# Glyph: Scaling Context Windows via Visual-Text Compression

**Jiale Cheng**[1,2*] , **Yusen Liu**[2*] , **Xinyu Zhang**[2*] , **Yulin Fei**[2*] , **Wenyi Hong**[2,3]
**Ruiliang Lyu**[2] , **Weihan Wang**[2] , **Zhe Su**[2] , **Xiaotao Gu**[2] , **Xiao Liu**[2,3] , **Yushi Bai**[2,3]
**Jie Tang**[3], **Hongning Wang**[1] , **Minlie Huang**[1†]

[1]The Conversational Artificial Intelligence (CoAI) Group, Tsinghua University

[2]Zhipu AI

[3]The Knowledge Engineering Group (KEG), Tsinghua University

chengjl23@mails.tsinghua.edu.cn, aihuang@tsinghua.edu.cn

## Abstract

Large language models (LLMs) increasingly rely on long-context modeling for tasks such as document understanding, code analysis, and multi-step reasoning. However, scaling context windows to the million-token level brings prohibitive computational and memory costs, limiting the practicality of long-context LLMs. In this work, we take a different perspective—visual context scaling—to tackle this challenge. Instead of extending token-based sequences, we propose Glyph, a framework that renders long texts into images and processes them with vision–language models (VLMs). This approach substantially compresses textual input while preserving semantic information, and we further design an LLM-driven genetic search to identify optimal visual rendering configurations for balancing accuracy and compression. Through extensive experiments, we demonstrate that our method achieves 3–4× token compression while maintaining accuracy comparable to leading LLMs such as Qwen3-8B on various long-context benchmarks. This compression also leads to around 4× faster prefilling and decoding, and approximately 2× faster SFT training. Furthermore, under extreme compression, a 128K-context VLM could scale to handle 1M-token-level text tasks. In addition, the rendered text data benefits real-world multimodal tasks, such as document understanding. Our code and model are released at https://github.com/thu-coai/Glyph.

## 1 Introduction

Recent advances in large language models (LLMs) have enabled remarkable progress across a wide spectrum of real-world tasks (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023; GLM et al., 2024; Yang et al., 2025). As LLMs become increasingly capable, the demand for long-context
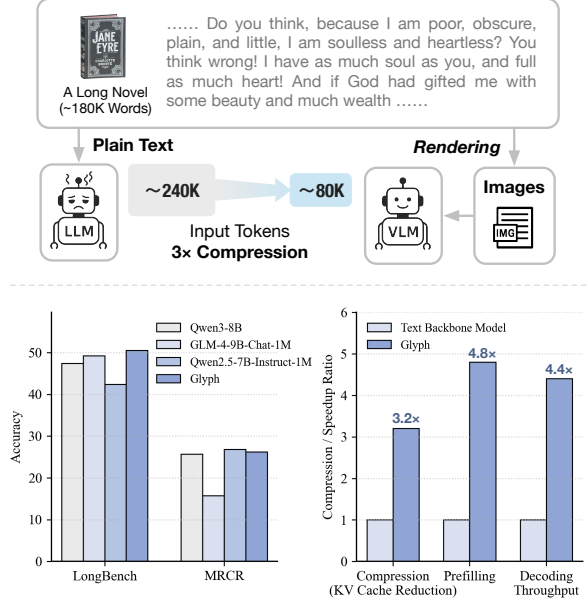


Figure 1: (Upper) Comparison of two paradigms for long-context tasks: conventional approaches directly feeding plain text into LLMs, and the proposed VLM-based paradigm, Glyph, which renders text as compact images to achieve substantial input-token compression. (Lower) Glyph attains competitive performance on LongBench and MRCR, while offering significant compression and inference speedup over its text backbone model on 128K-token inputs.

modeling has grown critical, especially for applications such as document understanding, code analysis, and multi-hop reasoning (Bai et al., 2024; Comanici et al., 2025). However, scaling context windows to hundreds of thousands or even millions of tokens poses prohibitive training and inference costs in both computation and memory, severely limiting the practicality of such models in real-world applications.

Recent work has explored two main directions to alleviate these costs. One line of work extends positional encodings, such as YaRN (Peng et al., 2023), allowing well-trained models to accept longer inputs without additional training. However, such methods neither accelerate inference nor maintain

---

[*] Core contributors.

[†] Corresponding author.

accuracy when extrapolated to much longer sequences (Wu et al., 2024). Another line focuses on modifying the attention mechanism, e.g., sparse or linear attention (Huang et al., 2023; Yang et al., 2024; Peng et al., 2025; Chen et al., 2025a), which reduces the quadratic complexity of self-attention and improves per-token efficiency. Yet, as context length grows to hundreds of thousands of tokens, the overall overhead remains substantial, since the number of tokens is unchanged. Retrieval-augmented approaches (Laban et al., 2024; Yu et al., 2025a) instead shorten the input length through external retrieval, but they risk missing important information and could introduce additional latency.

Distinct from the aforementioned approaches, we propose Glyph, a new paradigm that scales context length by rendering plain text into compact images and leveraging vision-language models (VLMs) to process the rendered inputs. In this way, the VLM operates directly on the glyphs of the text—treating each visual token as a compact carrier of multiple textual tokens—thereby increasing the information density without sacrificing semantic fidelity. This glyph-based visual representation allows a fixed-context VLM to process substantially longer texts than a text-only LLM with the same context length, thereby enabling long-context understanding without expanding the context window or relying on external retrieval mechanisms. For example, consider the novel "Jane Eyre" ($\approx$240K text tokens). A conventional 128K–context LLM cannot accommodate the entire book, and truncation easily leads to wrong answers for questions requiring global coverage, such as "Who supports Jane when she is in distress after leaving Thornfield?" In contrast, Glyph renders the book into compact images (e.g. $\approx$80K visual tokens), enabling a 128K–context VLM to process the full novel and answer such questions reliably.

Specifically, Glyph consists of three main stages, namely, continual pre-training, LLM-driven rendering search, and post-training. In the continual pre-training stage, we render large-scale long-context text into diverse visual forms, enabling the VLM to transfer its long-context capability from text tokens to visual tokens. Since the text-to-image conversion directly determines the trade-off between context compression and model performance, devising an optimal configuration of the conversion is crucial for downstream performance. To this end, we design an LLM-driven genetic search to auto-matically explore rendering parameters (e.g., font size, layout, resolution) to maximize compression while preserving long-context ability. The resulting configuration is then applied in the post-training stage, where we perform supervised fine-tuning and reinforcement learning to further improve the model's performance on visualized input. An auxiliary OCR task is applied to enhance the model's ability to recognize textual content within images, thereby better aligning its visual and textual representations, yielding the final Glyph model.

