# SampleAttention: Near-Lossless Acceleration of Long Context LLM Inference with Adaptive Structured Sparse Attention

Qianchao Zhu [*1]  Jiangfei Duan [*2]  Chang Chen [1]  Xiuhong Li [1]  Siran Liu [1]  Guanyu Feng [3]  Xin Lv [3]
Chuanfu Xiao [1]  Dahua Lin [2]  Chao Yang [1]

## Abstract

Large language models (LLMs) now support extremely long context windows, but the quadratic complexity of vanilla attention results in significantly long Time-to-First-Token (TTFT) latency. Exisiting sparse attention approaches employ either static sparse pattern or fixed sparsity ratio to utilize the high attention sparsity, failing to capture the adaptive sparsity ratio and dynamic sparse pattern across attention heads, input contents and model architectures. To balance accuracy and performance efficiently, we introduce a robust indicator for accuracy, Cumulative Residual Attention (CRA), which measures the percentage of attention recall. Leveraging this key insight, we present SampleAttention, which employs a novel two-stage query-guided key-value filtering approach to efficiently and dynamically select a minimal set of important column and slash strips to meet a desired CRA threshold, thus maximizing efficiency while preserving accuracy. Comprehensive evaluations show that SampleAttention can establish a new Pareto frontier in the accuracy-efficiency trade-off, and reduces TTFT by up to $5.29\times$ compared with FlashAttention2.

## 1 Introduction

Recent advances (Xiong et al., 2023; Liu et al., 2023a; Chen et al., 2023b; Li et al., 2023a; Chen et al., 2023a) race to scale the context window of large language models (LLMs) (Brown et al., 2020; Vaswani et al., 2017; Touvron et al., 2023) for more complex applications, including document analysis (Zhang et al., 2024a), code copilot (Chen et al., 2021c; Roziere et al., 2023), and prolonged conversations (Chiang et al., 2023; Taori et al., 2023). Popular LLMs like Gemini (Team et al., 2023), Claude (Anthropic, 2023) and Kimi (Moonshot, 2023) now support context lengths exceeding 1 million tokens. However, the increase in context length makes it challenging to support live interactions due to the quadratic complexity of attention mechanism. As illustrated in Figure 1, the attention computation time increases quadratically with sequence length, quickly dominating the *Time to First Token* (TTFT) latency (i.e. prefill latency). For example, in a 1 million token context, the attention of ChatGLM3-6B (Du et al., 2021) takes 1555 seconds, constituting over 90% of the TTFT when evaluated on an A100 GPU.

Prior work has consistently demonstrated that attention

---
[*]Equal contribution  [1]Peking University  [2]The Chinese University of Hong Kong  [3]Zhipu.AI. Correspondence to: Chao Yang <chao_yang@pku.edu.cn>.
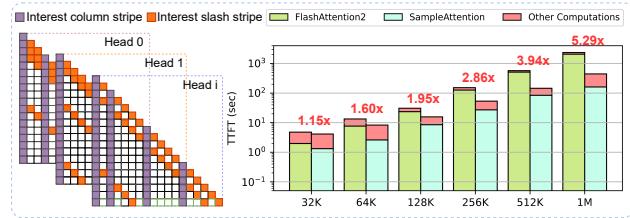
*Figure 1.* Compared to previous static and dynamic sparse attention methods, SampleAttention captures adaptive structured sparse patterns for each head. It achieves a significant reduction in TTFT compared to FlashAttention2.

scores exhibit high sparsity (Zaheer et al., 2020; Kitaev et al., 2020; Jiang et al., 2024; Li et al., 2024). This characteristic makes sparse attention a promising approach for reducing prefill latency, as it allows the model to compute attention selectively for only the most important query and key-value tokens rather than the entire sequence. Based on these observations, plenty of approaches propose to approximate the dense attention with static sparse pattern, like Long-Former (Beltagy et al., 2020), BigBird (Zaheer et al., 2020), LongNet (Ding et al., 2023) and StreamingLLM (Xiao et al., 2023b). However, these approaches fail to capture the dynamic sparse pattern across heads and inputs (Jiang et al., 2024; Likhosherstov et al., 2021) and cannot achieve the same accuracy of full attention.

Recent work proposes to address this dynamic sparse pattern through runtime attention index selection (Liu et al.,

2022; Han et al., 2023; Jiang et al., 2024). DSA (Liu et al., 2022) approximates attention patterns using low-rank hidden dimensions, but incurs significant computational overhead with long contexts. While HyperAttention (Han et al., 2023) employs Locality Sensitive Hashing (LSH) to identify important attention scores, its coarse-grained selection approach struggles to maintain model accuracy. MInference (Jiang et al., 2024) takes a hybrid approach: it predefines several sparse patterns and matches each attention head to its optimal pattern offline to achieve a target sparsification budget, then dynamically searches for sparse indices during the prefill phase. However, this approach's reliance on predefined patterns and fixed budgets fails to capture both the varying sparsity ratios across attention heads and the dynamic sparse patterns that adapt to different input contents.

In this paper, we find effectively exploiting the inherently-high attention sparsity to accelerate prefill computation remains challenging due to two key factors. First, the optimal sparsity ratio varies adaptively across attention heads, input contents, and model architectures, making it necessary to be determined at runtime. Second, attention patterns also varies across heads and contents, often combining typical column and slash patterns. Some attention heads even display intricate combinations of these patterns, further complicating sparse pattern selection. These dynamic characteristics create significant challenges for existing methods to achieve an optimal trade-off due to their lack of flexibility. This highlights the need for a more adaptable, runtime-efficient approach to determine both the sparsity ratio and the pattern.

To address these challenges, we propose a novel approach, SampleAttention, which can dynamically determine sparse ratios and patterns at runtime. To select a minimal set of significant column and slash patterns while maintaining accuracy, we introduce a robust metric for evaluating model accuracy called *Cumulative Residual Attention* (**CRA**), which measures the capability of attention recall. The details of this key insight will be presented in *Section 3*. Based on CRA, we elaborate on a two-stage query-guided key-value filtering method proposed by SampleAttention in *Section 4*. This implementation is designed to efficiently identify important columns and slashes at runtime. SampleAttention also develops an automated tuning method that uses a small profiling dataset to determine the optimal hyperparameter setting for each model across different length ranges. SampleAttention significantly accelerates vanilla attention by reducing both I/O and computation requirements. We also implement hardware-efficient kernels. Notably, SampleAttention aims to reduce the computation overhead of attention, and is orthogonal and can be combined with existing KV cache eviction approaches (Zhang et al., 2024c; Ribar et al., 2023; Mu et al., 2024) to further reduce memory consumption.

We evaluate SampleAttention on ChatGLM (GLM et al., 2024), YI (Young et al., 2024) and InternLM (Cai et al., 2024) with a suite of popular benchmarks covering various generative tasks across different sequence lengths. Experimental results show that SampleAttention achieves nearly no accuracy loss[1] for different LLMs, significantly outperforming prior works, and reduces the TTFT by up to $5.29\times$ compared with FlashAttention2.

## 2 BACKGROUND

### 2.1 LLM Inference

LLMs are built upon transformer architectures (Vaswani et al., 2017), which stacks multiple transformer decoder blocks. Each block consists of an attention layer (Self-Attention) followed by a feed-foward network (MLP).

The inference process of LLMs operates in two phases: *prefill* and *decoding*. During the prefill phase, the model processes the entire input prompt in parallel and generates the first output token. This phase also generates and stores the Key-Value (KV) cache for each token in the prompt, which will be used in subsequent computations. The decoding phase follows the prefill phase and sequentially generates each new token based on all previous tokens. The model takes one output token as input each time, and leverages the KV cache to generate the subsequent new token autoregressively. The KV cache of the input token will also be saved as the context for subsequent generation.

When processing long input sequences, the computational demands of handling lengthy prompts can result in significant *Time To First Token* (TTFT) latency (i.e. prefill latency), creating a substantial bottleneck for real-world applications. For example, the TTFT of 1 million sequence for ChatGLM-6B (Du et al., 2021) takes near 30 minutes. This prohibitive latency makes it impractical for applications requiring real-time responses. Therefore, reducing TTFT for long sequences becomes crucial for enabling practical deployment of LLMs in scenarios demanding both long context processing and responsive interaction.

### 2.2 Attention Computation

The attention layer enables the model to weigh the importance of different tokens in the input sequence and dynamically adjust their influence on the output. We start with a regular full attention for one attention head to examine the mechanism, while the following contents can be seamlessly applied to multiple attention heads.

In attention layer, each token in the input sequence is transformed into three vectors: query, key and value tensors. Let

---

[1]Near-lossless refers to that model accuracy stays above 99% of the baseline according to MLPerf (Reddi et al., 2020).
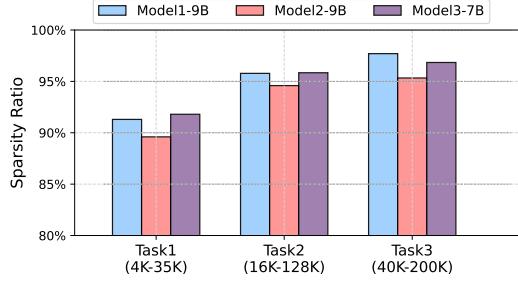
*Figure 2.* The average sparsity ratio of three different models with long-context window on tasks with varying length ranges.

