-text compression (VTC) (Wei et al., 2025; Cheng et al., 2025; Xing et al., 2025b; Zhao et al., 2025; Xing et al., 2025a). Rather than reducing fine-grained information, VTC adopts an alternative representation by transforming textual content into visual forms via lightweight rendering, enabling vision-language models (VLMs) to encode rich semantic information using substantially fewer vision tokens. This design is lightweight and model-free, as shown in Figure 1, introducing no additional training stages or reliance on external compression models. Prior works such as DeepSeek-OCR (Wei et al., 2025) and Glyph (Cheng et al., 2025) focus on text reconstruction or long-context understanding, showing that long text sequences can be represented with $3\times$ – $10\times$ token compression while maintaining high decoding precision. However, whether such high-density visual representations can preserve and support multi-step reasoning processes remains unclear. Notably, mathematical reasoning, with its symbolic structure and step-wise derivations, is naturally amenable to visual rendering, making it a suitable and principled testbed for studying reasoning-oriented vision-text compression.

To bridge this gap, we propose VTC-R1, a new efficient reasoning paradigm that iteratively integrates vision-text compression into long-context reasoning. VTC-R1 treats the reasoning process as multiple processes, where *the preceding process are regarded as long-context and rendered into compact images*, and performs iterative reasoning (Yan et al., 2025) with VLMs. As illustrated in Figure 2, the reasoning process is decomposed into a sequence of reasoning steps. Upon the completion of each step, it is rendered into an image. To proceed to the next step, the accumulated images of previous steps are fed back into the model alongside the question, functioning as a form of optical memory that compactly encodes previous reasoning using vision tokens.

We construct a training dataset based on OpenR1-Math-220K (Hugging Face, 2025), a large-scale long-context reasoning corpus generated by DeepSeek-R1 (Guo et al., 2025).

We segment each long reasoning trace into shorter reasoning segments and render the preceding segments into images, forming paired image-text reasoning data with up to $3.4\times$ token compression as shown in Table 1. We then fine-tune representative VTC-VLM (i.e., Glyph (Cheng et al., 2025)) and the state-of-the-art VLM (i.e., Qwen3-VL (Bai et al., 2025)), under this iterative reasoning framework. Extensive experiments on diverse mathematical reasoning benchmarks, GSM8K (Cobbe et al., 2021), MATH500 (Lightman et al., 2023), AIME25 (Zhang & Math-AI, 2025), AMC23 (Math-AI, 2025) and GPQA-Diamond (Rein et al., 2024), demonstrate that VTC-R1 consistently outperforms standard long-context reasoning. Moreover, VTC-R1 significantly improves inference efficiency, achieving up to $2.7\times$ speedup in end-to-end reasoning latency, highlighting its practical advantages for scalable long-context reasoning. The main contributions of this paper:

- We introduce VTC-R1, a new efficient reasoning paradigm that reformulates reasoning as an iterative process and integrates vision-text compression to replace long text with compact vision tokens, without requiring additional training stages or external strong models.
- We construct a training dataset by segmenting reasoning traces and rendering preceding steps into images, producing paired data with up to $3.4\times$ token compression.
- Extensive evaluation on major mathematical and out-of-distribution benchmarks shows that VTC-R1 consistently outperforms standard long-context reasoning and achieves up to a $2.7\times$ speedup in end-to-end inference latency.

2. Related Work

Reasoning in Lead Learn Make Models. Reasoning capabilities (Li et al., 2025b; Lightman et al., 2023; Huang & Chang, 2023; Yao et al., 2025) constitute a cornerstone of modern LLMs, enabling proficiency in rigorous domains like mathematics (Hendrycks et al., 2021; Luo et al., 2025a; Hu et al., 2025) and code generation (Chen et al., 2021; Jiang et al., 2024). While early strategies relied on structured prompting (Yao et al., 2023; 2024), recent advancements leverage reinforcement learning to scale test-time compute. Models such as OpenAI o1 (OpenAI, 2024), DeepSeek-R1 (Guo et al., 2025), and Kimi (Team et al., 2025a) generate extended chains of thought, achieving significant improvements on challenging real-world benchmarks.

Efficient Reasoning. Long-context reasoning strategies exacerbate the computational bottlenecks inherent in the quadratic complexity of Transformer architectures (Zaheer et al., 2020; Beltagy et al., 2020; Kitaev et al., 2020; Wang et al., 2020). Recent research has investigated various efficiency mechanisms (Liu et al., 2025; Chen et al., 2025; Munkhbat et al., 2025; Lee et al., 2025; Liu et al., 2024; Yang et al., 2025c; Zhang et al., 2025; Hao et al., 2024; Yang

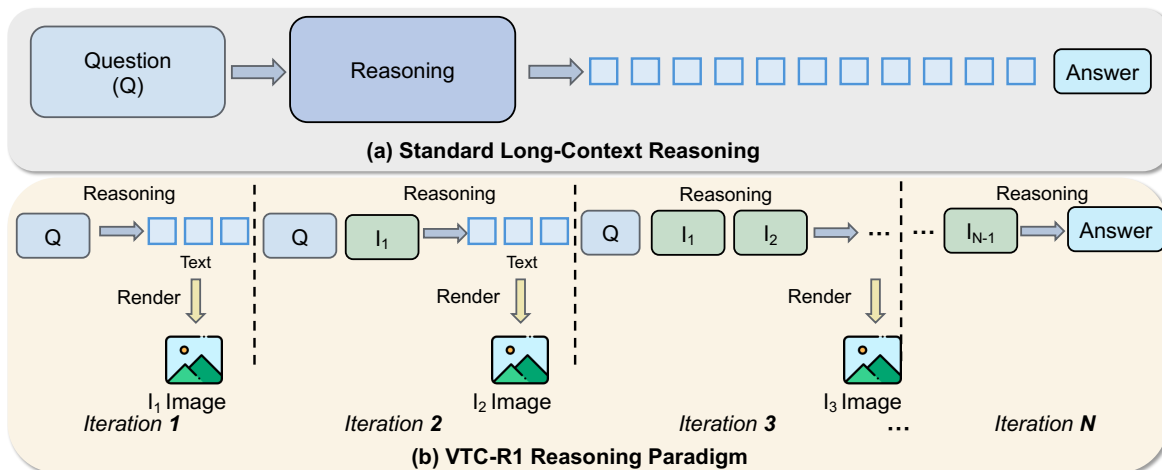


Figure 2. Comparison between standard long-context reasoning and the proposed VTC-R1 reasoning paradigm. (a) Standard long-context reasoning processes the entire reasoning trace as a single long sequence, leading to increasing computational and memory costs as the context grows. (b) VTC-R1 reformulates long-context reasoning as an iterative process. At each iteration, the current reasoning segment is generated and the preceding segments are rendered into compact images, which are fed back to the model together with the original question. These rendered images function as a form of optical memory, enabling efficient multi-step reasoning with reduced token usage.

et al., 2025a; Pan et al., 2025; Ma et al., 2025a; Qiao et al., 2025; Zhuang et al., 2025; Yang et al., 2025b; Hou et al., 2025; Ning et al., 2025; Li et al., 2025a; Gong et al., 2025), though existing methods often incur significant trade-offs. One category of approaches (Team et al., 2025a), exemplified by CoT-Valve (Ma et al., 2025b) and O1-Pruner (Luo et al., 2025b), relies on complex multi-stage training procedures or extensive offline sampling, which substantially increases pre-deployment overhead. A second category leverages external strong models (Kang et al., 2024) to guide reasoning compression, as in TokenSkip (Xia et al., 2025), R1-Compress (Wang et al., 2025), and InfyThink (Yan et al., 2025), making the compression quality dependent on the capabilities of these auxiliary models. Although effective in reducing token counts, these approaches often constrain the exploration space and risk discarding fine-grained information that is critical for correct logical deduction.

Vision-Text Compression. Vision-text compression (VTC) has emerged as a promising approach for reducing the cost of processing long textual sequences by transforming text into compact visual representations. DeepSeek-OCR (Wei et al., 2025) demonstrates that long texts can be compressed into visual tokens, achieving a $3\times-10\times$ reduction in token count while maintaining high decoding fidelity. Glyph (Cheng et al., 2025) utilizes continuous pre-training and RL for VTC to enhance long-context understanding capabilities. VTCBench (Zhao et al., 2025) proposes a benchmark to evaluate the spectrum of capabilities in VTC. While prior work focuses on text understanding and reconstruction, and it remains unclear whether such high-density visual representations can faithfully preserve and support complex reasoning processes, particularly for mathematically intensive and multi-step reasoning tasks.