We conduct extensive experiments to evaluate the performance of Glyph. Results demonstrate that Glyph achieves 3–4× token compression of long sequences while preserving accuracy comparable to state-of-the-art LLMs such as Qwen3-8B. This compression not only extends the effective context length but also improves both training and inference efficiency, yielding up to 4.8× faster pre-filling, 4.4× faster decoding, as well as about 2× faster SFT training. Moreover, we find that incorporating rendered text data effectively enhances performance on real-world multimodal long-context tasks, such as document understanding.

Our contributions can be summarized as follows:

- We introduce a novel framework, Glyph, which enables long-context modeling through visual-text compression using VLMs, providing an alternative route to scaling context windows without incurring prohibitive computational and memory costs.

- We propose an LLM-driven genetic search that automatically identifies the optimal configurations of text-to-image rendering, ensuring both task performance and effective compression.

- We demonstrate that Glyph can achieve 3-4× token compression for long text sequences while preserving performance, enabling substantial improvements in memory efficiency, training, and inference speed.

## 2 Related Work

### 2.1 Long-Context Modeling

Research on extending LLMs to long contexts mainly focuses on architectural and training methods. Architecturally, studies have proposed sparse and hierarchical attention (Yang et al., 2016; Beltagy et al., 2020; Huang et al., 2023; Yang et al.,
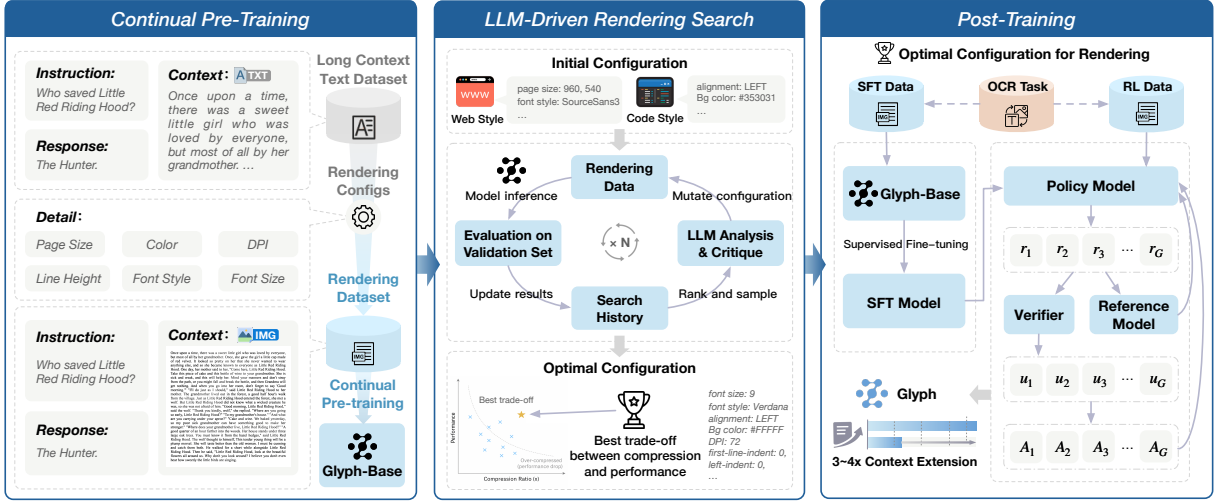
Figure 2: Glyph consists of three main stages: continual pre-training on rendered long-text data, LLM-driven genetic search for optimal rendering configurations, and post-training with SFT, RL. Together, these stages enable efficient long-context modeling with visual-text compression.

2024; Peng et al., 2025; Chen et al., 2025a), positional interpolation and extrapolation (Su et al., 2021; Press et al., 2021; Sun et al., 2022; Peng et al., 2023), and content-aware encodings (Chen et al., 2025b; Zhu et al., 2024). On the training side, LongAlign (Zhang et al., 2024) had built instruction datasets and loss-weighting strategies for sequences up to 100k tokens, while LongLoRA (Chen et al., 2024) had combined shifted sparse attention with parameter-efficient fine-tuning. LongRecipe (Wang et al., 2024b) had improved efficiency by integrating token analysis, index transformation, and optimization, scaling open-source models from 8k to 128k. ProLong (Liu et al., 2024b) had taken a data-centric view, selecting samples with long-range dependencies. In contrast, our method compresses text into visual tokens, which can be combined with existing techniques to reduce cost and extend context length.

## 2.2 Multimodal Large Language Model

Multimodal large language models (MLLMs) extend traditional LLMs to process and reason over text and visual inputs jointly. Early studies primarily focus on architectural design and effectively leveraging powerful language backbones, as exemplified by PALI (Chen et al., 2022), LLaVA (Liu et al., 2023), and CogVLM (Wang et al., 2024a). Subsequent work further enhances these models through improvements in both LLM backbones and large-scale vision-language pretraining (Hong et al., 2024a; Bai et al., 2025), while also expanding to additional modalities such as video and audio (Hurst et al., 2024). Notably, MLLMs demon-

strate strong capabilities in image perception and optical character recognition (OCR) (Hong et al., 2024b; Liu et al., 2024a), where multiple characters or words can be represented by a single visual token, highlighting the potential for effective context compression.

## 3 Method

We present Glyph, a novel paradigm for scaling long-context text understanding through visual compression. Unlike conventional long-context LLMs that extend token-based context windows, Glyph transforms ultra-long textual inputs into compact visual images and processes them with a vision–language model. This fundamentally different modeling method bypasses the prohibitive memory and computation costs of million-token sequences while preserving textual semantics. Furthermore, we introduce an LLM-driven genetic search to automatically discover optimal rendering configurations, ensuring the best trade-off between compression ratio and performance.

## 3.1 Overall Framework

As illustrated in Figure 2, Glyph consists of three tightly-coupled stages: (1) Continual Pre-Training, which teaches the VLM to understand and reason over rendered long texts with diverse visual styles; (2) LLM-Driven Rendering Search, automatically discovering the optimal rendering configuration for downstream tasks; and (3) Post-Training, including SFT and RL under the discovered configuration to further improve the model's long-context capabilities. Together, these stages enable Glyph to

achieve both high accuracy and significant gains in token compression, computational efficiency, and memory usage.