$\mathbf{Q} \in \mathbb{R}^{S_q \times d}$ and $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{S_k \times d}$ be the query and key-value tensor of one head, where $S_q, S_k$ is the sequence length respectively, and $d$ is the head dimension. The full attention output $\mathbf{O} \in \mathbb{R}^{S_q \times d}$ can be formulated as,

$$\mathbf{P} = \texttt{softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}) \in [0,1]^{S_q \times S_k} \quad (1)$$

$$\mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{S_q \times d} \quad (2)$$

where $\texttt{softmax}$ is applied in row-wise, and $\mathbf{P}$ is the attention score.

This attention mechanism poses a fundamental challenge during long-context inference: both the memory footprint of $\mathbf{P}$ and the computational complexity scale quadratically with sequence length. While FlashAttention (Dao et al., 2022) effectively addresses the memory bottleneck through online softmax computation, the quadratic computational complexity remains unresolved, resulting in substantial response delays. As previously discussed, attention computation can constitute over $90\%$ of the total TTFT latency, making its optimization crucial for achieving practical long-context inference.

### 2.3 Inherently-High Attention Sparsity

In attention computation, applying $\texttt{softmax}$ over long sequences tends to reduce the influence of smaller elements, making them less significant. This insight motivates us to investigate the inherent sparsity in the attention scores, which can potentially accelerate the attention mechanism without compromising accuracy. Formally, the full attention score matrix $\mathbf{P}$ can be approximated with sparse attention score $\hat{\mathbf{P}}$,

$$\hat{\mathbf{P}} = \texttt{softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} - c(1 - \mathbf{M})) \quad (3)$$

where $\mathbf{M} \in \{0,1\}^{S_q \times S_k}$ is a binary mask matrix that determines the sparse attention pattern, and $c$ is a large constant that effectively zeroes out masked attention scores after the $\texttt{softmax}$ operation. The sparsity ratio measures the percentage of attention scores that are masked.
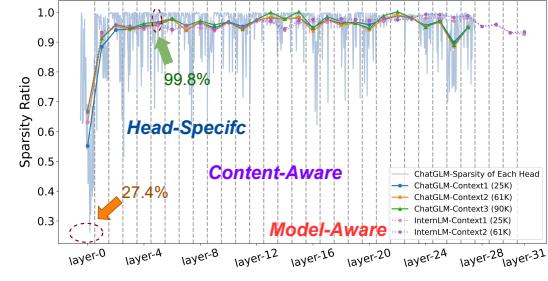


*Figure 3.* The sparsity ratio of ChatGLM3 (28 layers×32 heads) and InternLM2 (32 layers×32 heads), evaluated over different tasks during prefill. The sparsity ratio varies across different attention heads, input contents and model architectures.

Our observations reveal that LLMs inherently exhibit a significant sparsity ratio in their attention computations, even without explicit optimization for this property. This finding emerges from our comprehensive evaluation of attention patterns across different model architectures and input prompts, as shown in Figure 2. These observations suggest that attention sparsity is an intrinsic characteristic of how LLMs process information during the prefill phase, and the average sparsity ratio is high across different models and datasets, suppressing $89.6\%$.

## 3 MOTIVATION

While the high attention sparsity has been discussed in recent works (Jiang et al., 2024; Li et al., 2024; Xiao et al., 2024), we find that the adaptive sparsity ratio and dynamic sparse pattern make it challenging to effectively exploit this peoperty to accelerate prefill attention computation for long context LLMs.

### 3.1 C#1: Adaptive Sparsity Ratio

Our observations reveal that attention sparsity in LLMs exhibits adaptive sparsity ratio across three dimensions: attention heads, input contents, and model architectures (Figure 3).

- *Head-Specific:* different attention heads exhibit remarkably different sparsity ratios, even within the same layer. Speciafically, one head in the first layer has a sparsity ratio as low as $27.4\%$, while the highest can reach $99.8\%$.

- *Content-Aware:* the sparsity ratios of different input prompts and context lengths are different. Our observation indicates that as the context becomes longer, the sparsity ratio increases correspondingly (Appendix A.2).

- *Model-Aware:* different models exhibit distinct sparsity ratios (Figure 2, Figure 3, Appendix A.2).

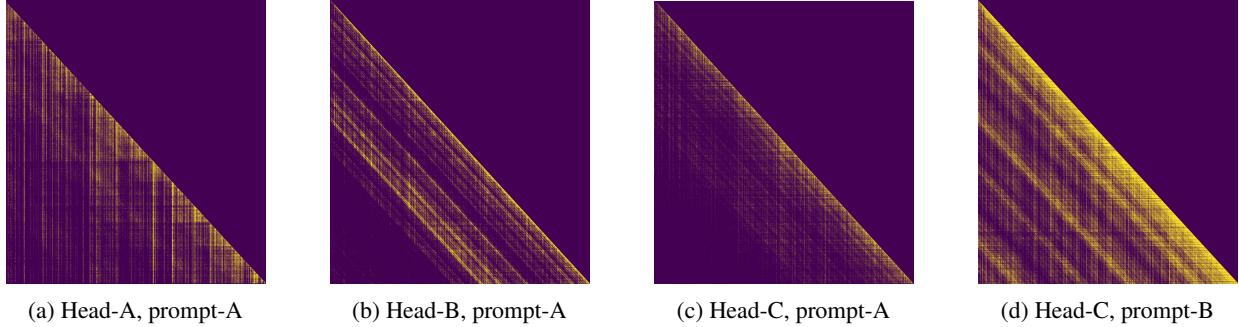| (a) Head-A, prompt-A | (b) Head-B, prompt-A | (c) Head-C, prompt-A | (d) Head-C, prompt-B |

*Figure 4.* The visualization of attention reveals diverse structured sparse patterns. Different heads with same prompt exhibit dynamic sparse indices and ratios, but the patterns can generally be categorized into (a) **column** , (b) **slash** , or (c) **composed** pattern with column and slash. These structured patterns generally extend across the entire head. However, in some heads, such as (c), a prominent column structure in the upper half gradually fades in the lower half. Additionally, as shown in (d), the same head exhibits significant pattern differences under different prompts, highlighting the content-aware nature.

This adaptive sparsity ratio demonstrates that applying a fixed sparsification budget across all attention heads and input contents, as employed by MInference (Jiang et al., 2024) and DuoAttention (Xiao et al., 2024), to search for sparse pattern is suboptimal. To maximize the efficiency of sparse attention while maintaining accuracy, the sparsity ratio need to be dynamically determined at runtime for each individual attention head and input prompt, allowing the model to adapt to the inherent variations in sparsity ratios.

### 3.2 C#2: Dynamic Sparse Pattern

Our observations also reveal that the attention pattern varies across different attention heads, input contents, and model architectures. Figure 4 visualizes distinct sparse patterns from different heads. Our analysis identifies two significant sparse patterns that substantially contribute to the attention score. The ***column pattern*** embodies crucial global contextual information (Figure 4(a)), with the attention sink (Xiao et al., 2023b) as a typical example. On the other hand, the ***slash pattern*** maintains connections between contexts of regular intervals (Figure 4(b)), such as the local window pattern which captures recent context information. These two patterns can be combined to cover diverse sparse patterns. For example, Figure 4(c-d) display a composition of column and slash patterns on a single head. Moreover, the concrete sparse pattern varies across different input contents.

These dynamic patterns present a significant challenge to determine the exact sparse pattern during inference. While MInference (Jiang et al., 2024) acknowledges similar column and slash patterns, its approach of classifying attention heads into fixed categories (A-shape, Vertical-Slash, or Block-Sparse) and determining optimal patterns offline fails to capture the content-dependent variations in sparse patterns. This limitation highlights the need for a more flexible, runtime-adaptive approach to pattern selection.
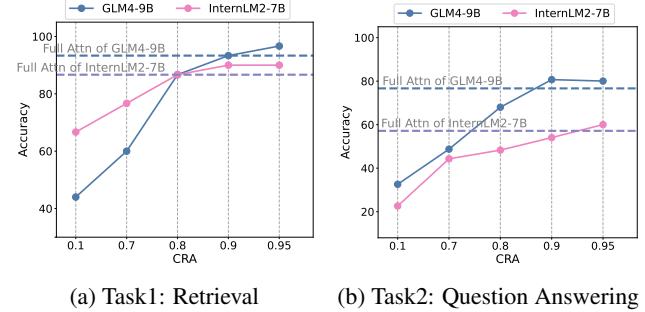


| (a) Task1: Retrieval | (b) Task2: Question Answering |

*Figure 5.* The curves of model accuracy and attention recall with changing **CRA** thresholds $\alpha$ for different tasks.

### 3.3 Insight: Cumulative Residual Attention

To effectively balance the trade-off between efficiency and accuracy when determining runtime sparsity ratios and patterns, we need a reliable metric to guide these decisions. Our finding reveals that the ***Cumulative Residual Attention*** (**CRA**), defined as the minimum sum of remaining attention probabilities per query after sparsification, serves as a robust indicator of model accuracy. As demonstrated in Figure 5, there exists a consistent positive correlation between the **CRA** threshold and model accuracy across different LLMs and tasks. This relationship provides a principled way to navigate the efficiency-accuracy trade-off: while lower **CRA** thresholds enable greater computational speedups, they should be carefully balanced against potential accuracy degradation.