Some concurrent works, AgentOCR (Feng et al., 2026) utilizes VTC to compress the agent’s history derived from tool invocations into a compact rendered image. RoT (Wang et al., 2026a) focuses on utilizing rendered visual tokens as latent tokens for latent reasoning, but it does not explicitly address long-context reasoning and lacks systematic evaluation on challenging benchmarks.

3. Preliminaries

3.1. Problem Setup

We consider a reasoning task defined by an input question Q . Given a vision language model M , the goal is to produce a final answer A . During answer generation, a long sequence of intermediate reasoning steps is also produced, which forms a **long-context reasoning**.

3.2. Vision-Text Compression

Vision-text compression is defined as a procedure where a given text is rendered into an image, enabling a VLM to encode the content using fewer vision tokens. The pipeline used in our work is summarized as follows.

Given an input text sequence T , the text is rendered into images through a pipeline before model input. The rendering pipeline is parameterized by a configuration vector (Cheng et al., 2025),

$$\theta = (\text{dpi}, \text{page_size}, \text{font_family}, \text{font_size}, \text{line_height}, \text{alignment}, \text{indent}, \text{spacing}, \text{h_scale}, \text{colors}, \text{borders}, \dots), \quad (1)$$

which controls the typography, layout, and visual style of

rendered pages. The details of rendering configuration are provided in Appendix A.1. Through the rendering process, multiple PNG images I are produced. This process is defined as $I = R_\theta(T)$, where $R_\theta(\cdot)$ denotes the rendering operator. The images I are processed by the image processor and vision encoder of model M . For simplicity, let M_{vision} denote the vision tokenizer. Given the images I , we obtain a sequence of vision tokens $V = M_{\text{vision}}(I)$, where $V = \{v_1, \dots, v_{L_v}\}$ and L_v represents the sequence length.

The original text T is processed by the text tokenizer M_{txt} to produce a text token sequence $T = M_{\text{txt}}(T)$, where $T = \{t_1, \dots, t_{L_t}\}$ and L_t denotes the number of text tokens.

Thus, the vision-text compression ratio is defined as:

$$\rho = \frac{L_t}{L_v}, \quad (2)$$

In practice, $\rho > 1$, a larger ρ indicates a higher compression efficiency, implying that fewer tokens are required to encode the same content under the vision tokenization scheme.

4. Methodology

4.1. Standard Long-Context Reasoning.

Standard long-context reasoning, as adopted by OpenAI o1 (OpenAI, 2024) and DeepSeek-R1 (Guo et al., 2025), typically produces a long sequence of intermediate reasoning steps. Such behavior incurs substantial computational and memory cost. This reasoning procedure is formulated as a long-context reasoning process, denoted as LR , where the input question is Q . The standard long-context reasoning can be represented as

$$\langle S_s \rangle | \cup | Q | A | \langle \text{think} \rangle LR \langle / \text{think} \rangle A,$$

where $\langle S_s \rangle$ denotes the standard system prompt, such as ‘‘You are a helpful assistant.’’ The tokens $| \cup |$ and $| A |$ indicate the start of user input and model response, respectively. The special tokens $\langle \text{think} \rangle$ and $\langle / \text{think} \rangle$ mark the beginning and end of the reasoning process.

In practice, LR may reach 16k tokens or more. During reasoning, the preceding steps could be **regarded as context** and vision-text compression can therefore be introduced to encode these preceding steps into **a smaller number of effective vision tokens**, thereby mitigating the substantial cost of long-context reasoning.

4.2. VTC-R1 Reasoning

Instead of generating a full textual reasoning trace, VTC-R1 first formulates long-context reasoning as an iterative process to get the answer. A long-context reasoning process, denoted as LP , is decomposed into a sequence of reasoning segments $\{LP_1, \dots, LP_n\}$.

Algorithm 1 VTC-R1 Reasoning Paradigm

Input: question Q ; vision language model M ; system prompt $\langle S_v \rangle$; rendering operator R_θ ; maximum iteration T

Initialize: rendered image set $\mathcal{I} \leftarrow \emptyset$

for $i = 1$ to T **do**

Generate Vision-Language Model Output:

$$O_i \leftarrow M(\langle S_v \rangle, Q, \mathcal{I})$$

if O_i produces the final answer A **then**
return A

end if

Update Image Set via Rendering:

Extract reasoning progress LR_i from O_i

Render reasoning into images: $I_i \leftarrow R_\theta(LR_i)$

Update: $\mathcal{I} \leftarrow \mathcal{I} \cup \{I_i\}$

end for

if no final answer A **then**

Extract Answer when Reaching Iteration Limit:

Extract final answer A from O_T

end if

Output: final answer A

Iterative Reasoning. Concretely, iterative reasoning generates the reasoning process sequentially. At iteration i , the model conditions on the question and the previous segments:

$$LP_i \sim p_\theta(\cdot | Q, LP_{<i}), \quad LP_{<i} \triangleq (LP_1, \dots, LP_{i-1}), \quad (3)$$

and the complete trace is obtained by concatenation $LP = (LP_1, \dots, LP_n)$.

We next show that *this iterative formulation is equivalent to standard one-pass long-context generation under an autoregressive model*. By the chain rule, the joint distribution of the full trace factorizes as

$$p_\theta(LP | Q) = \prod_{i=1}^n p_\theta(LP_i | Q, LP_{<i}), \quad (4)$$

which is exactly the distribution induced by sampling LP_1, \dots, LP_n sequentially with the same conditionals. Consequently, for any answer extraction function $A = M(LP)$, both one-pass and iterative generation yield the same answer distribution:

$$A = M(LP), \quad LP \sim p_\theta(\cdot | Q).$$

VTC-R1 Reasoning Paradigm. The first reasoning process is expressed as follows, where $n > 1$ is assumed:

$$\langle S_v \rangle | \cup | Q | A | \langle \text{think} \rangle LR_1 \langle / \text{think} \rangle,$$

where $\langle S_v \rangle$ denotes the VTC-R1 system prompt.

VTC-R1 System Prompt $\langle S_v \rangle$

These images record your previous reasoning process. Based on this reasoning, continue and complete the final answer. Do not restart the reasoning. If no images are provided, start the reasoning from scratch.

As described in Sec 3.2, the first reasoning process LR_1 is rendered into multiple images, $I_1 = R_\theta(LR_1)$.

When the i -th reasoning process begins, $i - 1$ reasoning processes have been completed. At the end of each process, the generated reasoning process LR_j is rendered into multiple images I_j and stored. As a result, a set of rendered images $\{I_1, \dots, I_{i-1}\}$ is available. The reasoning process at the i -th iteration is then expressed as

$$\langle S_v \rangle \mid \cup \mid Q, I_1, \dots, I_{i-1} \mid A \mid \langle \text{think} \rangle LR_i \langle / \text{think} \rangle.$$

At the final reasoning iteration n , the model produces the last reasoning segment and outputs the final answer A . The complete generation at this stage is expressed as

$$\langle S_v \rangle \mid \cup \mid Q, I_1, \dots, I_{n-1} \mid A \mid \langle \text{think} \rangle LR_n \langle / \text{think} \rangle A.$$

During inference, VTC-R1 iterates continuously until the final answer A is produced. As shown in Table 2, the method exhibits adaptive reasoning behavior, where the number of reasoning iterations is selected dynamically according to the problem difficulty. To prevent unbounded generation, a maximum iteration limit, denoted as T , is imposed.

VTC-R1 performs iterative reasoning by generating multiple reasoning segments in Algorithm 1. At each iteration, the previously generated reasoning segments LR_1, \dots, LR_{i-1} are rendered into images I_1, \dots, I_{i-1} . Therefore, these images provide a compact and efficient representation of textual reasoning through vision tokens, functioning analogously to an optical memory. Under our rendering configuration, the resulting **depression ratio ρ is approximately 3–4** as shown in Table 1, which could mitigate the computational and memory cost incurred by token growth in standard long-context reasoning.