## 3.2 Task Definition

**Task Formulation.** We formalize the standard long-context instruction following task as a triple $(\mathcal{I}, \mathcal{C}, \mathcal{R})$, where $\mathcal{I}$ is a concise user instruction specifying the core goal, $\mathcal{C} = \{c_1, \ldots, c_T\}$ is an ultra-long textual context, and $\mathcal{R}$ is the target response. The conventional learning objective is to maximize

$$P(\mathcal{R} \mid \mathcal{I}, \mathcal{C}),$$

i.e., to generate an accurate response conditioned on both the instruction and the long textual context.

Scaling this token-based formulation to million-token contexts, however, imposes prohibitive memory and computation costs. To overcome these limitations, we reformulate the input representation through *visual compression*. Instead of directly feeding $\mathcal{C}$ as text tokens, we render it into a sequence of visual pages $\mathcal{V} = \{v_1, \ldots, v_n\}$, each containing the glyphs of multiple text segments. This allows the model to reason over a compressed but semantically equivalent input:

$$P(\mathcal{R} \mid \mathcal{I}, \mathcal{V}).$$

Each training instance is thus represented as $(\mathcal{I}, \mathcal{V}, \mathcal{R})$.

**Rendering Pipeline.** The rendering pipeline parameterizes how text is visualized before being fed into the model. Each rendering is specified by a configuration vector:

$$\boldsymbol{\theta} = \big(\texttt{dpi}, \texttt{page\_size}, \texttt{font\_family}, \texttt{font\_size},$$
$$\texttt{line\_height}, \texttt{alignment}, \texttt{indent}, \texttt{spacing},$$
$$\texttt{h\_scale}, \texttt{colors}, \texttt{borders}, \dots \big),$$

which controls typography, layout, and visual style of the rendered pages. Given the context $\mathcal{C}$ and configuration $\boldsymbol{\theta}$, the pipeline produces a sequence of images that serve as the VLM's long-context input.

To quantify the degree of compression, we define the compression ratio:

$$\rho(\boldsymbol{\theta}) = \frac{|\mathcal{C}|}{\sum_{i=1}^n \tau(v_i)},$$

where $\tau(v_i)$ denotes the number of visual tokens consumed by page $v_i$. A higher $\rho$ indicates that each visual token encodes more textual information, thus achieving stronger compression.

In practice, $\boldsymbol{\theta}$ determines both information density (through font size, dpi) and visual clarity (through layout and spacing). By varying $\boldsymbol{\theta}$, we can continuously adjust the balance between compression and readability for the VLM.

## 3.3 Continual Pre-Training

The purpose of continual pre-training is to transfer long-context comprehension from the textual to the visual modality. This stage exposes the VLM to a wide range of rendering styles and tasks so that it can align the semantics between rendered images and their corresponding texts.

**Data Construction.** To enhance model robustness, better aligning long-text capability, we adapt diverse rendering configurations over a large amount of long-context text data. We also develop a series of rules to exclude the improper combination of rendering parameters, e.g., a smaller line height than font size. Moreover, with human prior, we define several style themes, including *document_style*, *web_style*, *dark_mode*, *code_style*, and *artistic_pixel*. These themes are designed to capture a wide range of document layouts and text styles, which can better exploit the knowledge that VLM has obtained in its pre-training stage.

We further introduce three families of continual pre-training tasks, including:

- **OCR Tasks**: the model reconstructs all text on one or multiple rendered pages.
- **Interleaved Language Modeling**: certain text spans are rendered as images, while the rest remain in text, training the model to switch seamlessly between modalities.
- **Generation Tasks**: given partial rendered pages (e.g., the beginning or end), the model completes the missing parts.

These tasks jointly teach the model to read, reason, and generate under visually compressed contexts.

**Loss Function.** We minimize the cross-entropy loss

$$\mathcal{L}_{\text{CPT}} = -\mathbb{E}_{(\mathcal{I}^*, \mathcal{V}, \mathcal{R})} \sum_t \log P_\phi(r_t \mid \mathcal{I}^*, \mathcal{V}, r_{<t}), \tag{1}$$

where $\mathcal{I}^*$ denotes an optional instruction (e.g., absent in interleaved language modeling tasks) and $\phi$ is initialized from the base VLM. This stage

produces a model capable of understanding rendered text and handling long contexts, referred to as Glyph-Base.

## 3.4 LLM-Driven Rendering Search

Although diverse rendering improves generalization, downstream tasks often require a specific trade-off between compression and visual clarity for the VLM. We therefore perform an LLM-driven genetic search after continual pre-training to automatically identify the optimal rendering configuration $\theta^*$ used in the post-training stage.

**Genetic Algorithm.** Starting from an initial population of candidate configurations $\{\theta_k\}$ sampled from pre-training configurations, we iteratively perform the following steps:

1. **Rendering Data**: render the validation set using each configuration $\theta_k$ to obtain visual inputs.
2. **Evaluation on Validation Set**: perform model inference on the rendered data, measure task accuracy and compression ratio, and update the results.
3. **LLM Analysis & Critique**: use an LLM to suggest promising mutations and crossovers based on the current population and validation results.
4. **Search History**: record all configurations and their performance; rank and sample promising candidates for the next iteration.

This process continues until the population converges, i.e., when no further improvement is observed in validation accuracy or compression over a pre-defined number of generations. The resulting configuration $\theta^*$ is then adopted for post-training.

## 3.5 Post-Training

With the optimal rendering configuration $\theta^*$ fixed, we further improve Glyph-Base through two complementary optimization stages—*supervised fine-tuning* and *reinforcement learning*—supplemented by an *auxiliary OCR alignment* task. Together, these components jointly enhance the model's ability to reason over visually compressed inputs and to recognize textual details.

**Supervised Fine-Tuning.** To endow the model with robust comprehension under visual inputs, we curate a high-quality text SFT corpus and render its long-context inputs using the optimal configuration. Each response adopts a thinking-style format,

in which each example contains explicit reasoning traces (e.g., "`<think>...</think>`"). This encourages the model to perform step-by-step reasoning when reading massive token contexts.

Formally, the loss function can be written as

$$\mathcal{L}_{\text{SFT}} = -\mathbb{E}_{(\mathcal{I}, \mathcal{V}, \mathcal{R})} \sum_t \log P_\phi(r_t \mid \mathcal{I}, \mathcal{V}, r_{<t}),$$
(2)

where $\phi$ is initialized from the continual pre-training checkpoint. This stage establishes a strong initialization for reinforcement learning.