Leveraging this insight, we can dynamically identify a minimal set of attention indices that satisfy the desired **CRA** threshold, thereby optimizing computational efficiency while preserving accuracy. This dynamic selection approach naturally accommodates varying sparsity ratios and sparse patterns across different attention heads, input contents, and model architectures.
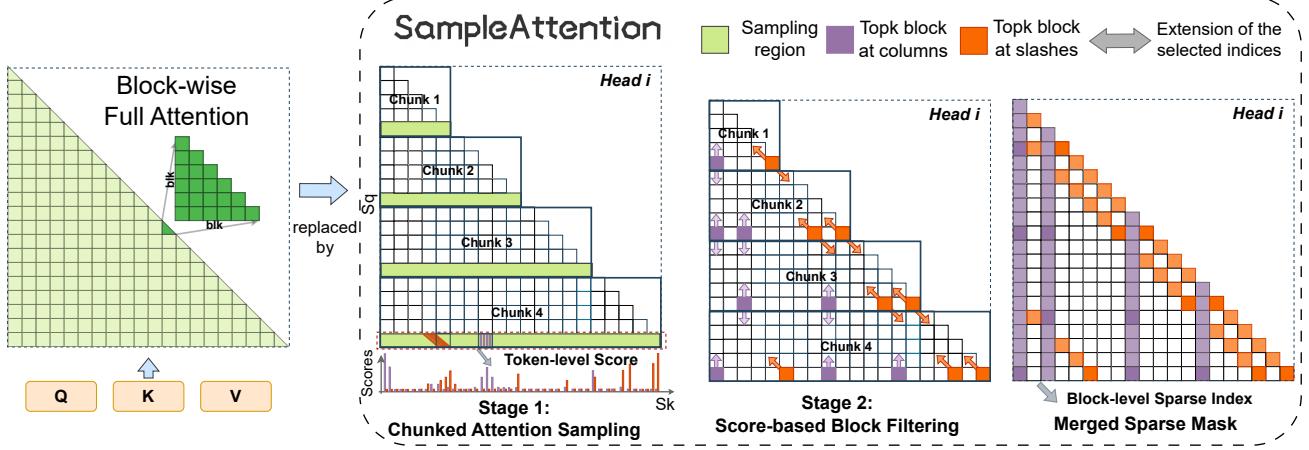
*Figure 6.* SampleAttention replaces the original full attention with a two-stage implementation. In the first stage, token-level attention scores are computed by performing chunked sampling across multiple query blocks and accumulating the scores along the column and slash direction. In the second stage, we determine the minimum quota of blocks required for each sampling region based on the block-reduced scores and thresholds $(\alpha_c, \alpha_s)$. Then, we perform top-k operation on each head to filter out the required block indices $I_c$ and $I_s$. These indices will be extended along the column and slash patterns and merged into $\hat{\mathbf{M}}$ to enable sparse computation in attention.
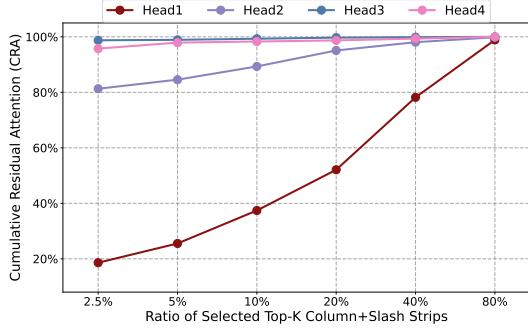


*Figure 7.* Relationship between the ratio of selected top-k columns-slashes and **CRA**. The high similarity of the numerical distribution enables a small amount of critical stripes-slashes to cover the majority values of the full attention score matrix.

The precise dynamic selection for a **CRA** threshold requires computing the full attention score, which is computationally expensive. Inspired by the observed significant column and slash pattern, we find that selecting an appropriate number of key column and slash strips can accurately approximate the **CRA** (Figure 7). The number of selected strips can be determined at runtime to vary the sparsity ratio across different heads and contents. Additionally, the combination of column and slash patterns provides sufficient flexibility to capture diverse attention distributions encountered.

## 4 SAMPLEATTENTION

In this section, we demonstrate how SampleAttention efficiently determines dynamic sparsity ratios and structured patterns at runtime, ensuring maximum efficiency through

sparse acceleration while maintaining nearly lossless accuracy. In Section 4.1, we provide an overview of the key considerations and the two-stage implementation of SampleAttention. Subsequently, in Section 4.2 and Section 4.3, we delve into the details of the two critical stages of SampleAttention: query-guided chunked sampling and score-based key-value filtering. Additionally, in Section 4.4, we explain how the introduced hyperparameters affect the tradeoff between accuracy and performance, and outline the method for tuning them. Finally, we detail the hardware-efficient implementation in Section 4.5.

### 4.1 Overview

Given a desired **CRA** threshold, SampleAttention dynamically selects a set of important column and slash strips to capture adaptive sparsity ratio and dynamic sparse pattern. However, the selection is non-trivial since SampleAttention requires fast and accurate estimation of attention scores. To address this, SampleAttention proposes a two-stage sampling algorithm.

The first stage is *Query-Guided Chunked Sampling*, where SampleAttention estimates the full attention score by computing the attention scores for a few queries, inspired by the observed column pattern. While prior work like MInference (Jiang et al., 2024) uses only last $last_q$ queries for estimation, we found this approach insufficient for capturing complex hybrid column and slash patterns. For instance, as shown in Figure 4(c), the column pattern is significant in the upper half region, but gradually fades in the lower half region. To capture such varying patterns effectively, SampleAttention partitions the queries into $chunk_n$ equal segments and perform separate attention sampling for each

chunk, enabling more accurate pattern detection across the entire attention matrix.

The second stage is *Score-Based Key-Value Filtering*, where SampleAttention filters key columns and slashes that meet a specified **CRA** threshold $\alpha$ based on the sampled attention scores. A naive approach of jointly selecting column and slash strips would require evaluating $n_{column} \times n_{slash}$ combinations, making the selection process computationally expensive. To improve efficiency, we decompose the single threshold $\alpha$ into separate thresholds: $\alpha_c$ for columns and $\alpha_s$ for slashes. This decomposition reduces the computational complexity from $n_{column} \times n_{slash}$ to $n_{column} + n_{slash}$. The filtering stage independently selects column and slash strips before merging them into the final pattern.

SampleAttention introduces several tunable hyperparameters to control the trade-off between efficiency and accuracy. To facilitate rapid adoption of SampleAttention, we provide a method for automatic hyperparameter tuning on a compact validation dataset, which maximizes computational efficiency while maintaining accuracy.

---

**Algorithm 1** Two-stage Implementation of SampleAttention

**Input:** $Q, K, V, \alpha_c, \alpha_s \in [0, 1], chunk_n$

*# Stage1: Query-Guided Chunked Attention Sampling*
$itv \leftarrow \frac{S_q}{chunk_n}$, $blk \leftarrow 128$
$Q_{slice} \leftarrow [Q_{[i*itv-blk:i*itv]}$ for $i$ in range$(1, chunk_n+1)]$
$\hat{A} \leftarrow \text{softmax}\left(Q_{slice}K^\top/\sqrt{d} + m_{\text{casual}}\right)$
$\hat{A}_c, \hat{A}_s \leftarrow \text{block\_reduction}\left(\hat{A}, blk\right)$

*# Stage2: Score-Based Key-Value Block Filtering*
$k_c \leftarrow \text{find\_k}\left(\text{cumsum}(\text{sort}(\hat{A}_c)), \alpha_c\right)$
$I_c \leftarrow \text{arg\_topk}\left(\hat{A}_c, k_c\right)$
$k_s \leftarrow \text{find\_k}\left(\text{cumsum}(\text{sort}(\hat{A}_s)), \alpha_s\right)$
$I_s \leftarrow \text{arg\_topk}\left(\hat{A}_s, k_s\right)$

*# Extend and Merge Block-sparse Mask across Each Head*
$\hat{M} \leftarrow \text{merge\_index}(I_c, I_s, itv)$

*# Final Sparse FlashAttention with Block Index*
$O \leftarrow \text{sparse\_flash\_attn}\left(Q, K, V, \hat{M}\right)$
return $O$

---

### 4.2 Stage1: Query-Guided Chunked Sampling

As discussed in Section 4.1, we need to divide the queries into chunks to ensure that sampling can more accurately determine the sparse structure within each head. Compared to random sampling or bottom sampling, this equidistant sampling technique is low-overhead and more stable. Our experiments show that this straightforward approach is effective: sampling a small number of query blocks can accurately approximate the actual **CRA** (Further details can be found in Appendix A.3). It should be noted that existing methods (Jiang et al., 2024; Li et al., 2024) commonly use bottom sampling, which fixes $chunk_n$ at 1. This means that only query blocks at the bottom of the score matrix are sampled. Due to the overly concentrated sample locations, this approach may lead to biases in index selection.