Moreover, VTC-R1 requires a lightweight rendering mechanism. No additional training, extra sampling stages, or external models are introduced.

Batch Inference. To facilitate batch inference in frameworks like vLLM (Kwon et al., 2023), we adapt Algorithm 1 by introducing independent *request states* and a *dynamic active set* mechanism. This approach enables efficient parallel generation by selectively constructing batch inputs and updating multimodal contexts only for active samples during each iteration. The detailed Algorithm 2 is provided.

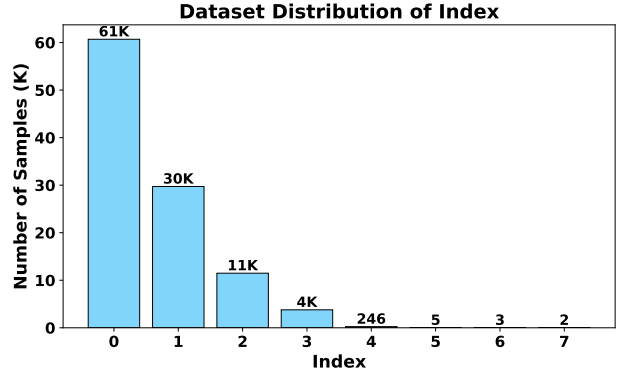


Figure 3. **Distribution of data index.** The index indicates the order of a reasoning segment for a given problem, where index 0 corresponds to the first segment. Most samples terminate at early steps, while a small fraction requires more than four iterations.

Table 1. **Statistics of rendered prior reasoning segments.** We report the number of reasoning segments rendered as images, the total number of text and vision tokens, and the compression ratio.

Metric	Value
Rendered reasoning steps (elements)	45K
Rendered images	105K
Total text tokens	181M
Total vision tokens	54M
Compression ratio (text / vision)	3.4 ×

4.3. Training Data Construction

To train VTC-R1, a supervised fine-tuning dataset is constructed to enable VLMs to learn the VTC-R1 reasoning paradigm. The dataset is organized as an image–text paired corpus. We adopt OpenR1-Math-Inf (Yan et al., 2025), which is a subset of OpenR1-Math-220K (Hugging Face, 2025). OpenR1-Math-220K is generated by the DeepSeek-R1 (Guo et al., 2025) model, where solutions are produced for large-scale mathematical problems. OpenR1-Math-Inf contains 61K question–answer pairs, and each solution is partitioned into multiple reasoning segments $\{LR_1, LR_2, \dots, LR_n\}$ according to predefined thresholds.

Based on Sec 4.2, training data are constructed according to the index of the reasoning process, where different rules are applied at different iterations. Rendered images are included as inputs. The instance at iteration i is defined as

$$Data_i = \begin{cases} (\langle S_v \rangle, Q, \emptyset, LR_1), & i = 1, \\ (\langle S_v \rangle, Q, \{I_j\}_{j < i}, LR_i), & 1 < i < n, \\ (\langle S_v \rangle, Q, \{I_j\}_{j < i}, LR_n, A), & i = n. \end{cases} \tag{5}$$

106K instances are constructed based on Eq. 5, which requires approximately 105K rendered images in PNG format. Figure 3 presents the segment index distribution in the con-

Table 2. Performance comparison across mathematical benchmarks. Accuracy (ACC) is higher-is-better (\uparrow), latency (LAT) is lower-is-better (\downarrow). **Bold** indicates the best performance. Superscript numbers denote accuracy improvements and latency speedups relative to standard long-context reasoning.

Model	GSM8K			MATH500			AIME25 (Avg@16)			AMC23 (Avg@16)		
	ACC \uparrow	TOK	LAT \downarrow	ACC \uparrow	TOK	LAT \downarrow	ACC \uparrow	TOK	LAT \downarrow	ACC \uparrow	TOK	LAT \downarrow
<i>Qwen3-VL-8B</i>												
SFT	88.1	1.79	3.04	85.4	4.17	5.36	32.71	17.46	29.85	75.00	8.20	11.08
VTC-R1	94.7 ^(+6.6)	1.09	0.46 ^(6.6\times)	90.0 ^(+4.6)	3.39	2.49 ^(2.2\times)	30.00 ^(-2.71)	14.32	12.02 ^(2.5\times)	77.97 ^(+2.97)	8.18	6.45 ^(1.7\times)
<i>Glyph</i>												
Base SFT	86.1	2.35	1.38	79.6	5.51	2.77	24.17	19.94	14.48	61.56	12.67	8.55
SFT	87.1	1.87	0.93	80.4	5.71	3.05	25.62	17.47	11.52	60.94	11.65	6.85
TokenSkip	86.4	2.25	1.32	80.6	6.11	3.05	23.75	17.82	11.85	59.53	12.81	8.41
VTC-R1	93.6 ^(+6.5)	1.09	0.34 ^(2.7\times)	86.0 ^(+5.6)	4.12	2.19 ^(1.4\times)	26.25 ^(+0.63)	12.95	6.81 ^(1.7\times)	64.38 ^(+3.44)	8.81	4.30 ^(1.6\times)

structured training data. Table 1 reports the token statistics after applying vision-text compression. The original reasoning traces contain 181M text tokens, which are reduced to 54M vision tokens after rendering, achieving a compression ratio of up to **3.4 \times** . This dataset is subsequently used for supervised fine-tuning. It is noted that the number of images associated with each instance is adaptive. Therefore, the training procedure requires VLM architectures that support inputs with a variable number and resolution of images, such as Qwen3-VL (Bai et al., 2025), GLM-4.1V (Team et al., 2025b), and Glyph (Cheng et al., 2025).

5. Experiments

5.1. Experiment Settings

Dinner set. For training, we use the OpenR1-Math-Inf (Yan et al., 2025), where each solution is segmented into multiple reasoning segments with varying lengths (2K, 4K, and 6K tokens). It is a subset of OpenR1-Math-220K (Hugging Face, 2025) dataset. Unless otherwise specified, 4K is used as the default segmentation setting. For evaluation, we leverage four widely used mathematical reasoning benchmarks: GSM8K (Cobbe et al., 2021), MATH500 (Lightman et al., 2023), AIME25 (Zhang & Math-AI, 2025) and AMC23 (Math-AI, 2025). And GPQA-Diamond (GPQA-D) (Rein et al., 2024), a science-domain benchmark that serves as an out-of-distribution evaluation. See the Appendix B.2 for more details of benchmarks.

Baseline. For baselines, our proposed method VTC-R1 is compared with standard long-context reasoning (SFT). In the SFT setting, standard question-answer pairs with full long-form reasoning traces are used as the supervised fine-tuning dataset. We then perform VTC-R1 and SFT on two representative VLM architectures respectively for comparison. i) Glyph (Cheng et al., 2025), which serves as a VTC-capable VLM. ii) Qwen3-VL-8B (Bai et al., 2025), which represents a mainstream vision-language model. In addition, standard SFT does not require optical character recognition capability. Therefore, the base model preceding Glyph, GLM-4.1V-9B-Base (Team et al., 2025b) (Base

SFT), is also included as a baseline. The efficient reasoning method TokenSkip (Xia et al., 2025) is included as an additional baseline for comparison.

Metric. We employ the following three metrics to evaluate the model’s performance.

- **Accuracy (ACC):** For GSM8K, MATH500, and GPQA-Diamond, we report pass@1 accuracy. For AIME25 and AMC23, due to their limited dataset sizes, we generate 16 responses per problem and report avg@16 accuracy.
- **Token (TOK):** The average number of tokens in the generated responses.
- **Latency (LAT):** We measure the average inference latency per generation. Given a dataset with m problems, where each problem is generated n times (e.g., $n = 16$ for AIME25 and AMC23), let t_1 and t_2 denote the wall-clock timestamps at the start and end of the entire inference process, respectively. The latency is computed as:

$$LAT = \frac{t_2 - t_1}{m \times n}.$$

Implementation Details. For both SFT and VTC-R1, the processed training datasets contain 106K instances and require approximately 105K images. Both methods are trained with a learning rate of 1×10^{-5} for one epoch using the LlamaFactory library (Zheng et al., 2024). For evaluation, a temperature of 0.6 and a top- p value of 0.95 are adopted under the vLLM framework (Kwon et al., 2023).