**Reinforcement Learning.** After SFT, we further refine the policy using Group Relative Policy Optimization (GRPO). For each input $(\mathcal{I}, \mathcal{V})$, we sample a group of candidate responses $\{r_1, \ldots, r_G\}$ from the old policy $\pi_{\phi_{\text{old}}}$. We first define the importance sampling weight:

$$w_i = \frac{\pi_\phi(r_i \mid \mathcal{I}, \mathcal{V})}{\pi_{\phi_{\text{old}}}(r_i \mid \mathcal{I}, \mathcal{V})}.$$
(3)

Each sampled response $r_i$ receives a reward score $u(r_i) \in \{0, 1\}$, which integrates:

- **Verifiable rewards** from an external LLM judge, scoring based on the accuracy of the answer, which is a reference-based LLM-as-a-judge with the reference being the ground truth.

- **Format rewards** that ensure the response correctly follows the defined thinking style.

The group-normalized advantage is computed as:

$$A_i = \frac{u(r_i) - \text{mean}(\{u(r_j)\}_{j=1}^G)}{\text{std}(\{u(r_j)\}_{j=1}^G)},$$
(4)

and the GRPO objective is

$$\mathcal{J}_{\text{GRPO}}(\phi) = \mathbb{E}_{(\mathcal{I}, \mathcal{V}) \sim P, \{r_i\}_{i=1}^G \sim \pi_{\phi_{\text{old}}}} \left[ \frac{1}{G} \sum_{i=1}^G \Big( \right.$$
$$\min\big(w_i A_i, \ \text{clip}(w_i, 1 - \epsilon_l, 1 + \epsilon_h) A_i\big)$$
$$\left. - \beta D_{\text{KL}}\big(\pi_\phi \| \pi_{\text{SFT}}\big) \Big) \right],$$
(5)

where $\epsilon$ and $\beta$ are hyperparameters.

**Auxiliary OCR Alignment.** A persistent challenge of visual compression is the faithful recovery of fine-grained text from rendered images.

| Model | Single-Doc QA | | Multi-Doc QA | | Summarization | | Few-shot | | Synthetic | | Code | | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | QP | NQA | HQA | 2QA | QSUM | GovRep | TREC | TriQA | PR Zh | PR En | RB | LCC | |
| GPT-4.1 | 51.60 | 35.73 | 69.10 | 74.15 | 23.50 | 33.36 | 77.00 | 93.36 | 100.00 | 100.00 | 67.94 | 68.43 | 56.03 |
| LLaMA-3.1-8B-Instruct | 44.56 | 26.34 | 56.88 | 46.67 | **23.28** | **32.36** | 19.25 | <u>89.12</u> | 62.20 | <u>99.50</u> | 42.81 | 46.35 | 41.34 |
| Qwen2.5-7B-Instruct-1M | **45.29** | 25.61 | 60.70 | 40.51 | <u>22.95</u> | <u>29.97</u> | 59.37 | 86.93 | <u>98.5</u> | **100.00** | 29.80 | 21.72 | 42.42 |
| Qwen3-8B | <u>44.67</u> | 26.13 | <u>65.83</u> | **73.92** | 19.60 | 26.85 | <u>70.50</u> | 87.98 | **100.00** | 97.26 | 40.89 | 44.87 | 47.46 |
| GLM-4-9B-Chat-1M | 43.75 | <u>26.72</u> | 58.98 | 50.89 | 22.84 | 27.60 | 61.50 | **90.07** | **100.00** | <u>99.50</u> | <u>55.64</u> | <u>59.54</u> | <u>49.27</u> |
| Glyph | 40.64 | **28.45** | **66.42** | <u>72.98</u> | 19.78 | 25.53 | **82.62** | 88.54 | 89.03 | <u>99.50</u> | **60.80** | <u>48.85</u> | **50.56** |

Table 1: Performance comparison of Glyph with leading LLMs on LongBench (%). Our model achieves competitive results in the overall average score. Best results are **bolded**, and second-best are <u>underlined</u>. Refer to Table 10 for the rest of the results.

| Model | 4 Needle | | | | | | 8 Needle | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0k-8k | 8k-16k | 16k-32k | 32k-64k | 64k-128k | Avg | 0k-8k | 8k-16k | 16k-32k | 32k-64k | 64k-128k | Avg |
| GPT-4.1 | 50 | 38 | 29 | 42 | 38 | 39.4 | 33 | 26 | 17 | 22 | 19 | 23.4 |
| LLaMA-3.1-8B-Instruct | <u>33.42</u> | <u>25.97</u> | <u>22.73</u> | **26.97** | 12.68 | <u>24.35</u> | <u>23.80</u> | 17.69 | **19.85** | <u>17.72</u> | 11.79 | **18.17** |
| Qwen2.5-7B-Instruct-1M | 25.96 | 20.13 | 19.93 | 24.25 | <u>17.29</u> | 21.51 | 17.64 | 19.48 | 12.41 | 14.80 | <u>14.24</u> | 15.71 |
| Qwen3-8B | 29.34 | 22.67 | 20.34 | 23.63 | **19.11** | 23.02 | 18.75 | <u>19.69</u> | <u>16.81</u> | **17.86** | **15.00** | 17.62 |
| GLM-4-9B-Chat-1M | 15.17 | 13.78 | 9.18 | 20.27 | 15.05 | 14.69 | 14.55 | 9.65 | 9.34 | 9.47 | 8.97 | 10.40 |
| Glyph | **35.44** | **26.82** | **24.15** | <u>25.69</u> | 16.37 | **25.81** | **25.12** | **21.22** | 16.43 | 13.91 | 13.51 | <u>18.14</u> |