Subsequently, we accurately compute the attention scores for the query blocks at the bottom of each chunk. These token-level scores are then reduced at the block size ($blk$) granularity in both column and slash directions. The resulting block-level scores help us dynamically determine the required sparsity ratios and indices in the second stage.

### 4.3 Stage 2: Score-Based Key-Value Filtering

In the second stage, our approach performs a refined filtering of key-value indices based on the sample attention scores. A significant challenge remains in efficiently and dynamically selecting the minimal set of key-value indices of interest under different patterns, denoted as $I_c$ and $I_s$, that align with the input prompt and satisfy the CRA threshold $\alpha$. Determining a fixed quota of critical indices does not achieve the optimal balance between accuracy and performance, especially given the extensive sequence lengths. To address this inefficiency, we employ an adaptive approach that filters block-level indices based on blocked attention scores in both the column and slash directions, comparing them with their respective thresholds.

In detail, as shown in Algorithm 1, for each sampled block within a chunk, SampleAttention first accumulates the attention scores along the column and slash directions and then reduces them at the block granularity. This accumulation serves as a statistical approximation of the overall attention scores. Using these block-level scores in both directions, SampleAttention can effectively select the minimal number of key-value blocks $k_c$ and $k_s$ that meet the CRA thresholds $\alpha_c$ and $\alpha_s$ for each attention head, respectively. Finally, based on the derived minimal quota, SampleAttention performs a top-k operation in each direction to filter out the essential block indices $I_c$ and $I_s$. It is noteworthy that sampled chunks from different positions may filter out different sets of column and slash indices. SampleAttention extends these indices obtained from each sampled block according to the patterns so that they cover the entire attention matrix. This ensures that the final merged mask $\hat{M}$ achieves an almost lossless sparse approximation. This approach also enables the identification of critical attention sinks and local window masks, thereby maintaining stable accuracy.

### 4.4 Hyperparameter Tuning

Due to the introduction of three hyperparameters in SampleAttention, as shown in Table 1, it is crucial to analyze their impact on the accuracy and performance. Here, we

briefly discuss the influence of these parameters and outline the approach to tuning them. Detailed results of modifying these hyperparameters are studied in Section 5.3.

*Thresholds for Columns and Slashes.* The most critical hyperparameters in SampleAttention are the CRA thresholds $\alpha_c$ and $\alpha_s$. Generally, larger threshold values can reduce the speedup but enhance the model's accuracy. Therefore, it is necessary to predetermine cost-effective threshold values for columns and slashes offline using a compact dataset. We can use the accuracy and latency metrics of FlashAttention2 (Dao, 2023) as reference standards for this process. Moreover, since long-context tasks span a wide range, further segmenting the context by length and tuning each segment individually can more effectively leverage sparsity at different lengths, thereby achieving more cost-effective threshold values.

*Sampling Positions and Ratios.* Additionally, the number of sampling chunks significantly affects SampleAttention's performance by influencing the sampling positions and the ratio of selected indices. For instance, too few sampling samples may fail to capture the full sparse structures in the head, thereby reducing accuracy. Conversely, excessive sampling can increase overhead and introduce redundant computations in attention. Therefore, during the tuning, we introduce multiple values for $chunk_n$ to expand the search space and identify more efficient configurations.

*Table 1.* The meaning of hyperparameters and they will be tuned offline for different length ranges.

| Hyperparameter | Description |
| --- | --- |
| $\alpha_c$ | The desired **CRA** threshold for columns |
| $\alpha_s$ | The desired **CRA** threshold for slashes |
| $chunk_n$ | The number of sampling chunks |

### 4.5  Hardware-efficient Implementation

To achieve substantial speedup in wall-clock time, SampleAttention is implemented with IO-awareness to maximize hardware-efficiency. First, the query-guided key-value filtering involves a series of small operators (`bmm`, `mask_fill`, `softmax`, `reduction`) that read and write large intermediate results. SampleAttention significantly reduces IO overhead by fusing these operators. Second, SampleAttention implements an efficient adaptive structured sparse attention kernel by modifying FlashAttention2 (Dao, 2023). These hardware-aware optimizations enhance speed performance significantly.

## 5  EXPERIMENTS

### 5.1  Setup

**Backbones.** We evaluate our method on three widely used open-source LLM variants: **ChatGLM4-9B** with a 1M context window based on GLM (Du et al., 2021; GLM et al., 2024); ; **YI-9B**, featuring a 200K context window (Young et al., 2024); and **InternLM2-7B**, also with a 200K context window (Cai et al., 2024). All utilized models are decoder-only transformers (Radford et al., 2018), and are pre-trained via causal language modeling. They encompass similar architectural components, such as rotary positional encoding (Su et al., 2024), and grouped-query attention (Ainslie et al., 2023). Simultaneously, there are notable differences, e.g., the former augments the context window capacity via continued training with an extended sequence length, whereas the latter achieves length extrapolation through rope scaling. It is important to note that we specifically replace the full-attention implementation during the prompt prefill stage with SampleAttention and baseline methods, while preserving an uncompressed key-value (KV) cache and dense attention computation in the decoding phase.

**Tasks.** We evaluate SampleAttention and other methods' understanding capabilities in long-context scenarios on three distinct tasks: RULER (Hsieh et al., 2024), LongBench (Bai et al., 2023) and Infinite Bench (Zhang et al., 2024b). **RULER** provides a comprehensive evaluation of long-context language models through flexible configurations for sequence lengths. Unlike the traditional needle-in-a-haystack (Kamradt, 2023) test, RULER extends it by incorporating diverse types and quantities of "needles" and introduces new tasks such as multi-hop tracing and aggregation, which evaluate more complex behaviors beyond simple retrieval. RULER encompasses 13 tasks, making it an excellent tool for testing long-context understanding. **LongBench**, a multi-task benchmark, comprises single and multi-document QA, summarization, few-shot learning, synthetic tasks, and code completion. It offers over 4,750 test cases with task lengths from 4K-35K. **InfiniteBench**, a benchmark specifically designed to evaluate language models' ability in handling, understanding, and reasoning in contexts exceeding an average length of 200K. It comprises 10 unique tasks, each crafted to assess different aspects of language processing and comprehension in extended contexts.

**Baselines and settings.** We conducted all experiments on a single NVIDIA-A100 GPU (80GB) to evaluate accuracy and performance of attention operation during the prefill stage. We consider the full attention (as the gold baseline), Minference (Jiang et al., 2024), BigBird (Zaheer et al., 2020), Streaming-LLM (Xiao et al., 2023b), HyperAttention (Han et al., 2023) and Hash-Sparse (Pagliardini et al., 2023) as
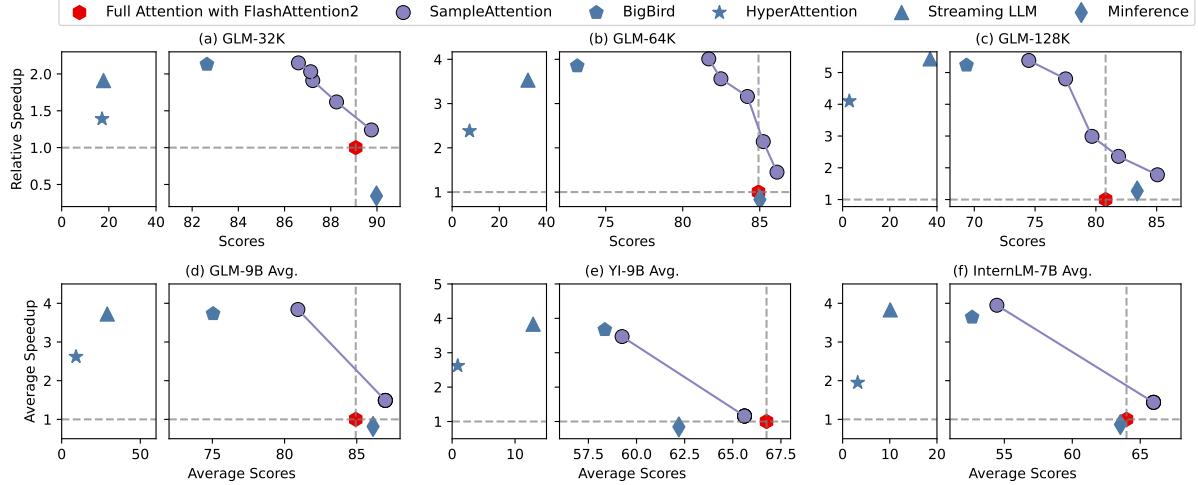
*Figure 8.* The trade-off between accuracy (measured via RULER benchmark scores) and speedup (relative to FlashAttention2) is analyzed across various sparse attention methods. Figures (a)–(c) illustrate this trade-off under the GLM4-9B architecture for sequence lengths of 32K, 64K, and 128K, respectively. Meanwhile, Figures (d)–(f) present averaged results across different models for each of the three sequence lengths, focusing exclusively on configurations demonstrating optimal performance in both acceleration and precision. Notably, our proposed method consistently outperforms MInference in both accuracy and speedup across all evaluated scenarios, showcasing superior efficiency-accuracy trade-offs.

baselines to compare model accuracy across different tasks. Minference requires pre-profiling to determine the optimal pattern. And both BigBird and StreamingLLM were assigned a window size ratio of $8\%$. BigBird retains a global ratio of $8\%$. StreamingLLM sets its initail attention sink at 4 tokens. HyperAttention set both bucket size and the number of sampled columns to 256. For SampleAttention, we generate small-scale tasks using the RULER benchmark at specific lengths (e.g., 16K, 32K, 64K, 128K) to tune hyperparameters and subsequently apply these parameters to tasks across different length ranges. This approach enables SampleAttention to achieve an optimal trade-off balance within each length range.