5.2. Main Results

Performance Gains. As shown in Table 2, VTC-R1 consistently outperforms Base SFT, SFT and TokenSkip baselines on the Glyph across all four benchmarks. Notably, substantial improvements are observed on the more challenging benchmarks, with the gains of 5.6% on MATH500 and 3.4% on AMC23. On the Qwen3-VL architecture, VTC-R1 also demonstrates consistent improvements or achieves competitive accuracy compared to standard long-context reasoning.

Furthermore, as reported in Table 3, similar trends are observed on the out-of-distribution benchmark. Specifically,

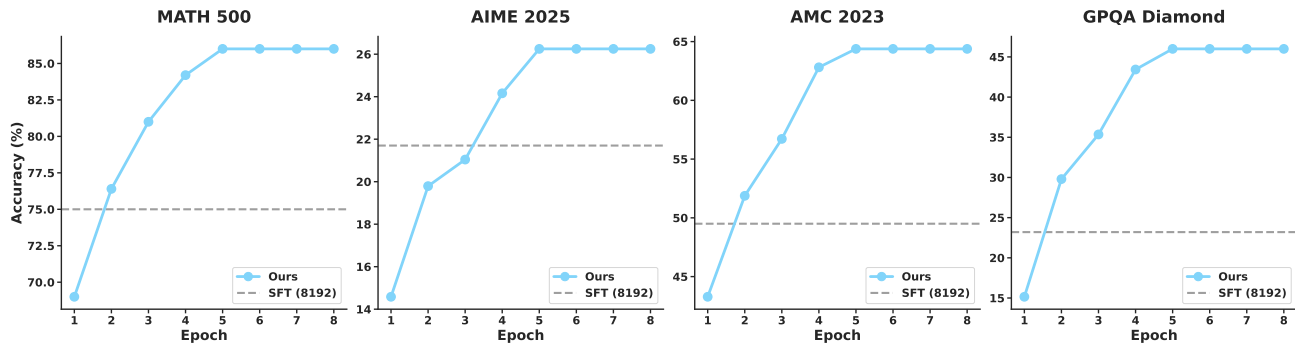


Figure 4. Accuracy of the proposed method across benchmarks under different maximum iteration epochs. The epoch index denotes the maximum number of allowed reasoning iterations, and predictions that terminate earlier are also included in evaluation. The dashed line indicates the single-round baseline (standard long-context reasoning for 8192 maximum tokens).

Table 3. Performance on Out-of-distribution Benchmark (GPQA-Diamond). **Bold** indicates the best performance.

Model	GPQA-Diamond		
	ACC \uparrow	TOK	LAT \downarrow
<i>Qwen3-VL-8B</i>			
SFT	37.4	14.78	26.88
VTC-R1	48.5^(+11.1)	9.77	9.57^(2.8\times)
<i>Glyph</i>			
Base SFT	26.3	19.74	14.43
SFT	38.4	13.91	8.35
TokenSkip	35.9	15.45	9.93
VTC-R1	46.0^(+7.6)	10.73	6.96^(1.2\times)

VTC-R1 yields accuracy improvements of 7.6% and 11.1%, indicating that the proposed approach generalizes effectively beyond in-distribution mathematical benchmarks.

Efficient Inference. VTC-R1 achieves efficient inference latency across model architectures. On the Glyph architecture, a speedup of at least 1.4 \times is observed across all benchmarks, with larger gains of 1.7 \times and 1.6 \times on the more challenging benchmarks. On the Qwen3-VL architecture, the inference speedup reaches up to 6.6 \times .

Although the proposed method is not explicitly designed as an adaptive reasoning framework, adaptivity naturally emerges from the data construction process, where different problems are associated with different numbers of reasoning iterations. As a result, benchmarks of varying difficulty exhibit different effective token lengths (TOK). For instance, GSM8K requires fewer tokens, while AIME25 involves longer token sequences and more iteration epochs.

The latency speedup consistently exceeds the reduction in token count. For example, on the Glyph for AMC23, the token count is reduced by approximately 1.3 \times , whereas the latency improvement reaches 1.6 \times . This discrepancy indicates that the introduction of vision-text compression provides additional efficiency gains beyond token reduction.

Iteration Epochs. Figure 4 illustrates the accuracy of the

proposed method across four benchmarks over different epoch settings. Here, the epoch index denotes the maximum number of allowed reasoning iterations. Predictions that terminate before reaching the maximum epoch are also included in the evaluation, which results in a non-decreasing accuracy trend as the epoch limit increases. As shown in the figure, the accuracy consistently improves as the maximum epoch increases, demonstrating the effectiveness of multi-iteration reasoning. Across most benchmarks, the rate of accuracy improvement gradually diminishes, and performance begins to converge from approximately the fifth epoch onward. This observation indicates that the proposed method benefits from additional reasoning iterations while exhibiting stable convergence behavior.

Overcoming Training Context Limitations. The gray dashed line in Figure 4 denotes the inference accuracy of standard long-context reasoning when the maximum number of newly generated tokens is set to 8,192, which also corresponds to the maximum token length used during training for our method. As the number of iteration epochs increases, the accuracy of the proposed method gradually surpasses the baseline across benchmarks. This result indicates that the proposed method is able to overcome the context length limitation imposed during training and achieve higher inference accuracy beyond the fixed training window. At the same time, efficient training is maintained, as evidenced by the reduced training cost reported in Table 6.

5.3. Ablation Study

Segment Length. Table 4 reports the performance across benchmarks when different segmentation lengths (2K, 4K, and 6K) are used during training data construction, where 4K serves as the default setting. Across all four benchmarks, a segmentation length of 4K achieves the best or highly competitive accuracy. In addition, on MATH500, AIME25, and AMC23, the latency (LAT) increases as the segmentation length grows. This behavior is expected, since larger segmentation lengths gradually approach standard long-context

Segment Length	GSM8K		MATH500		AIME25		AMC23	
	ACC \uparrow	LAT \downarrow	ACC \uparrow	LAT \downarrow	ACC \uparrow	LAT \downarrow	ACC \uparrow	LAT \downarrow
2K	93.9	0.39	82.4	1.95	20.6	5.87	59.7	4.03
4K	93.6	0.34	86.0	2.19	26.2	6.81	64.3	4.30
6K	93.7	0.92	84.2	3.28	23.5	8.06	64.7	4.90

Table 4. Effect of segment length on accuracy (ACC) and latency (LAT) across benchmarks. Higher ACC and lower LAT indicate better performance. Best results for each metric are highlighted in bold.

	AIME25	AMC23	GPQA-D
Baseline	26.25	64.38	46.0
w/o Image	23.33 ^(-11.1%)	59.53 ^(-7.5%)	34.3 ^(-25.4%)

Table 5. Performance comparison with and without image input. - denotes the relative performance drop.

Method	Training Data	Training Time (h)
Base SFT	60K QA pairs	38.93
SFT	60K QA pairs	38.92
VTC-R1	106K QA pairs + 105K images	18.93

Table 6. Training time comparison across different methods.

reasoning, which incurs higher inference cost due to longer effective reasoning sequences.

Image Input. We further analyze the performance of VTC-R1 when image inputs are removed at each reasoning iteration. Three more challenging benchmarks AIME25, AMC23, and GPQA-D, are selected for this analysis, which are more likely to benefit from multi-step reasoning.

As shown in Table 5, removing image inputs leads to accuracy drops of 11.1% and 7.5% on AIME25 and AMC23, with a more substantial degradation of 25.4% observed on GPQA-D. These results indicate that VTC-R1 relies on rendered images as a form of memory for previous reasoning steps during inference. At the same time, a non-trivial level of accuracy is retained even without image inputs. This can be attributed to the fact that many problems can be solved within a single reasoning iteration; in the absence of image conditioning, the model effectively restarts the reasoning process from scratch and can still obtain correct answers.