Table 2: Performance comparison of our model against leading LLMs on the 4-needle and 8-needle sub-tasks of the MRCR benchmark (%). Our method consistently ranks first or second across most settings while preserving about 3× compression ratio. Performance on the 2-needle task is deferred to the Appendix.
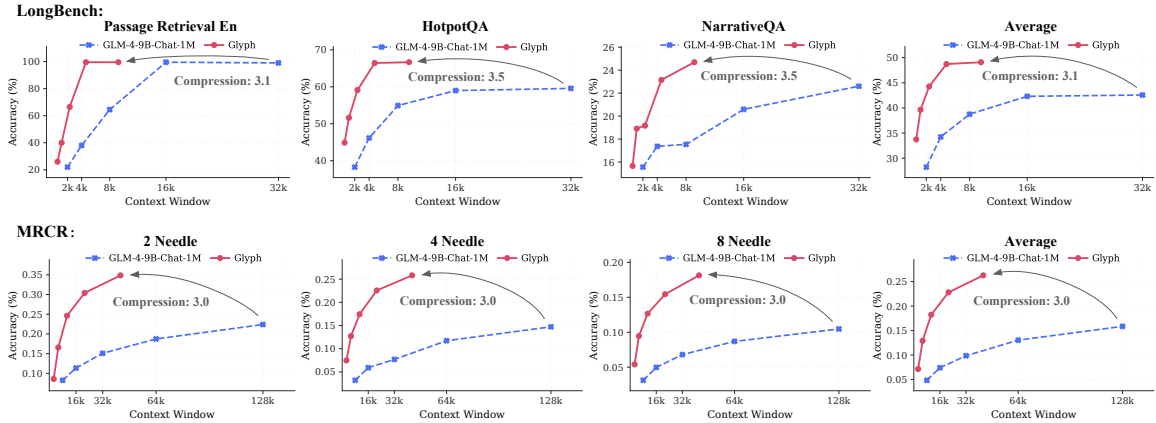


Figure 3: Performance comparison of Glyph and the baseline across different context windows, demonstrating that Glyph achieves performance equivalent to longer contexts with substantially shorter context windows.

Throughout both SFT and RL, we therefore incorporate an auxiliary OCR alignment task that encourages the model to correctly read and reproduce low-level textual details. The form of the OCR task is the same as in the continual pre-training stage. In the RL stage, the reward for the OCR task is given by the Levenshtein distance.

By integrating structured SFT supervision, RL optimization, and continuous OCR-aware alignment, Glyph acquires both powerful long-context reasoning ability and stable low-level text recognition, achieving strong downstream performance under highly compressed visual contexts.

## 4 Experiments

### 4.1 Experimental Setup

To comprehensively evaluate the effectiveness of our method, we have conducted extensive experiments covering long-context understanding, efficiency, cross-modal generalization, and several ablations and analysis. Implementation details, descriptions of baselines and benchmarks are provided in Appendix B.

### 4.2 Main Results on Performance

#### 4.2.1 Results on LongBench & MRCR

Tables 1 and 2 summarize overall results. Glyph achieves performance on par with or surpassing state-of-the-art text-only LLMs of similar size,

| Model | Niah-S1 | Niah-S2 | Niah-M1 | Niah-M2 | Niah-V | Niah-Q | VT | CWE | FWE | QA-1 | QA-2 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPT-4.1 | 100.0 | 98.85 | 100.0 | 100.0 | 99.67 | 100.0 | 100.0 | 97.87 | 98.66 | 86.82 | 77.47 | 96.30 |
| LLaMA-3.1-8B-Instruct | 99.33 | 99.33 | 99.33 | **99.00** | 98.17 | <u>99.67</u> | 87.07 | 57.30 | 81.85 | **84.00** | 58.00 | 87.55 |
| Qwen2.5-7B-Instruct-1M | **100.00** | <u>99.67</u> | <u>99.67</u> | **99.00** | 93.83 | 98.75 | 85.40 | 72.10 | 85.67 | <u>80.00</u> | 60.67 | 88.61 |
| Qwen3-8B | **100.00** | **100.00** | 95.33 | 84.67 | 97.42 | 99.33 | <u>98.47</u> | 74.67 | 86.67 | 70.33 | 53.33 | 87.29 |
| GLM-4-9B-Chat-1M | **100.00** | **100.00** | 92.67 | **99.00** | 95.00 | **100.00** | 98.20 | 49.50 | 83.22 | 72.67 | 56.67 | 86.08 |
| **DPI: 72 / Compression rate: average 4.0, up to 7.7** | | | | | | | | | | | | |
| Glyph | 73.33 | 64.67 | 67.33 | 56.00 | 73.42 | 71.42 | 77.93 | 94.40 | 92.67 | 59.33 | 63.33 | 72.17 |
| **DPI: 96 / Compression rate: average 2.2, up to 4.4** | | | | | | | | | | | | |
| Glyph | 98.00 | 95.33 | 95.67 | 85.00 | 96.33 | 95.83 | 94.93 | <u>94.80</u> | <u>98.00</u> | 79.00 | <u>70.67</u> | <u>91.23</u> |
| **DPI: 120 / Compression rate: average 1.2, up to 2.8** | | | | | | | | | | | | |
| Glyph | <u>99.67</u> | 99.00 | **100.00** | <u>93.67</u> | 99.00 | 99.58 | **99.33** | **98.97** | **99.11** | 79.00 | **74.00** | **94.67** |

Table 3: Performance on the Ruler benchmark (%). We demonstrate the impact of different DPI settings on our model's performance and the resulting compression ratios. For each configuration, the table includes both the average compression ratio across all sub-tasks and the maximum compression achieved for specific sub-task types.



Figure 4: Speedup ratios of Glyph over the text backbone model for prefill, decoding, and training across different sequence lengths.



Figure 5: Model performance degradation across different sequence lengths on the Ruler benchmark.

including Qwen3-8B and GLM-4-9B-Chat-1M, demonstrating that Glyph remains effective on long-context tasks with a large reduction in input tokens.

Figure 3 further illustrates the context scaling behavior of Glyph. For LongBench, we report the results with truncated contexts; for MRCR, we utilize the original dataset split. On LongBench, our model achieves an average effective compression ratio of 3.3×, with certain tasks reaching up to around 5×. On MRCR, the average compression ratio is 3.0×. This means that within the same token budget, Glyph can effectively utilize several times more original context than text-only models.

More importantly, as the input length grows, this advantage scales up. When a text-only model extends its window from 32k to 64k tokens, it gains 32k additional tokens of usable context. Under the same expansion, Glyph —with a compression ratio of around 3×—effectively gains about 96k tokens' worth of original text. This advantage translates into a faster improvement as the context length increases.

### 4.2.2 Results on Ruler

On the Ruler benchmark, Glyph also achieves performance comparable to leading LLMs across most categories (Table 3). We exclude the UUID task from this benchmark due to its huge difficulty for VLMs, which is further discussed in the limitations section.

Beyond raw scores, we demonstrate the advantage of test-time scaling. When we increase the rendering resolution (DPI) at inference time, our model shows substantial gains: at higher DPI settings, it even surpasses strong text-only baselines. This demonstrates that the performance of VLMs on text-only long-context tasks has a high ceiling, and that Glyph still holds considerable potential.