## 5.2 Trade-off between Accuracy and Efficiency

**Main results.** Figure 8 compares the trade-offs between accuracy and the speedup relative to FlashAttention2 among different sparse methods and SampleAttention under various hyperparametres on the RULER benchmark. Table 2, on the other hand, compares the accuracy of different sparse methods on LongBench and InfiniteBench tasks. The results show that:

- The performance in accuracy of SampleAttention is consistently robust across all benchmarks (including subdomains), various models, and diverse sequence lengths. When compared to full attention, which serves as the gold standard, SampleAttention consistently achieves scores above 99% of full attention,

demonstrating near-lossless efficiency. Furthermore, our approach establishes a new Pareto frontier in the accuracy-efficiency trade-off, surpassing existing methods in both dimensions.

- While Minference achieves accuracy comparable to FullAttention at lengths of 32K and 64K, it fails to provide any speedup benefits. In contrast, sample attention maintains nearly lossless performance and achieves speedup ranging from $1.24\times$ to $2.36\times$ compared to FlashAttention2 across different lengths.

- BigBird exhibits varying degrees of performance degradation across different models and lengths. Nonetheless, on average, BigBird still achieves approximately 86% of the scores achieved by full attention and provides a relatively stable speedup due to the nature of its static pattern.

- StreamingLLM and HyperAttention result in performance degradation across all tasks, demonstrating that these techniques fail to capture critical KV elements in long sequences at the prefill stage.

## 5.3 Ablation Study and Tuning for Hyperparameter

As introduced in Section 4.4, the sparse properties of different models exhibit significant variability across sequence lengths. Therefore, we implement an automated, offline hyperparameter tuning method for the specified models on the small-scale datasets. This approach discretizes sequence lengths into distinct intervals and performs multi-task tuning

*Table 2.* Accuracy comparison across various sparse methods on LongBench and InfiniteBench under GLM4-9B model. The hyperparameters applied are the optimal accuracy configurations, tuned for different sequence length ranges. The best results are highlighted in **Bold** while the second best results are marked with an <u>Underline</u>.

| Benchmark | Baseline | Task Type | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Single-Doc QA | Multi-Doc QA | Summari-zation | Few-shot Learning | Synthetic Tasks | Code Completion | Total Score | | | |
| LongBench | Full Attention | <u>213.12</u> | <u>174.35</u> | <u>109.69</u> | 273.87 | 231.49 | 121.52 | 1124.04 | | | |
| | Ours | **214.53** | **174.42** | 108.92 | **278.33** | **234.55** | **125.18** | **1135.93** | | | |
| | Minference | 212.14 | 173.37 | **110.02** | <u>274.45</u> | <u>231.87</u> | <u>124.37</u> | <u>1126.22</u> | | | |
| | BigBrid | 207.57 | 146.45 | 95.64 | 272.17 | 161.60 | 117.38 | 1000.81 | | | |
| | StreamingLLM | 142.79 | 129.36 | 89.71 | 168.13 | 19.70 | 98.43 | 648.12 | | | |
| | HyperAttention | 125.74 | 119.08 | 88.05 | 206.35 | 32.69 | 86.35 | 658.26 | | | |
| | | En.Sum | En.QA | En.MC | En.Dia | Zh.QA | Code. Debug | Math. Find | Retr. PassKey | Retr. Number | Retr. KV |
| InfiniteBench | Full Attention | **28.30** | <u>12.17</u> | 58.95 | **34.00** | 13.22 | 30.71 | <u>37.71</u> | **100** | **100** | <u>44.0</u> |
| | Ours | **28.30** | **16.52** | **61.57** | <u>31.50</u> | <u>14.28</u> | 31.40 | 37.14 | **100** | **100** | **49.6** |
| | Minference | 28.00 | 11.39 | <u>60.26</u> | 28.70 | **14.81** | **31.70** | **39.43** | **100** | **100** | 43.0 |

within each segment. Such a strategy ensures an optimal balance between accuracy and efficiency over the entire range of sequence lengths. Figure 9 demonstrates how varying the CRA thresholds for column and slash patterns influences both accuracy and the sparsity ratio under different conditions. Table 3, on the other hand, investigates how different numbers of sampling chunks impact accuracy and computational speedup.

**CRA threshold $\alpha_c$ and $\alpha_s$:** Tuning $\alpha_c$ and $\alpha_s$ within a given model is crucial for finding a cost-effective configuration. Generally, as shown in Figures 9 (a) and (c), increasing either $\alpha_c$ or $\alpha_s$ alone can improve accuracy, but at the expense of increased computational load. However, there are differences between them. The YI model in Figure (c) at 128K shows significantly greater sensitivity to changes in the slash threshold $\alpha_s$, whereas Figure (a) demonstrates a more balanced response to changes in both thresholds. Additionally, Figure (b) shows that smaller thresholds can still deliver sufficiently good scores, enabling higher speedup while maintaining accuracy.

**Number of Sampling Chunks:** We further evaluate the impact of different numbers (1, 2, 4, 6) of sampling chunks on accuracy, given a fixed threshold. Table 3 demonstrates that an appropriate $chunk_n$ size can yield more cost-effective results. For instance, increasing the number from 1 to 2 helps SampleAttention enhance accuracy without significantly impacting the speedup ratio. However, an excessively large number of chunks may not improve accuracy and can decrease the speedup. Thus, selecting an appropriate $chunk_n$ size is essential.

**Cross-Task Robustness** The hyperparameters of SampleAttention are inherently model-specific due to architectural design and intrinsic sparsity patterns. To evaluate cross-task robustness, we tested shared hyperparameter configurations

across three distinct benchmarks under the same model. First, we performed tuning on GLM4 and InternLM2 using subsets of the RULER benchmark at various sequence lengths. We selected optimal configurations under two criteria: accuracy-optimized settings with negligible performance loss compared to full-attention baselines, and configurations maximizing acceleration gains. Detailed hyperparameters are provided in Appendix 6. Experimental results in Figure 10 demonstrate that hyperparameters tuned on a subset of tasks generalize effectively across diverse benchmarks. For instance, GLM4's accuracy-optimized hyperparameters maintained near-lossless performance on different tasks from LongBench and InfiniteBench. This indicates robust cross-domain adaptability without significant performance degradation. Additionally, consistent speedup gains were observed for the same model under identical sequence lengths across tasks.

*Table 3.* The impact of changing $chunk_n$ on *scores/speedup* in different cases. The scores above are based on RULER, while the speedup below are relative to FlashAttention2. The best score results are highlighted in **bold**, while the best speedup results are marked with <u>**underline**</u>.

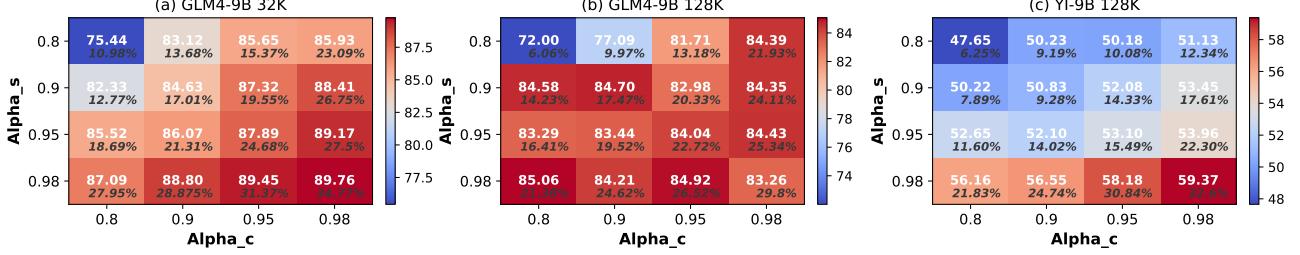| Model | $(\alpha_c, \alpha_s)$ | $chunk_n$ | | | |
|---|---|---|---|---|---|
| | | 1 | 2 | 4 | 6 |
| GLM (128K) | (0.90,0.90) | 82.89/ <u>**1.92**</u> | **84.70**/ 1.89 | 84.02/ 1.70 | 83.62/ 1.53 |
| | (0.95,0.95) | 84.17/ 1.64 | 84.04/ 1.60 | 83.14/ 1.46 | 83.73/ 1.33 |
| Yi (128K) | (0.95,0.95) | 52.81/ <u>**2.17**</u> | 54.54/ 2.12 | 53.10/ 1.89 | 54.58/ 1.71 |
| | (0.98,0.98) | 56.24/ 1.29 | 58.58/ 1.25 | **59.37**/ 1.21 | 59.36/ 1.13 |

*Figure 9.* The heatmaps under different cases illustrate the impact of choosing different values of $\alpha_c$ and $\alpha_s$ on the accuracy of SampleAttention (the white text represents the scores under the RULER benchmark, while the black italic text denotes the actual ratio of calculated blocks). The $chunk_n$ values for the GLM and YI models are set to 2 and 4, respectively.
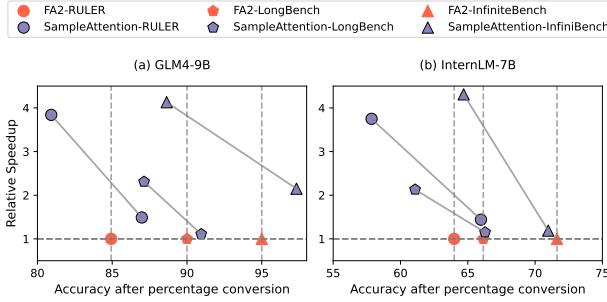


*Figure 10.* Results from offline tuning and evaluation of (a) GLM4-9B and (b) InternLM2-7B across RULER, LongBench, and InfiniteBench benchmarks. Different tasks share the same hyperparameters from offline tuning when sequence lengths fall within the same range.