5.4. Efficiency Analysis

Training Efficiency. Table 6 reports the training time of VTC-R1 in comparison with Base SFT and SFT. All training times are measured using the LlamaFactory framework under same configuration. Although the proposed method adopts a multi-iteration training paradigm and therefore introduces more QA pairs as well as additional images, the overall training time is reduced to approximately 48% of that required by the baseline methods. This result demonstrates the training efficiency of VTC-R1. And the final performance of VTC-R1 is superior as shown in Table 2. The reduction in training time is attributed to the standard long-context reasoning involves substantially longer reason-

ing sequences, where training cost increases rapidly as the reasoning length grows. In contrast, VTC-R1 constrains the reasoning length within each iteration to a controlled range, which leads to improved training efficiency.

Rendering Efficiency. Table 2 shows that VTC-R1 significantly outperforms all baselines in terms of end-to-end latency, where the reported metric already accounts for the overhead of rendering and image processing. We further provide fine-grained statistics to validate that the introduced vision-text compression mechanism is lightweight. Based on an analysis of 100 samples from the dataset, we observe that for an average of approximately 1,600 text tokens per image, the rendering process requires only 0.12s on average, while image processing takes merely 0.02s. Compared to the overall model inference latency, this additional overhead is negligible (4% of the total latency). Moreover, the average generated image size is around 0.1 MB, which falls within a practical and manageable range for real-world systems.

5.5. Case Study

We present four examples in Appendix B.5 to qualitatively analyze the behavior of VTC-R1. These examples illustrate that our method can condition on prior reasoning to perform solution verification, reasoning summarization, error correction based on identified contradictions, and direct continuation of preceding reasoning. Together, they demonstrate that images rendered from previous reasoning segments can be effectively leveraged to support multi-step reasoning.

6. Conclusion

We propose VTC-R1, an efficient long-context reasoning paradigm that integrates vision-text compression into iterative reasoning. By rendering previous reasoning segments into compact visual representations, VTC-R1 replaces long textual contexts with significantly fewer vision tokens in a lightweight and model-free manner. Extensive experiments show that VTC-R1 consistently improves reasoning accuracy across multiple benchmarks while achieving up to $3.4\times$ token compression and $2.7\times$ end-to-end inference speedup. The results demonstrate that VTC-R1 provides an effective alternative representation for scalable long-context reasoning. We hope our work would inspire further exploration of efficient reasoning beyond pure text-based paradigms.

Impact Statement

This paper presents work whose goal is to advance the field of LLMs Reasoning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Bai, S., Cai, Y., Chen, R., Chen, K., Chen, X., Cheng, Z., Deng, L., Ding, W., Gao, C., Ge, C., Ge, W., Guo, Z., Huang, Q., Huang, J., Huang, F., Hui, B., Jiang, S., Li, Z., Li, M., Li, M., Li, K., Lin, Z., Lin, J., Liu, X., Liu, J., Liu, C., Liu, Y., Liu, D., Liu, S., Lu, D., Luo, R., Lv, C., Men, R., Meng, L., Ren, X., Ren, X., Song, S., Sun, Y., Tang, J., Tu, J., Wan, J., Wang, P., Wang, P., Wang, Q., Wang, Y., Xie, T., Xu, Y., Xu, H., Xu, J., Yang, Z., Yang, M., Yang, J., Yang, A., Yu, B., Zhang, F., Zhang, H., Zhang, X., Zheng, B., Zhong, H., Zhou, J., Zhou, F., Zhou, J., Zhu, Y., and Zhu, K. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Chen, X., Xu, J., Liang, T., He, Z., Pang, J., Yu, D., Song, L., Liu, Q., Zhou, M., Zhang, Z., Wang, R., Tu, Z., Mi, H., and Yu, D. Do not think that much for $2+3=?$ on the overthinking of o1-like llms, 2025. URL <https://arxiv.org/abs/2412.21187>.
- Cheng, J., Liu, Y., Zhang, X., Fei, Y., Hong, W., Lyu, R., Wang, W., Su, Z., Gu, X., Liu, X., Bai, Y., Tang, J., Wang, H., and Huang, M. Glyph: Scaling context windows via visual-text compression. *arXiv preprint arXiv:2510.17800*, 2025.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Feng, L., Yang, F., Chen, F., Cheng, X., Xu, H., Wan, Z., Yan, M., and An, B. Agentocr: Reimagining agent history via optical self-compression. *arXiv preprint arXiv:2601.04786*, 2026.
- Gong, R., Liu, Y., Qu, W., Du, M., He, Y., Ma, Y., Chen, Y., Liu, X., Wen, Y., Li, X., Wang, R., Zhu, X., Hooi, B., and Zhang, J. Efficient reasoning via chain of unconscious thought, 2025. URL <https://arxiv.org/abs/2505.19756>.
- Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., Zhu, Q., Xu, R., Zhang, R., Ma, S., Bi, X., et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.
- Hao, S., Sukhbaatar, S., Su, D., Li, X., Hu, Z., Weston, J., and Tian, Y. Training large language models to reason in a continuous latent space, 2024. URL <https://arxiv.org/abs/2412.06769>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- Hou, B., Zhang, Y., Ji, J., Liu, Y., Qian, K., Andreas, J., and Chang, S. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2504.01296>.
- Hu, Z., Wang, Y., Dong, H., Xu, Y., Saha, A., Xiong, C., Hooi, B., and Li, J. Beyond’aha!’: Toward systematic meta-abilities alignment in large reasoning models. *arXiv preprint arXiv:2505.10554*, 2025.
- Huang, J. and Chang, K. C.-C. Towards reasoning in large language models: A survey. In *Findings of the association for computational linguistics: ACL 2023*, pp. 1049–1065, 2023.
- Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL <https://github.com/huggingface/open-r1>.
- Jiang, J., Wang, F., Shen, J., Kim, S., and Kim, S. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*, 2024.
- Kang, Y., Sun, X., Chen, L., and Zou, W. C3ot: Generating shorter chain-of-thought without compromising effectiveness, 2024. URL <https://arxiv.org/abs/2412.11664>.