Furthermore, we analyze performance under dif-

7

| Model | SP | CP | UA | Acc | F1 |
|---|---|---|---|---|---|
| GLM-4.1V-9B-Base | 36.76 | 23.41 | 21.52 | 29.18 | 28.78 |
| Glyph-Base | 47.91 | 22.24 | 14.80 | 32.48 | 34.44 |
| Glyph | **57.73** | **39.75** | **27.80** | **45.57** | **46.32** |

Table 4: Results on MMLongBench-Doc (%). SP, CP, UA, and Acc denote Single-page, Cross-page, Unanswerable, and Overall Accuracy, respectively.

| Configuration | LongBench | MRCR | Ruler | Avg. |
|---|---|---|---|---|
| Random Config | 41.78 | 15.82 | 65.13 | 40.91 |
| Manual Config | **43.45** | 19.33 | 68.09 | 43.62 |
| Search-based Config | **43.45** | **22.10** | **71.24** | **45.60** |

Table 5: Ablation study comparing randomly combined, manually designed, and search-based configurations on three benchmarks under SFT setting. The search-based configuration achieves the best overall performance.

| Model | LongBench | MRCR | Ruler |
|---|---|---|---|
| Glyph | 50.56 | 26.27 | 72.17 |
| – w/o OCR (in RL) | -1.40 | -2.00 | -0.35 |
| – w/o RL | -7.11 | -4.17 | -0.93 |
| – w/o OCR (in SFT) | -8.12 | -8.42 | -1.23 |

Table 6: Ablation study showing the performance drop (%) relative to the final Glyph model when components are progressively removed.

| Model | 2 Needle | 4 Needle | 8 Needle |
|---|---|---|---|
| GLM-4-9B-Chat-1M | 10.08 | 6.19 | 2.26 |
| Qwen2.5-7B-Instruct-1M | **11.36** | 7.34 | **7.77** |
| Glyph | 9.36 | **7.62** | 7.64 |

Table 7: Average MRCR performance (%) across 128K–1M context lengths under different needle counts.

ferent sequence lengths (Figure 5). At short contexts, text-only models such as LLaMA-3.1-8B-Instruct maintain a slight edge. However, as the input length grows, Glyph exhibits obviously slower degradation. This aligns with the earlier observations on LongBench and MRCR: thanks to compression, an increase in the nominal text context window translates to a much smaller increase in the effective length the Glyph model actually needs to handle. Consequently, our model maintains accuracy more stably as the context grows.

## 4.3 Efficiency Evaluation

We further evaluate the efficiency of our method in both training and inference, comparing Glyph with the text backbone model. The evaluation setting is detailed in Appendix B. As shown in Figure 4, Glyph provides clear speedups in both metrics, demonstrating significant gains at the inference stage and SFT training stage. As the sequence length grows from 8k to 128k, our model demonstrates markedly better scalability, achieving stable SFT training throughput speedup and growing inference speedup.

## 4.4 Cross-Modal Generalization

Although our training data mainly consists of rendered text images rather than natural multimodal inputs, we are interested in whether such training can generalize to real-world multimodal tasks, like long document understanding. To this end, we evaluate Glyph on the MMLongBench-Doc benchmark, which contains 130 long PDF documents with diverse layouts and embedded images. As shown in

Table 4, Glyph achieves clear improvements over our backbone model GLM-4.1V-9B-Base, confirming its ability to generalize across modalities.

## 4.5 Ablation Study & Analysis

We conduct a series of ablations and analyses to better understand our method.

**Configuration Search.** We compare three types of rendering configurations for SFT: (i) randomly sampled configuration from the pre-training sets, (ii) manually designed settings based on prior knowledge, and (iii) the configuration obtained from our search procedure. While all settings achieve comparable compression ratios, Table 5 shows that the searched configuration consistently outperforms the other two, both on average and across most individual tasks. This demonstrates the importance of systematic exploration for finding appropriate rendering strategies.

**OCR Auxiliary Tasks.** We also test the impact of adding OCR auxiliary tasks during both SFT and RL training. As shown in Table 6, including OCR objectives yields consistent performance gains across benchmarks. This suggests that explicitly reinforcing low-level text recognition helps the model build stronger representations, which in turn improves long-context understanding ability.

**Extreme Compression Exploration** To further examine the potential of our approach, we explore more aggressive compression settings. We apply a configuration with an effective $8\times$ compression ratio during post-training, and evaluate the resulting

model on MRCR with sequence lengths extended from 128k to 1024k. The results (Table 7) show that Glyph successfully demonstrates the potential for $8\times$ effective context expansion, achieving performance on par with GLM-4-9B-Chat-1M and Qwen2.5-1M. This experiment highlights that our method can indeed be pushed to more extreme compression regimes while retaining performance, suggesting substantial headroom for extending usable context far beyond current limits, like a model that can deal with 4M, even 8M context tokens.

## 5 Conclusion

In this work, we present Glyph, an efficient long-context modeling framework that renders long texts into compact images and processes them with vision-language models. With continual pre-training, an LLM-driven genetic rendering search and targeted post-training, Glyph achieves 3–4× context compression while maintaining competitive performance with similar size leading LLMs such as Qwen3-8B. Extensive experiments further demonstrate substantial gains in inference speed and memory efficiency, and show that our method demonstrates cross-modal benefits, enhancing multimodal long-context tasks like document understanding. Our findings demonstrate that enhancing token information density constitutes a promising new paradigm for scaling long-context LLMs, orthogonal to existing attention-based approaches, and there remains great room for further exploration in depth.

### Limitations and Future Work

Despite the effectiveness of Glyph and its strong potential for broader applications, we want to discuss several limitations of the current work that are worth further exploration.

**Sensitivity to rendering parameters.** Our method relies on rendering textual inputs into images before processing. We find that performance can be noticeably affected by rendering configurations such as resolution, font, and spacing. Although our search procedure allows us to identify a configuration that performs well on downstream tasks, how to make the model more robust across various rendering settings remains an open problem.

**OCR-related challenges.** As discussed in the Ruler benchmark, UUID recognition remains par-

ticularly challenging for current VLMs, and even the strongest models (e.g., Gemini-2.5-Pro) often fail to reproduce them correctly. Such rare alphanumeric sequences frequently result in misordered or misclassified characters, which may stem from their distributional sparsity in training data or from architectural limitations of visual encoders. While these cases have little impact on most tasks, improving OCR fidelity could push the upper bound of our approach.