*Figure 11.* (a) The percentage of time spent on sampling and sparse computation in SampleAttention. (b) Comparison of the TTFT metric using FlashAttention2 as the baseline.

## 5.4 Acceleration Speedup Benchmarking

We conducted micro-benchmarks on a single A100 to evaluate the time breakdown and TTFT metrics. The baseline selected is FlashAttention2. All tests were conducted using the configuration from ChatGLM4-9B: 32 heads, and $d = 128$, with synthetic data from the RULER benchmark as input. We standardized the batch size of the input data to 1 to support longer sequence lengths. The hyperparameters are tuned to attain the best possible speedup without accuracy loss compared to FlashAttention2, ensuring the optimal acceleration-accuracy trade-off. For example, the hyperparameters for 32K are set to $\alpha_c = 0.95$, $\alpha_s = 0.95$, and $chunk_n = 1$. For lengths above 128K, the parameters tuned for 128K are reused.

**Sampling overhead.** Figure 11(a) presents the time breakdown for a full 40-layer model using synthetic data spanning sequence lengths from 32K to 128K. The results reveal that as sequence lengths increase, the relative contribution of sampling overhead diminishes. This trend underscores the potential of SampleAttention to deliver substantial acceleration advantages for longer sequences. However, in short-sequence scenarios, the performance gains are less
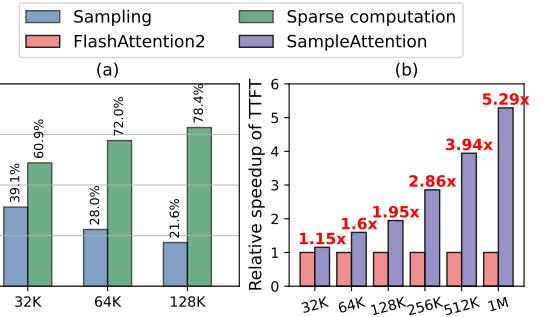
pronounced due to the computational overhead associated with dynamic sampling and index construction. For instance, in the GLM4-9B model at an 8K sequence length, the latency of SampleAttention remains nearly identical to FlashAttention2 without compromising accuracy.

**Scaling the sequence length to 1M.** We also conducted GPU performance evaluations scalable to a sequence length of 1 million. To avoid memory issues, for sequence lengths greater than 256K, we divide the input sequence into multiple chunks and partition them by head during attention. We also further reduced a large number of intermediate results through operator fusion. Figure 11(b) shows that when the sequence length scales to 1M, the TTFT metric can be significantly reduced by $5.29\times$.

## 6 RELATED WORK

**Approximate Attention.** Plenty of works have been proposed to approximate quadratic attention with lower complexity(Ainslie et al., 2020; Beltagy et al., 2020; Zaheer et al., 2020; Kitaev et al., 2020; Ding et al., 2023; Child et al., 2019; Pagliardini et al., 2023; Roy et al., 2021; Han et al., 2023; Wang et al., 2020; Choromanski et al., 2020; Katharopoulos et al., 2020; Chen et al., 2021b;a; Zhu

et al., 2023; Ribar et al., 2023; Roy et al., 2021; Liu et al., 2023b). For example, BigBird (Zaheer et al., 2020) combines window-, global- and random-attention to capture long range dependency. Reformer (Kitaev et al., 2020) reduces computional cost via locality-sensitive hashing. LongNet (Ding et al., 2023) replaces full attention with dilated attention. Linformer (Wang et al., 2020) employs low-rank matrix to approximate attention. HyperAttention (Han et al., 2023) utilizes locality sensitive hashing to identify important entries on attention map. However, these approaches uses either static or coarse-grained sparse pattern, and often overlook the head-specific sparsity pattern. They cannot be losslessly applied in pretrained LLMs without additional finetuning or training. Additionally, although the recent work Minference (Jiang et al., 2024) addresses the varying sparsity among attention heads by employing three distinct sparse attention patterns for long sequences, it relies on predefined fixed ratios for these sparse patterns and fails to account for the dynamic nature of the actual model and its prompts.

**KV Cache Compression.** Long sequence comes with substantial KV cache memory consumption. StreamingLLM (Xiao et al., 2023b) keeps attention sinks and several recent tokens for infinite length generation. H2O (Zhang et al., 2024c) dynamically retains a balance of recent and heavy hitter tokens according to attention score during decoding. FastGen (Ge et al., 2023) adaptively construct KV cache according to observed head-specific policies. SnapKV (Li et al., 2024) strategically compresses the KV cache by selecting clustered critical positions for each attention head, leading to improvements in memory and latency efficiency during decoding. CHAI (Agarwal et al., 2024) exemplifies head pruning methods that target reducing the KV cache overhead at the attention head level, thereby accelerating decoding. Recent efforts also quantize KV cache to lower precision to reduce memory consumption (Duanmu et al., 2024; Xiao et al., 2023a; Zhao et al., 2023). These works target on reducing the memory consumption of KV cache, while SampleAttention focuses on mitigating the long context computation overhead. SampleAttention can be combined with these approaches to further reduce memory consumption of KV cache.

## 7 CONCLUSION

In this paper, we identify the challenge of effectively exploiting the inherent high attention sparsity to accelerate prefill attention, due to the highly dynamic structured patterns and optimal sparsity ratios exhibited by the attention mechanism over long contexts. To address this, we propose SampleAttention, which utilizes **CRA** as a robust inficator of model accuracy and adaptively determines the sparsity ratio and pattern at runtime. Through an innovative two-stage

query-guided filtering approach and hyperparameter tuning, it dynamically selects the minimal set of critical key-values, maximizing efficiency while maintaining accuracy. Experimental results demonstrate that SampleAttention consistently maintains robust accuracy across various benchmarks, models, and sequence lengths, and significantly reduces the TTFT metric.

## REFERENCES

Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

Agarwal, S., Acun, B., Hosmer, B., Elhoushi, M., Lee, Y., Venkataraman, S., Papailiopoulos, D., and Wu, C.-J. Chai: Clustered head attention for efficient llm inference. arXiv preprint arXiv:2403.08058, 2024.

Agrawal, A., Panwar, A., Mohan, J., Kwatra, N., Gulavani, B. S., and Ramjee, R. Sarathi: Efficient llm inference by piggybacking decodes with chunked prefills. arXiv preprint arXiv:2308.16369, 2023.

Ainslie, J., Ontanon, S., Alberti, C., Cvicek, V., Fisher, Z., Pham, P., Ravula, A., Sanghai, S., Wang, Q., and Yang, L. ETC: Encoding long and structured inputs in transformers. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 268–284, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. emnlp-main.19. URL https://aclanthology.org/2020.emnlp-main.19.

Ainslie, J., Lee-Thorp, J., de Jong, M., Zemlyanskiy, Y., Lebrón, F., and Sanghai, S. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. arXiv preprint arXiv:2305.13245, 2023.

Anthropic. Claude. 2023.

Bai, Y., Lv, X., Zhang, J., Lyu, H., Tang, J., Huang, Z., Du, Z., Liu, X., Zeng, A., Hou, L., et al. Longbench: A bilingual, multitask benchmark for long context understanding. arXiv preprint arXiv:2308.14508, 2023.

Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. arXiv preprint arXiv:2004.05150, 2020.

Bommasani, R., Hudson, D. A., Adeli, E., Altman, R. B., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N. S., Chen, A. S., Creel,

K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N. D., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M. S., Krishna, R., Kuditipudi, R., and et al. On the opportunities and risks of foundation models. CoRR, abs/2108.07258, 2021. URL https://arxiv.org/abs/2108.07258.

Borgeaud, S., Mensch, A., Hoffmann, J., Cai, T., Rutherford, E., Millican, K., Van Den Driessche, G. B., Lespiau, J.-B., Damoc, B., Clark, A., et al. Improving language models by retrieving from trillions of tokens. In International conference on machine learning, pp. 2206–2240. PMLR, 2022.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html.

Bulatov, A., Kuratov, Y., and Burtsev, M. Recurrent memory transformer. Advances in Neural Information Processing Systems, 35:11079–11091, 2022.