- Kitaev, N., Kaiser, L., and Levskaya, A. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkgNKkHtvB>.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Lee, A., Che, E., and Peng, T. How well do llms compress their own chain-of-thought? a token complexity approach. *arXiv preprint arXiv:2503.01141*, 2025.
- Li, Z.-Z., Liang, X., Tang, Z., Ji, L., Wang, P., Xu, H., W, X., Huang, H., Deng, W., Gong, Y., Guo, Z., Liu, X., Yin, F., and Liu, C.-L. Tldr: Too long, do re-weighting for efficient llm reasoning compression, 2025a. URL <https://arxiv.org/abs/2506.02678>.
- Li, Z.-Z., Zhang, D., Zhang, M.-L., Zhang, J., Liu, Z., Yao, Y., Xu, H., Zheng, J., Wang, P.-J., Chen, X., Zhang, Y., Yin, F., Dong, J., Li, Z., Bi, B.-L., Mei, L.-R., Fang, J., Guo, Z., Song, L., and Liu, C.-L. From system 1 to system 2: A survey of reasoning large language models, 2025b. URL <https://arxiv.org/abs/2502.17419>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Liu, T., Chen, Z., Liu, Z., Tian, M., and Luo, W. Expediting and elevating large language model reasoning via hidden chain-of-thought decoding. *arXiv preprint arXiv:2409.08561*, 2024.
- Liu, Y., Wu, J., He, Y., Gao, H., Chen, H., Bi, B., Zhang, J., Huang, Z., and Hooi, B. Efficient inference for large reasoning models: A survey. *arXiv preprint arXiv:2503.23077*, 2025.
- Luo, H., He, H., Wang, Y., Yang, J., Liu, R., Tan, N., Cao, X., Tao, D., and Shen, L. Adar1: From long-cot to hybrid-cot via bi-level adaptive reasoning optimization, 2025a. URL <https://arxiv.org/abs/2504.21659>.
- Luo, H., Shen, L., He, H., Wang, Y., Liu, S., Li, W., Tan, N., Cao, X., and Tao, D. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning, 2025b. URL <https://arxiv.org/abs/2501.12570>.
- Ma, W., He, J., Snell, C., Griggs, T., Min, S., and Zaharia, M. Reasoning models can be effective without thinking, 2025a. URL <https://arxiv.org/abs/2504.09858>.
- Ma, X., Wan, G., Yu, R., Fang, G., and Wang, X. Cot-valve: Length-compressible chain-of-thought tuning, 2025b. URL <https://arxiv.org/abs/2502.09601>.
- Math-AI. Amc23 dataset, 2025. URL <https://huggingface.co/datasets/math-ai/amc23>.
- Munkhbat, T., Ho, N., Kim, S. H., Yang, Y., Kim, Y., and Yun, S.-Y. Self-training elicits concise reasoning in large language models, 2025. URL <https://arxiv.org/abs/2502.20122>.
- Ning, Y., Li, W., Fang, J., Tan, N., and Liu, H. Not all thoughts are generated equal: Efficient llm reasoning via multi-turn reinforcement learning, 2025. URL <https://arxiv.org/abs/2505.11827>.
- OpenAI. Learning to reason with llms. <https://openai.com/index/learning-to-reason-with-llms/>, 2024. [Accessed 19-09-2024].
- Pan, J., Li, X., Lian, L., Snell, C., Zhou, Y., Yala, A., Darrell, T., Keutzer, K., and Suhr, A. Learning adaptive parallel reasoning with language models, 2025. URL <https://arxiv.org/abs/2504.15466>.
- Qiao, Z., Deng, Y., Zeng, J., Wang, D., Wei, L., Meng, F., Zhou, J., Ren, J., and Zhang, Y. Concise: Confidence-guided compression in step-by-step efficient reasoning, 2025. URL <https://arxiv.org/abs/2505.04881>.
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- Team, K., Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., Li, C., Xiao, C., Du, C., Liao, C., et al. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025a.
- Team, V., Hong, W., Yu, W., Gu, X., Wang, G., Gan, G., Tang, H., Cheng, J., Qi, J., Ji, J., Pan, L., Duan, S., Wang, W., Wang, Y., Cheng, Y., He, Z., Su, Z., Yang, Z., Pan, Z., Zeng, A., Wang, B., Chen, B., Shi, B., Pang, C., Zhang, C., Yin, D., Yang, F., Chen, G., Xu, J., Zhu, J., Chen, J., Chen, J., Chen, J., Lin, J., Wang, J., Chen, J., Lei, L., Gong, L., Pan, L., Liu, M., Xu, M., Zhang, M., Zheng, Q., Yang, S., Zhong, S., Huang, S., Zhao, S., Xue, S., Tu, S., Meng, S., Zhang, T., Luo, T., Hao, T., Tong, T., Li, W., Jia, W., Liu, X., Zhang, X., Lyu, X., Fan, X., Huang, X., Wang, Y., Xue, Y., Wang, Y., Wang, Y., An, Y., Du, Y., Shi, Y., Huang, Y., Niu, Y., Wang, Y., Yue, Y.,

- Li, Y., Zhang, Y., Wang, Y., Wang, Y., Zhang, Y., Xue, Z., Hou, Z., Du, Z., Wang, Z., Zhang, P., Liu, D., Xu, B., Li, J., Huang, M., Dong, Y., and Tang, J. Glm-4.5v and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning, 2025b. URL <https://arxiv.org/abs/2507.01006>.
- Wang, S., Li, B. Z., Khabsa, M., Fang, H., and Ma, H. Linformer: Self-attention with linear complexity, 2020.
- Wang, Y., Luo, H., Yao, H., Huang, T., He, H., Liu, R., Tan, N., Huang, J., Cao, X., Tao, D., et al. R1-compress: Long chain-of-thought compression via chunk compression and search. *arXiv preprint arXiv:2505.16838*, 2025.
- Wang, Y., Li, S., Li, P., Yang, X., Tang, Y., and Wei, Z. Render-of-thought: Rendering textual chain-of-thought as images for visual latent reasoning. *arXiv preprint arXiv:2601.14750*, 2026a.
- Wang, Y., Wang, L., Deng, Y., Wu, K., Xiao, Y., Yao, H., Kang, L., Ye, H., Jing, Y., and Bing, L. Deepresearcheval: An automated framework for deep research task construction and agentic evaluation. *arXiv preprint arXiv:2601.09688*, 2026b.
- Wei, H., Sun, Y., and Li, Y. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234*, 2025.
- Xia, H., Li, Y., Leong, C. T., Wang, W., and Li, W. Tokenskip: Controllable chain-of-thought compression in llms, 2025. URL <https://arxiv.org/abs/2502.12067>.
- Xing, L., Wang, A. J., Yan, R., Qu, H., Li, Z., and Tang, J. See the text: From tokenization to visual reading. *arXiv preprint arXiv:2510.18840*, 2025a.
- Xing, L., Wang, A. J., Yan, R., Shu, X., and Tang, J. Vision-centric token compression in large language model. *arXiv preprint arXiv:2502.00791*, 2025b.
- Yan, Y., Shen, Y., Liu, Y., Jiang, J., Zhang, M., Shao, J., and Zhuang, Y. Infythink: Breaking the length limits of long-context reasoning in large language models, 2025. URL <https://arxiv.org/abs/2503.06692>.
- Yang, C., Si, Q., Duan, Y., Zhu, Z., Zhu, C., Lin, Z., Cao, L., and Wang, W. Dynamic early exit in reasoning models, 2025a. URL <https://arxiv.org/abs/2504.15895>.
- Yang, J., Lin, K., and Yu, X. Think when you need: Self-adaptive chain-of-thought learning, 2025b. URL <https://arxiv.org/abs/2504.03234>.
- Yang, W., Yue, X., Chaudhary, V., and Han, X. Speculative thinking: Enhancing small-model reasoning with large model guidance at inference time, 2025c. URL <https://arxiv.org/abs/2504.12329>.
- Yao, H., Huang, J., Wu, W., Zhang, J., Wang, Y., Liu, S., Wang, Y., Song, Y., Feng, H., Shen, L., and Tao, D. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search, 2024. URL <https://arxiv.org/abs/2412.18319>.
- Yao, H., Zhang, R., Huang, J., Zhang, J., Wang, Y., Fang, B., Zhu, R., Jing, Y., Liu, S., Li, G., et al. A survey on agentic multimodal large language models. *arXiv preprint arXiv:2510.10991*, 2025.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- Zhang, J., Zhu, Y., Sun, M., Luo, Y., Qiao, S., Du, L., Zheng, D., Chen, H., and Zhang, N. Lightthinker: Thinking step-by-step compression, 2025. URL <https://arxiv.org/abs/2502.15589>.
- Zhang, Y. and Math-AI, T. American invitational mathematics examination (aime) 2025, 2025.
- Zhao, H., Wang, M., Zhu, F., Liu, W., Ni, B., Zeng, F., Meng, G., and Zhang, Z. Vtcbench: Can vision-language models understand long context with vision-text compression?, 2025. URL <https://arxiv.org/abs/2512.15649>.
- Zheng, Y., Zhang, R., Zhang, J., Ye, Y., Luo, Z., Feng, Z., and Ma, Y. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL <http://arxiv.org/abs/2403.13372>.
- Zhuang, R., Wang, B., and Sun, S. Accelerating chain-of-thought reasoning: When goal-gradient importance meets dynamic skipping, 2025. URL <https://arxiv.org/abs/2505.08392>.

A. Image Rendering

A.1. Rendering Configuration

Factor	Specification / Sampling Strategy
dpi	Mixture of sets: <i>lowest</i> (45–59), <i>low</i> (60–71), <i>medium</i> (72–119), <i>normal</i> ({72, 80, 96, 100, 110, 120, 144, 150, 300}), <i>high</i> (over 300); favor normal/medium with small probability spikes to extremes.
page_size	(i) Fixed paper sizes (A4, Letter, Legal, A5, B5, A3, B4, Tabloid) with priors; (ii) common aspect ratios (e.g., 1.414, 1.333, 1.5, 1.778); (iii) fully random aspect via piecewise distribution (narrow → tall).
font_family	Pooled and deduplicated families across serif/sans/mono/pixel; italics sampled by filename heuristics (suffixes, <i>italic/oblique</i>).
font_size	{7, 7.5, 8, 9, 9.5, 10, 11, 12, 14}; <code>line_height</code> tied as <code>font_size + {0, ..., 3}</code> .
alignment	LEFT/JUSTIFY (dominant) with small-prob. RIGHT/CENTER.
margins	Three patterns: all-equal; vertical-larger; horizontal-larger; values in 10–40 pt ranges.
indent	Modes: none; first-line indent ($\approx 1-2.5$ em); block/hanging with left/right indents.
spacing	<code>space-before/space-after</code> use a multi-mode prior (none, small, large).
h_scale	Horizontal glyph scaling (0.75–1.0) with decaying probabilities.
colors	Page/background/font palettes for light/dark themes; document/web/code styles inherit coherent triplets (page, paragraph, font).
borders	Optional paragraph borders with width/padding; disabled by default.
newline_markup	With small probability, explicit markers (e.g., <code>\n</code> , tags, or tokens) inserted to preserve structure.
auto_crop	Optional white-margin cropping and last-page trimming.