**Task diversity.** The benchmarks in this work mainly focus on long-context understanding. While these tasks provide a strong proof of concept, they do not fully capture the diversity of real-world applications, such as agentic or reasoning-heavy tasks. We also observe that, compared with textual models, the visual-text model tends to generalize less effectively across tasks. Extending the scope of evaluation and training to a wider range of tasks will help better assess and improve the robustness and generality of our approach.

**Future directions.** Building upon the current study, several directions could further advance the proposed visual–text compression paradigm. First, rather than using a fixed rendering strategy, one promising avenue is to train adaptive rendering models that condition on the task type or user query, producing tailored visualizations that balance compression and performance. Second, enhancing the visual encoder's capability for fine-grained text recognition and alignment with language representations could improve robustness and transferability across tasks. Third, improving the alignment between visual–text and purely textual models, for instance, through knowledge distillation or cross-modal supervision, could narrow the performance gap in generalization. Fourth, our approach could be extended to broader applications, such as agent memory systems capable of managing long-term conversations or agentic contexts, and tasks that can leverage structured visual layouts for reasoning and retrieval. From the perspective of context engineering, this method offers a new way to optimize how contextual information is represented and managed. With further advances along this line, future models could go beyond the current limits of context length, effectively scaling from 1M to 10M input tokens.

# References

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, and 1 others. 2024. Longbench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. In *Proceedings of ACL*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA. Curran Associates Inc.

Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu, Chao Wang, Cheng Zhu, and 1 others. 2025a. Minimax-m1: Scaling test-time compute efficiently with lightning attention. *arXiv preprint arXiv:2506.13585*.

Xi Chen, Xiao Wang, Soravit Changpinyo, Anthony J Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, and 1 others. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.

Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024. Longlora: Efficient fine-tuning of long-context large language models. In *The International Conference on Learning Representations (ICLR)*.

Yutao Chen, Yiren Wang, and 1 others. 2025b. Cope: Complex positional encoding for long context extrapolation. *arXiv preprint arXiv:2508.18308*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, and 1 others. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, and 38 others. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *Preprint*, arXiv:2406.12793.

Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, and 1 others. 2024a. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*.

Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, and 1 others. 2024b. Cogagent: A visual language model for gui agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14281–14290.

Yunpeng Huang, Jingwei Xu, Junyu Lai, Zixu Jiang, Taolue Chen, Zenan Li, Yuan Yao, Xiaoxing Ma, Lijuan Yang, Hao Chen, and 1 others. 2023. Advancing transformer architecture in long-context large language models: A comprehensive survey. *arXiv preprint arXiv:2311.12351*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Philippe Laban, Alexander Richard Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. Summary of a haystack: A challenge to long-context llms and rag systems. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9885–9903.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. Llava-next: Improved reasoning, ocr, and world knowledge.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.

Yang Liu, Wei Gao, and 1 others. 2024b. Long context is not long at all: A prospector of long-dependency data for large language models. *arXiv preprint arXiv:2407.11234*.

Baolin Peng, Zhuohan Li, and 1 others. 2023. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*.

Bo Peng, Ruichong Zhang, Daniel Goldstein, Eric Alcaide, Xingjian Du, Haowen Hou, Jiaju Lin, Jiaxing Liu, Janna Lu, William Merrill, and 1 others. 2025.

Rwkv-7" goose" with expressive dynamic state evolution. *arXiv preprint arXiv:2503.14456*.

Ofir Press, Noah A. Smith, and Mike Levy. 2021. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*.

Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*.

Zhen Sun, Peng Cheng, Wei He, and 1 others. 2022. Xpos: Improving position interpolation with extrapolation. *arXiv preprint arXiv:2212.10554*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Kiran Vodrahalli, Santiago Ontanon, Nilesh Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna, Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, and 1 others. 2024. Michelangelo: Long context evaluations beyond haystacks via latent structure queries. *arXiv preprint arXiv:2409.12640*.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Song XiXuan, and 1 others. 2024a. Cogvlm: Visual expert for pretrained language models. *Advances in Neural Information Processing Systems*, 37:121475–121499.

Yizhi Wang, Fan Yang, and 1 others. 2024b. Longrecipe: Recipe for efficient long context generalization in large language models. *arXiv preprint arXiv:2406.12345*.

Yingsheng Wu, Yuxuan Gu, Xiaocheng Feng, Weihong Zhong, Dongliang Xu, Qing Yang, Hongtao Liu, and Bing Qin. 2024. Extending context window of large language models from a distributional perspective. *arXiv preprint arXiv:2410.01490*.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. 2024. Gated linear attention transformers with hardware-efficient training. In *International Conference on Machine Learning*, pages 56501–56523. PMLR.

Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proceedings of NAACL*.

Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiying Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, and 1 others. 2025a. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. *arXiv preprint arXiv:2507.02259*.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025b. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.

Tianle Zhang, Zhuohan Li, and 1 others. 2024. Longalign: Instruction-tuning long-context llms. *arXiv preprint arXiv:2401.10968*.

Wei Zhu, Ziheng Wang, and 1 others. 2024. Data-adaptive positional encoding for length generalization. *NeurIPS*.

# A Rendering Parameters

Our rendering parameters are detailed in Table 8. The best result, obtained through our LLM-driven genetic search, is presented in Figure 6, which shows the detailed configuration and its corresponding rendering.

# B Implementation Details

**Training Details** For continual pre-training of the 9B long-context backbone, the model is initialized from the released GLM-4.1V-9B-Base checkpoint and trained on a diverse mixture of rendered long-context data and vision-language corpora (e.g., OCR task) within 128k context length. The training uses a global batch size of 170 and a learning rate of 2e-6 with cosine decay for around 4000 steps.

For the rendering search, we run for 5 times with 200 steps in each round, to find the optimal configuration that maximizes the compression ratio while maintaining good performance.