Cai, Z., Cao, M., Chen, H., Chen, K., Chen, K., Chen, X., Chen, X., Chen, Z., Chen, Z., Chu, P., et al. Internlm2 technical report. arXiv preprint arXiv:2403.17297, 2024.

Chen, B., Dao, T., Liang, K., Yang, J., Song, Z., Rudra, A., and Re, C. Pixelated butterfly: Simple and efficient sparse training for neural network models. In International Conference on Learning Representations, 2021a.

Chen, B., Dao, T., Winsor, E., Song, Z., Rudra, A., and Ré, C. Scatterbrain: Unifying sparse and low-rank attention. Advances in Neural Information Processing Systems, 34:17413–17426, 2021b.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374, 2021c.

Chen, S., Wong, S., Chen, L., and Tian, Y. Extending context window of large language models via positional interpolation. arXiv preprint arXiv:2306.15595, 2023a.

Chen, Y., Qian, S., Tang, H., Lai, X., Liu, Z., Han, S., and Jia, J. Longlora: Efficient fine-tuning of long-context large language models. In The Twelfth International Conference on Learning Representations, 2023b.

Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.

Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers. arXiv preprint arXiv:1904.10509, 2019.

Choromanski, K. M., Likhosherstov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J. Q., Mohiuddin, A., Kaiser, L., et al. Rethinking attention with performers. In International Conference on Learning Representations, 2020.

Dao, T. Flashattention-2: Faster attention with better parallelism and work partitioning. arXiv preprint arXiv:2307.08691, 2023.

Dao, T., Fu, D., Ermon, S., Rudra, A., and Ré, C. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in Neural Information Processing Systems, 35:16344–16359, 2022.

Ding, J., Ma, S., Dong, L., Zhang, X., Huang, S., Wang, W., Zheng, N., and Wei, F. Longnet: Scaling transformers to 1,000,000,000 tokens. arXiv preprint arXiv:2307.02486, 2023.

Du, Z., Qian, Y., Liu, X., Ding, M., Qiu, J., Yang, Z., and Tang, J. Glm: General language model pretraining with autoregressive blank infilling. arXiv preprint arXiv:2103.10360, 2021.

Duanmu, H., Yuan, Z., Li, X., Duan, J., Zhang, X., and Lin, D. Skvq: Sliding-window key and value cache quantization for large language models. arXiv preprint arXiv:2405.06219, 2024.

Ge, S., Zhang, Y., Liu, L., Zhang, M., Han, J., and Gao, J. Model tells you what to discard: Adaptive kv cache compression for llms. arXiv preprint arXiv:2310.01801, 2023.

GLM, T., Zeng, A., Xu, B., Wang, B., Zhang, C., Yin, D., Rojas, D., Feng, G., Zhao, H., Lai, H., et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. arXiv preprint arXiv:2406.12793, 2024.

Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752, 2023.

Han, I., Jayaram, R., Karbasi, A., Mirrokni, V., Woodruff, D., and Zandieh, A. Hyperattention: Long-context attention in near-linear time. In The Twelfth International Conference on Learning Representations, 2023.

Hooper, C., Kim, S., Mohammadzadeh, H., Mahoney, M. W., Shao, Y. S., Keutzer, K., and Gholami, A. Kvquant: Towards 10 million context length llm inference with kv cache quantization. arXiv preprint arXiv:2401.18079, 2024.

Hsieh, C.-P., Sun, S., Kriman, S., Acharya, S., Rekesh, D., Jia, F., and Ginsburg, B. Ruler: What's the real context size of your long-context language models? arXiv preprint arXiv:2404.06654, 2024.

Jiang, H., Li, Y., Zhang, C., Wu, Q., Luo, X., Ahn, S., Han, Z., Abdi, A. H., Li, D., Lin, C.-Y., et al. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. arXiv preprint arXiv:2407.02490, 2024.

Kamradt, G. Needle in a haystack–pressure testing llms, 2023.

Katharopoulos, A., Vyas, A., Pappas, N., and Fleuret, F. Transformers are rnns: Fast autoregressive transformers with linear attention. In International conference on machine learning, pp. 5156–5165. PMLR, 2020.

Kitaev, N., Kaiser, Ł., and Levskaya, A. Reformer: The efficient transformer. arXiv preprint arXiv:2001.04451, 2020.

Kressner, D. and Tobler, C. Low-rank tensor Krylov subspace methods for parametrized linear systems. SIAM Journal on Matrix Analysis and Applications, 32(4): 1288–1316, 2011.

Kuratov, Y., Bulatov, A., Anokhin, P., Sorokin, D., Sorokin, A., and Burtsev, M. In search of needles in a 10m haystack: Recurrent memory finds what llms miss. arXiv preprint arXiv:2402.10790, 2024.

Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In Proceedings of the 29th Symposium on Operating Systems Principles, pp. 611–626, 2023.

Li, D., Shao, R., Xie, A., Sheng, Y., Zheng, L., Gonzalez, J. E., Stoica, I., Ma, X., and Zhang, H. How long can open-source llms truly promise on context length?, June 2023a. URL https://lmsys.org/blog/2023-06-29-longchat.

Li, Y., Dong, B., Lin, C., and Guerin, F. Compressing context to enhance inference efficiency of large language models. arXiv preprint arXiv:2310.06201, 2023b.

Li, Y., Huang, Y., Yang, B., Venkitesh, B., Locatelli, A., Ye, H., Cai, T., Lewis, P., and Chen, D. Snapkv: Llm knows what you are looking for before generation. arXiv preprint arXiv:2404.14469, 2024.

Likhosherstov, V., Choromanski, K., and Weller, A. On the expressive power of self-attention matrices. arXiv preprint arXiv:2106.03764, 2021.

Liu, L., Qu, Z., Chen, Z., Tu, F., Ding, Y., and Xie, Y. Dynamic sparse attention for scalable transformer acceleration. IEEE Transactions on Computers, 71(12):3165–3178, 2022.

Liu, X., Yan, H., Zhang, S., An, C., Qiu, X., and Lin, D. Scaling laws of rope-based extrapolation. arXiv preprint arXiv:2310.05209, 2023a.

Liu, Z., Wang, J., Dao, T., Zhou, T., Yuan, B., Song, Z., Shrivastava, A., Zhang, C., Tian, Y., Re, C., et al. Deja vu: Contextual sparsity for efficient llms at inference time. In International Conference on Machine Learning, pp. 22137–22176. PMLR, 2023b.

Moonshot. Kimi chat. 2023.

Mu, J., Li, X., and Goodman, N. Learning to compress prompts with gist tokens. Advances in Neural Information Processing Systems, 36, 2024.

Pagliardini, M., Paliotta, D., Jaggi, M., and Fleuret, F. Faster causal attention over large sequences through sparse flash attention. arXiv preprint arXiv:2306.01160, 2023.

Patel, P., Choukse, E., Zhang, C., Shah, A., Goiri, Í., Maleki, S., and Bianchini, R. Splitwise: Efficient generative llm inference using phase splitting. In 2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA), pp. 118–132. IEEE, 2024.

Peng, B., Alcaide, E., Anthony, Q., Albalak, A., Arcadinho, S., Cao, H., Cheng, X., Chung, M., Grella, M., GV, K. K., et al. Rwkv: Reinventing rnns for the transformer era. arXiv preprint arXiv:2305.13048, 2023.

Qin, R., Li, Z., He, W., Zhang, M., Wu, Y., Zheng, W., and Xu, X. Mooncake: Kimi's kvcache-centric architecture for llm serving. arXiv preprint arXiv:2407.00079, 2024.

Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. Improving language understanding by generative pre-training. 2018.

Reddi, V. J., Cheng, C., Kanter, D., Mattson, P., Schmuelling, G., Wu, C.-J., Anderson, B., Breughe, M., Charlebois, M., Chou, W., et al. Mlperf inference benchmark. In 2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA), pp. 446–459. IEEE, 2020.

Ribar, L., Chelombiev, I., Hudlass-Galley, L., Blake, C., Luschi, C., and Orr, D. Sparq attention: Bandwidth-efficient llm inference. arXiv preprint arXiv:2312.04985, 2023.

Roy, A., Saffar, M., Vaswani, A., and Grangier, D. Efficient content-based sparse attention with routing transformers. Transactions of the Association for Computational Linguistics, 9:53–68, 2021.

Roziere, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X. E., Adi, Y., Liu, J., Remez, T., Rapin, J., et al. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950, 2023.

Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. arXiv preprint arXiv:1701.06538, 2017.

Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., and Liu, Y. Roformer: Enhanced transformer with rotary position embedding. Neurocomputing, 568:127063, 2024.

Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.

Tay, Y., Dehghani, M., Bahri, D., and Metzler, D. Efficient transformers: A survey. ACM Computing Surveys, 55 (6):1–28, 2022.

Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.

Tillet, P., Kung, H.-T., and Cox, D. Triton: an intermediate language and compiler for tiled neural network computations. In Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages, pp. 10–19, 2019.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. Advances in neural information processing systems, 30, 2017.

Wang, S., Li, B. Z., Khabsa, M., Fang, H., and Ma, H. Linformer: Self-attention with linear complexity. arXiv preprint arXiv:2006.04768, 2020.