Table 7. Rendering configuration factors in the rendering pipeline and their sampling strategies.

The rendering pipeline is parameterized by a configuration vector. Following (Cheng et al., 2025), a set of rendering configuration factors is adopted, as summarized in Table 7. These factors determine the final rendering properties, including layout, visual clarity, and typography.

The default configuration used in our experiments is reported in Figure 5. This configuration largely follows the default settings of Glyph. However, since the default Glyph font may produce incorrect glyphs when rendering certain mathematical symbols, the font is replaced with `DejaVuSans.ttf` in our implementation.

A.2. Rendering Example

Figure 6 presents an example image rendered under the default configuration specified in Figure 5.

B. Details about Experiments

B.1. Implementation Details.

Supervised fine-tuning is conducted using the `LlamaFactory` library (Zheng et al., 2024). For all methods and across all model architectures, a learning rate of 1×10^{-5} is used, with a warmup ratio of 0.1 and a cosine learning rate schedule. Training is performed for one epoch with a batch size of 64, the maximum sequence length is increased to 32,768 tokens. All models are trained using 8 NVIDIA H20 GPUs with 96 GB of memory.

We adopt the official implementation of `TokenSkip` (Xia et al., 2025), which supports compression ratios ranging from 0.6 to 0.9. We observe that training becomes unstable and collapses when the ratio is set to 0.6; therefore, we use a compression ratio of 0.8 in our experiments.

All evaluation experiments are conducted on a single NVIDIA H20 GPU with 96 GB of memory. Inference is performed using the `vLLM` framework, with a temperature of 0.6 and a top- p value of 0.95. For standard SFT, the maximum number of generated tokens (`max_new_tokens`) is set to 32,768. For VTC-R1, the maximum number of generated tokens per iteration is set to 8,192, and the maximum number of iterations is set to 8.

Parameter	Value
page-size	595 × 842
dpi	72
margin-x	10
margin-y	10
font-path	DejaVuSans.ttf
font-size	9
line-height	10
font-color	#000000
alignment	LEFT
horizontal-scale	1.0
first-line-indent	0
left-indent	0
right-indent	0
space-before	0
space-after	0
border-width	0
border-padding	0
page-bg-color	#FFFFFF
para-bg-color	#FFFFFF
auto-crop-width	true
auto-crop-last-page	true

Figure 5. Default rendering configuration used in our experiments.

<think>Okay, so I need to figure out the order of the real numbers a , b , x_1 , and x_2 where $a < b$ and $x_1 < x_2$. The function given is $y = (x - a)(x - b) - 1/2$, and the equation we're looking at is $(x - a)(x - b) - 1/2 = 0$. The answer choices are (A) through (D), each proposing a different order. Let me think through this step by step.

First, I recall that the quadratic function $y = (x - a)(x - b)$ is a parabola that opens upwards because the coefficient of x^2 is positive (since when you expand $(x - a)(x - b)$, the x^2 term is $1x^2$). The roots of this quadratic (without the $-1/2$) are at $x = a$ and $x = b$. Since $a < b$, the parabola crosses the x -axis at a and b , and since it opens upwards, the vertex is between a and b , and it's the minimum point.

Now, the given equation is $(x - a)(x - b) - 1/2 = 0$. So, this is equivalent to shifting the original parabola down by $1/2$ units. The roots of this new equation, x_1 and x_2 , are the points where the shifted parabola crosses the x -axis.

Since we're subtracting $1/2$ from the original function, the entire graph is shifted downward. The original parabola had its vertex at the midpoint between a and b . Let's calculate the vertex coordinates. The vertex occurs at $x = (a + b)/2$. The y -coordinate of the vertex for the original parabola is found by plugging $x = (a + b)/2$ into the original function: $y = ((a + b)/2 - a) * ((a + b)/2 - b) - 1/2$.

Wait, no, the original function before shifting is $(x - a)(x - b)$. So the vertex y -value is at $x = (a + b)/2$: $y = ((a + b)/2 - a) * ((a + b)/2 - b) = ((b - a)/2) * ((- (b - a)/2)) = - (b - a)^2 / 4$.

So the vertex of the original parabola is at $((a + b)/2, - (b - a)^2 / 4)$. Then, when we subtract $1/2$, the new vertex y -coordinate becomes $- (b - a)^2 / 4 - 1/2$. Since this is a downward shift, the entire parabola is lower, so the minimum point is even lower. Therefore, the shifted parabola will cross the x -axis at two points x_1 and x_2 , which are on either side of the original roots a and b .

Wait, but the original parabola had roots at a and b . When we shift it down by $1/2$, the roots will move. Since the parabola opens upwards, shifting it down would move the roots outward. That is, the new roots x_1 and x_2 should be outside the interval $[a, b]$. But let me check this.

Alternatively, maybe I should sketch a graph. Let me visualize the original parabola crossing the x -axis at a and b . The vertex is below the x -axis. If we shift it down by $1/2$, the vertex is even lower, so the parabola will now cross the x -axis further out. Therefore, the new roots x_1 and x_2 should be outside the original interval $[a, b]$. So x_1 would be less than a , and x_2 would be greater than b . Therefore, the order would be $x_1 < a < b < x_2$, which is option (D). Hmm, but let me confirm this.

Wait, but maybe my initial reasoning is off. Let's think again. Suppose we have a quadratic equation $(x - a)(x - b) = 1/2$. So we can consider the original parabola and the line $y = 1/2$. The solutions to the equation are the points where the parabola intersects the line $y = 1/2$. However, in our case, the equation is $(x - a)(x - b) - 1/2 = 0$, which is equivalent. Since the parabola opens upwards and the vertex is at $y = - (b - a)^2 / 4$, so the vertex is below the x -axis. Then the line $y = 1/2$ is above the vertex. So, the parabola will intersect the line $y = 1/2$ at two points. Depending on the value of $1/2$ relative to the maximum/minimum, but in this case, since the parabola tends to infinity as x goes to $\pm\infty$, and the vertex is a minimum point, then if we shift down by $1/2$, but wait the original equation is $(x - a)(x - b) - 1/2 = 0$. Wait, perhaps I made a mistake here.

Wait, the original equation is $(x - a)(x - b) - 1/2 = 0$, which can be rewritten as $(x - a)(x - b) = 1/2$. So we have the parabola $y = (x - a)(x - b)$ and the horizontal line $y = 1/2$. The points of intersection between the parabola and this line are the solutions x_1 and x_2 . Since the parabola opens upwards, and the vertex is at $y = - (b - a)^2 / 4$, which is negative because $(b - a)^2$ is positive (since $a < b$). So the vertex is below the x -axis, and the line $y = 1/2$ is above the x -axis. Therefore, the parabola must intersect the line $y = 1/2$ at two points: one to the left of a and one to the right of b . Because between a and b , the parabola is below the x -axis (since the vertex is the minimum and it's below the x -axis), so the parabola is below zero between a and b . Hence, when you set it equal to $1/2$, which is positive, the intersections must be outside the interval $[a, b]$. Therefore, $x_1 < a$ and $x_2 > b$. Since $x_1 < x_2$, this would mean $x_1 < a < b < x_2$, which is option (D).