After this, we conduct further SFT and RL training. For SFT, we train for 1.5k steps with a batch size of 32. The Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.95$) is used with cosine decay and 160 warm-up steps, where the learning rate decays from 5e-6 to 2e-6. For reinforcement learning, we adopt the GRPO algorithm. Each training group samples 16 candidate responses, and degenerate samples with all-zero or all-one rewards are discarded. We apply the clip-higher trick from DAPO (Yu et al., 2025b) with $\epsilon_l$ being 0.2 and $\epsilon_h$ being 0.28. Training runs for 500 iterations with a batch size of 32. We also use the Adam optimizer with a constant learning rate of 1e-6.

| Factor | Specification / Sampling Strategy |
|---|---|
| dpi | Mixture of sets: *lowest* (45–59), *low* (60–71), *medium* (72–119), *normal* ({72,80,96,100,110,120,144,150,300}), *high* (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, `italic`/`oblique`). |
| font_size | $\{7, 7.5, 8, 9, 9.5, 10, 11, 12, 14\}$; line_height tied as font_size + $\{0, \dots, 3\}$. |
| alignment | LEFT/JUSTIFY (dominant) with small-prob. RIGHT/CENTER. |
| margins | Three patterns: all-equal; vertical-larger; horizontal-larger; values in 10–40pt ranges. |
| indent | Modes: none; first-line indent ($\approx$1–2.5 em); block/hanging with left/right indents. |
| spacing | space-before/space-after 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., \n, tags, or tokens) inserted to preserve structure. |
| auto_crop | Optional white-margin cropping and last-page trimming. |

Table 8: Controllable factors in the rendering pipeline and their sampling strategies. The mixture design yields broad yet realistic typography/layout coverage and tunable compression $\rho(\boldsymbol{\theta})$.

| Model | 2 Needle | | | | | |
|---|---|---|---|---|---|---|
| | 0k-8k | 8k-16k | 16k-32k | 32k-64k | 64k-128k | Avg |
| GPT-4.1 | 83 | 72 | 67 | 62 | 59 | 68.6 |
| LLaMA-3.1-8B-Instruct | <u>54.27</u> | **53.21** | **51.05** | <u>29.81</u> | <u>24.98</u> | <u>42.66</u> |
| Qwen3-8B | **58.95** | 41.18 | 36.18 | 24.99 | 20.89 | 36.44 |
| GLM-4-9B-Chat-1M | 39.77 | 15.87 | 18.42 | 18.63 | 18.42 | 22.22 |
| Qwen2.5-7B-Instruct-1M | 45.92 | <u>51.07</u> | <u>46.97</u> | **34.67** | **37.57** | **43.24** |
| Glyph | 41.51 | 40.78 | 39.58 | 29.67 | 22.41 | 34.85 |

Table 9: Performance of various models on the MRCR task (%) with the 2 Needle setting across different context length intervals (0k–8k, 8k–16k, 16k–32k, 32k–64k, 64k–128k) and the average score.

**Baselines.** We compare Glyph with leading open-sourced LLMs of similar size:

- **Qwen3-8B** achieves state-of-the-art performance across reasoning and a wide range of tasks.

- **Qwen2.5-7B-Instruct-1M** excels at long-context understanding and achieves strong performance across diverse benchmarks.

- **LLaMA-3.1-8B-Instruct** is a widely used model with strong instruction-following and multilingual capabilities.

- **GLM-4-9B-Chat-1M** delivers powerful long-context tasks and overall high performance.

**Backbone Model.** Our method relies on a strong VLM to process long-context tasks. Considering the impressive performance of GLM-4.1V-9B, especially in OCR and long document tasks, we have chosen GLM-4.1V-9B-Base as our backbone model.

**Evaluation Benchmarks.** To conduct a comprehensive analysis of the long-context performance, we have adopted three popular benchmarks, including LongBench, MRCR, and Ruler. Long-Bench consists of 21 datasets in total in 6 categories, covering diverse long-context tasks. MRCR is a task proposed by Vodrahalli et al. (2024). We use the OpenAI version, which consists of multi-turn conversations about writing, asking models to recall one of the contexts in dialogue history. Ruler is a widely used synthetic benchmark with 11 NIAH tasks. To validate cross-model benefits, we choose the MMLongBench-Doc, which involves 130 lengthy PDF with diverse layout and images, and 1062 questions.

**Efficiency Evaluation Setting.** For training, we focus on SFT since RL involves rollout time, making it difficult to compare fairly. Moreover, running RL at very long context lengths (e.g., 64k or beyond) requires excessive memory resources. This again highlights the advantage of our approach: through compression, we can conduct RL training at 32k context length while effectively covering over 100k tokens of raw text input, where RL is prohibitively difficult for LLMs due to memory and computation demands. For SFT, we measure per-sample training time under the same number of training data using 8×80G H100 GPUs. For inference, we deploy both models on a single 80G

| Model | Single-Doc QA | | Multi-Doc QA | | Summarization | | Few-shot | | Synthetic |
|---|---|---|---|---|---|---|---|---|---|
| | QA Zh | QA En | Mus | Dur | News | VcSum | Sam | Lsht | Pa C |
| GPT-4.1 | 63.90 | 51.27 | 55.63 | 24.58 | 23.70 | 14.66 | 41.25 | 50.00 | 26.5 |
| LLaMA-3.1-8B-Instruct | 62.20 | **54.98** | 31.61 | **33.75** | **24.21** | **16.23** | 7.61 | 0.00 | 7.13 |
| Qwen3-8B | 60.98 | 49.78 | <u>45.54</u> | 16.69 | 18.55 | 12.08 | <u>36.47</u> | 42.00 | <u>12.81</u> |
| GLM-4-9B-Chat-1M | **63.17** | 52.88 | 39.14 | <u>28.27</u> | <u>23.90</u> | <u>16.21</u> | 36.15 | **47.38** | 2.39 |
| Qwen2.5-7B-Instruct-1M | <u>62.98</u> | <u>53.62</u> | 34.72 | 21.85 | 21.02 | 12.20 | **39.17** | 28.68 | 3.50 |
| Glyph | 37.23 | 45.89 | **56.18** | 26.87 | 21.52 | 12.43 | 32.49 | <u>44.43</u> | **30.50** |

Table 10: The rest of the results on LongBench benchmark (%), which encompasses Single-Document QA, Multi-Document QA, Summarization, Few-shot Learning, and Synthetic task.

H100 and measure efficiency along two axes: (i) prefill latency at batch size 1, and (ii) per-sample inference time at the maximum feasible batch size with output length set to 256 tokens. We omit the KV cache testing, because KV cache scales linearly with the sequence length, and compression translates almost directly into savings of about 67% memory usage.

| Best Config | |
|---|---|
| page-size | 595,842 |
| dpi | 72 |
| margin-x | 10 |
| margin-y | 10 |
| font-path | Verdana.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-after | 0 |
| space-before | 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 6: The optimal parameter setting. The left column lists the values for page layout, font, and spacing, while the right column provides an example of the rendered text.