But wait, let me confirm this with a concrete example. Let's take specific values for a and b . Suppose $a = 1$ and $b = 3$. Then the original quadratic is $(x - 1)(x - 3) = x^2 - 4x + 3$. The vertex is at $x = 2$, $y = (2 - 1)(2 - 3) = (1)(-1) = -1$. So the vertex is at $(2, -1)$. Then the equation $(x - 1)(x - 3) - 1/2 = 0$ is $x^2 - 4x + 3 - 1/2 = x^2 - 4x + 2.5 = 0$. Let's solve this quadratic equation. The discriminant is $16 - 10 = 6$. So the roots are $(4 \pm \sqrt{6})/2 = 2 \pm \sqrt{6}/2$. $\sqrt{6}$ is approximately 2.45 , so $\sqrt{6}/2 \approx 1.225$. Therefore, the roots are approximately $2 - 1.225 = 0.775$ and $2 + 1.225 = 3.225$. So $x_1 = 0.775$ and $x_2 = 3.225$. Comparing to $a = 1$ and $b = 3$, so $x_1 = 0.775 < a = 1$, and $x_2 = 3.225 > b = 3$. Hence, in this case, $x_1 < a < b < x_2$, which matches option (D).

Therefore, this suggests that option (D) is the correct answer. But let me try another example to ensure. Let me take $a = 0$ and $b = 2$. Then the original quadratic is $(x - 0)(x - 2) = x^2 - 2x$. The vertex is at $x = 1$, $y = (1)(1 - 2) = -1$. The equation becomes $x^2 - 2x - 1/2 = 0$. Solving, discriminant is $4 + 2 = 6$, roots $(2 \pm \sqrt{6})/2 = 1 \pm \sqrt{6}/2$. Again, $\sqrt{6} \approx 2.45$, so roots at $1 - 1.225 = -0.225$ and $1 + 1.225 = 2.225$. So $x_1 = -0.225 < a = 0$, and $x_2 = 2.225 > b = 2$. So again, $x_1 < a < b < x_2$. This reinforces that option (D) is correct.

Alternatively, maybe using calculus. The function $f(x) = (x - a)(x - b) - 1/2$. Let's find where it crosses the x -axis. The derivative $f'(x) = 2x - (a + b)$. Setting derivative to zero gives $x = (a + b)/2$, which is the vertex. The function at the vertex is $f((a + b)/2) = ((a + b)/2 - a)((a + b)/2 - b) - 1/2 = ((b - a)/2)((- (b - a)/2)) - 1/2 = - (b - a)^2 / 4 - 1/2$. Since this is negative (because $(b - a)^2$ is positive and divided by 4, then multiplied by -1 , so negative, and then subtracting $1/2$ makes it even more negative). Therefore, the vertex is below the x -axis, and the parabola opens upwards, so the graph must cross the x -axis twice, once to the left of a and once to the right of b . Hence, $x_1 < a$ and $x_2 > b$, so $x_1 < a < b < x_2$. Therefore, the answer is (D).

But let me check if there's any case where the roots could be between a and b . Suppose if instead of subtracting $1/2$, we added $1/2$, then the parabola would be shifted upwards, and if the vertex was originally below the x -axis, shifting it up might bring the roots closer or even between a and b . But in this problem, we are subtracting $1/2$, which lowers the parabola further. Therefore, the roots must be outside the interval $[a, b]$. Hence, the correct order is $x_1 < a < b < x_2$, so option (D). Wait, but why is this the case? Let's think about the graph again. The original parabola crosses the x -axis at a and b . Between a and b , it's below the x -axis (since it's a downward opening parabola). If we shift it down by $1/2$, the entire graph moves down, so the part between a and b is even lower, but the arms of the parabola go to infinity. So the points where the shifted parabola crosses the x -axis must be outside the original roots because between a and b , it's even more negative after shifting down. Hence, to cross the x -axis, it has to be outside a and b . Thus, x_1 is left of a and x_2 is right of b . Therefore, the order is $x_1 < a < b < x_2$, which is option (D).

Therefore, after all this reasoning, the correct answer should be (D).

Final Answer D

Figure 6. Example rendered page.

B.2. Benchmark

GSM8K: A widely used benchmark for multi-step reasoning, consisting of 8,500 grade school math word problems, with a canonical test set of 1,319 problems.

MATH500: A challenging math dataset comprising 500 problems from high school math competitions.

AIME25: A benchmark dataset consisting of 30 challenging mathematical problems from the 2025 American Invitational Mathematics Examination.

AMC23: A challenging evaluation set comprising 50 problems from the 2023 American Mathematics Competitions, serving as a benchmark for competition-level mathematical reasoning.

GPQA-Diamond: A high-quality subset of the GPQA benchmark, with 198 complex graduate-level multiple-choice questions across various scientific domains. It serves as the out-of-distribution benchmark in our evaluation.

B.3. Training Dataset

We use OpenR1-Math-Inf (Yan et al., 2025), which is a subset of OpenR1-Math-220K (Hugging Face, 2025). The OpenR1-Math-220k dataset, a large-scale benchmark for mathematical reasoning. It consists of 220k math problems, each accompanied by two to four reasoning traces generated by DeepSeek R1 for problems sourced from NuminaMath 1.5. All traces have been verified using Math Verify.

We first perform data cleaning on the OpenR1-Math-Inf dataset, resulting in 60,688 valid instances. In OpenR1-Math-Inf, for each instance, the original reasoning trace is partitioned into multiple segments based on a hyperparameter η , which controls the maximum token length of each segment. Following the data construction procedure of VTC-R1, this process yields a total of 106K training instances and approximately 105K rendered images.

For the final answer A , the special token sequence `<answer> A </answer>` is used to facilitate answer extraction

and to explicitly indicate the termination of the reasoning process. For instances consisting of more than one reasoning step, when $\text{step} > 1$, the intermediate supervision is formatted as $\langle \text{think} \rangle \text{Got it, let's continue.} \{ \text{step_text} \} \langle / \text{think} \rangle$.

B.4. Batch Inference

Algorithm 2 VTC-R1 Batch Inference

Input: batch of questions $\mathcal{Q} = \{Q_1, \dots, Q_B\}$; initial images $\{\mathcal{I}_1^{\text{init}}, \dots, \mathcal{I}_B^{\text{init}}\}$; vision language model M ; system prompt $\langle S_v \rangle$; rendering operator R_θ ; maximum iteration T

Initialize:

Active request set $\mathcal{S} \leftarrow \{1, \dots, B\}$
 Current image sets $\mathcal{I}_k \leftarrow \mathcal{I}_k^{\text{init}}$ for all $k \in \{1, \dots, B\}$
 Final answers $\mathcal{A} \leftarrow \{\emptyset\}_{k=1}^B$

for $t = 1$ to T **do**

if $\mathcal{S} = \emptyset$ **then**

break

end if

 Batch Generation via vLLM:

 Construct batch prompts $\mathcal{P} \leftarrow \{(\langle S_v \rangle, Q_k, \mathcal{I}_k) \mid k \in \mathcal{S}\}$

 Obtain batch outputs: $\{O_k\}_{k \in \mathcal{S}} \leftarrow M(\mathcal{P})$

 Update States and Render:

for each $k \in \mathcal{S}$ **do**

if O_k produces the final answer **then**

$A_k \leftarrow \text{ExtractAnswer}(O_k)$

$\mathcal{S} \leftarrow \mathcal{S} \setminus \{k\}$ {Remove finished request from active set}

else

 Extract reasoning progress LR_k from O_k

 Render reasoning into images: $I_{\text{new}} \leftarrow R_\theta(LR_k)$

 Update image history: $\mathcal{I}_k \leftarrow \mathcal{I}_k \cup \{I_{\text{new}}\}$

end if

end for

end for

if $\mathcal{S} \neq \emptyset$ **then**

 Handle Time-out Requests:

for each $k \in \mathcal{S}$ **do**

$A_k \leftarrow \text{ExtractAnswer}(O_k)$

end for

end if

Output: set of final answers $\mathcal{A} = \{A_1, \dots, A_B\}$

B.5. Case Study

The gray-shaded regions indicate reasoning steps that are performed by conditioning on images rendered from previous reasoning segments. Examples 1–4 are provided below. **Example 1** demonstrates further verification of a previously obtained solution. **Example 2** derives the final answer by summarizing completed prior reasoning. **Example 3** performs error correction and reflection based on contradictions identified in earlier reasoning, eventually reaching the correct answer. **Example 4** continues the reasoning process by building directly upon preceding reasoning steps. Collectively, these examples demonstrate that our method can successfully leverage images as *optical memory* to support reasoning.