Wu, Y., Rabe, M. N., Hutchins, D., and Szegedy, C. Memorizing transformers. arXiv preprint arXiv:2203.08913, 2022.

Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., and Han, S. Smoothquant: Accurate and efficient post-training quantization for large language models. In International Conference on Machine Learning, pp. 38087–38099. PMLR, 2023a.

Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks. arXiv preprint arXiv:2309.17453, 2023b.

Xiao, G., Tang, J., Zuo, J., Guo, J., Yang, S., Tang, H., Fu, Y., and Han, S. Duoattention: Efficient long-context llm inference with retrieval and streaming heads. arXiv preprint arXiv:2410.10819, 2024.

Xiong, W., Liu, J., Molybog, I., Zhang, H., Bhargava, P., Hou, R., Martin, L., Rungta, R., Sankararaman, K. A., Oguz, B., et al. Effective long-context scaling of foundation models. arXiv preprint arXiv:2309.16039, 2023.

Young, A., Chen, B., Li, C., Huang, C., Zhang, G., Zhang, G., Li, H., Zhu, J., Chen, J., Chang, J., et al. Yi: Open foundation models by 01. ai. arXiv preprint arXiv:2403.04652, 2024.

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, T., Ladhak, F., Durmus, E., Liang, P., McKeown, K., and Hashimoto, T. Benchmarking large language models for news summarization. Transactions of the Association for Computational Linguistics, 12, 2024a.

Zhang, X., Chen, Y., Hu, S., Xu, Z., Chen, J., Hao, M., Han, X., Thai, Z., Wang, S., Liu, Z., et al. Infinitebench: Extending long context evaluation beyond 100k tokens. In Proceedings of the 62nd Annual Meeting of

the Association for Computational Linguistics (Volume 1: Long Papers), pp. 15262–15277, 2024b.

Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Ré, C., Barrett, C., et al. H2o: Heavy-hitter oracle for efficient generative inference of large language models. Advances in Neural Information Processing Systems, 36, 2024c.

Zhao, Y., Lin, C.-Y., Zhu, K., Ye, Z., Chen, L., Zheng, S., Ceze, L., Krishnamurthy, A., Chen, T., and Kasikci, B. Atom: Low-bit quantization for efficient and accurate llm serving. arXiv preprint arXiv:2310.19102, 2023.

Zhong, Y., Liu, S., Chen, J., Hu, J., Zhu, Y., Liu, X., Jin, X., and Zhang, H. Distserve: Disaggregating prefill and decoding for goodput-optimized large language model serving. arXiv preprint arXiv:2401.09670, 2024.

Zhu, L., Wang, X., Ke, Z., Zhang, W., and Lau, R. W. Biformer: Vision transformer with bi-level routing attention. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10323–10333, 2023.

## A APPENDIX

### A.1 Visualization of Attention Score

Figures 12 and Figures 13 present the sparse patterns across various heads in the ChatGLM3-6B model (28 layers x 32 heads) under a sequence length of 61K. We conducted row-by-row filtering based on the full attention softmax weight, using a **CRA** threshold of $\alpha = 0.95$, and randomly selected four heads from different layers for display.

According to the visualization results on the majority of heads, we observed two distinct and prominent patterns prevalent in the heatmap of attention weight: column stripes and slash stripes. Column stripe patterns embody the global contextual information whereas slash stripe capture local information.

### A.2 Sparsity Analysis

To further quantify the degree of sparsity exposed as sequence lengths increase, we conducted scalability tests on the ChatGLM3-6B model using the "Needle-in-a-Haystack" task to evaluate sparsity. The results are presented in Table 4. According to the results, the increase in sequence length introduces more apparent sparsity. With each doubling of length, the proportion of KV elements needed to maintain the same threshold $\alpha$ decreases by approximately 20%. Concurrently, a smaller threshold results in the filtering of more KV elements, which may also lead to a decline in task performance in accuracy.

*Table 4.* Sparsity analysis for the "Needle in a Haystack" task

| Sequence Length | Average Sparsity in ChatGLM-6B | Average Sparsity in InternLM-7B |
|:---:|:---:|:---:|
| 4K | 88.00% | 91.13% |
| 8K | 90.74% | 92.72% |
| 16K | 92.52% | 93.89% |
| 32K | 93.88% | 94.83% |
| 64K | 94.89% | 95.89% |
| 128K | 95.84% | 96.67% |

### A.3 Effectiveness of sampling

To validate the efficiency of this chunked sampling method, we conducted tests using two different sampling ratios on three different heads. One sampling rate was set to full sampling, while the other was configured with a query block size of 128 and a sampling rate where $chunk_n$ was 2. We applied varying proportions of `top-k` columns and slash stripes to the full attention matrix as masks to observe the changes in **CRA**. The results, shown in Table 5, indicate that the CRA obtained by selecting `top-k` stripes at a 0.4% sampling rate is very close to the CRA obtained from the complete attention scores, with the difference decreasing as the proportion increases. This demonstrates that the chun-

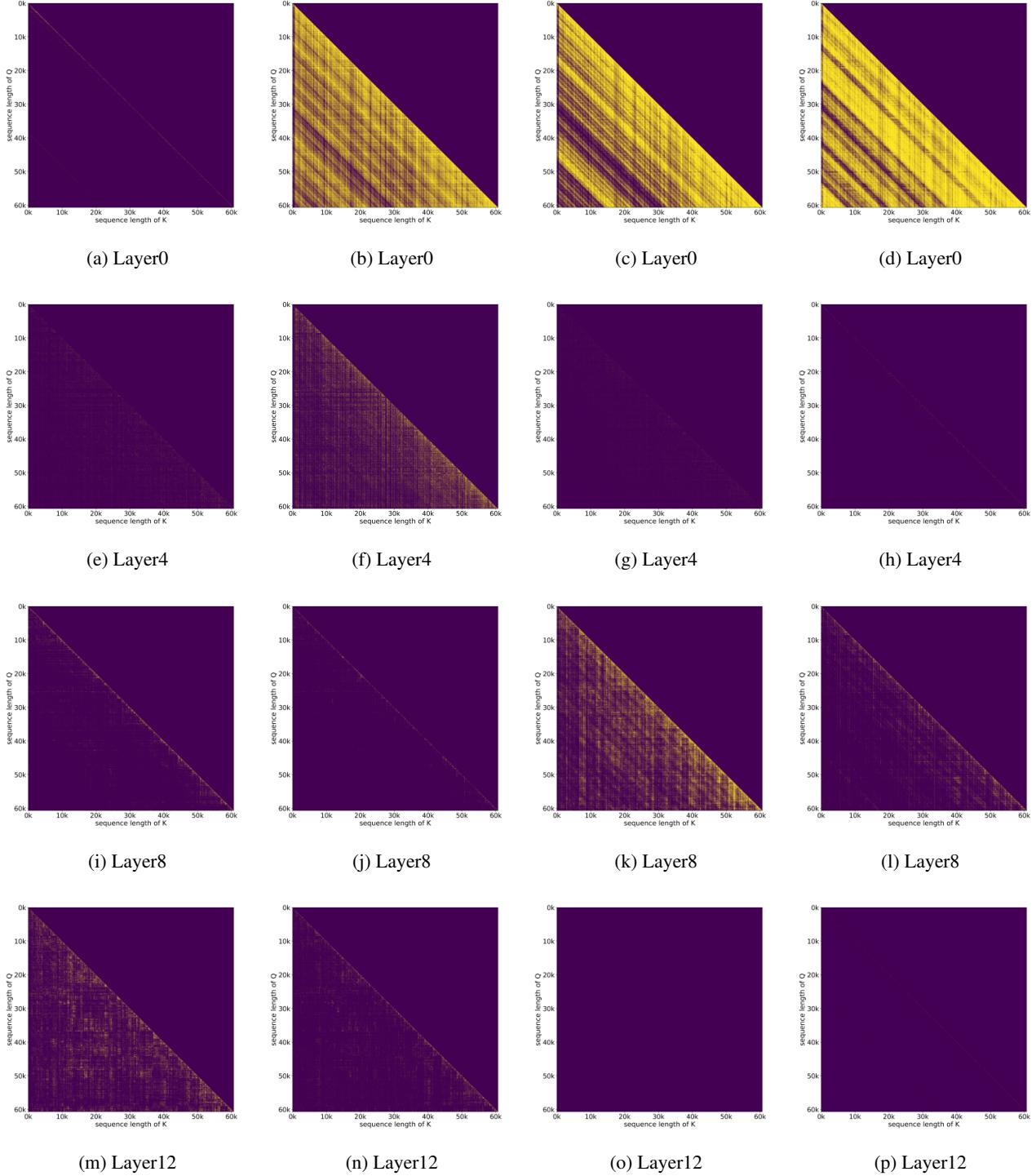ked sampling method in SampleAttention is both simple and efficient.

*Figure 12.* The visualization attention based on a content length of 61K, displays the sparse patterns for randomly chosen heads from layers 0, 4, 8 and 12.

(a) Layer16     (b) Layer16     (c) Layer16     (d) Layer16

(e) Layer20     (f) Layer20     (g) Layer20     (h) Layer20

(i) Layer24     (j) Layer24     (k) Layer24     (l) Layer24
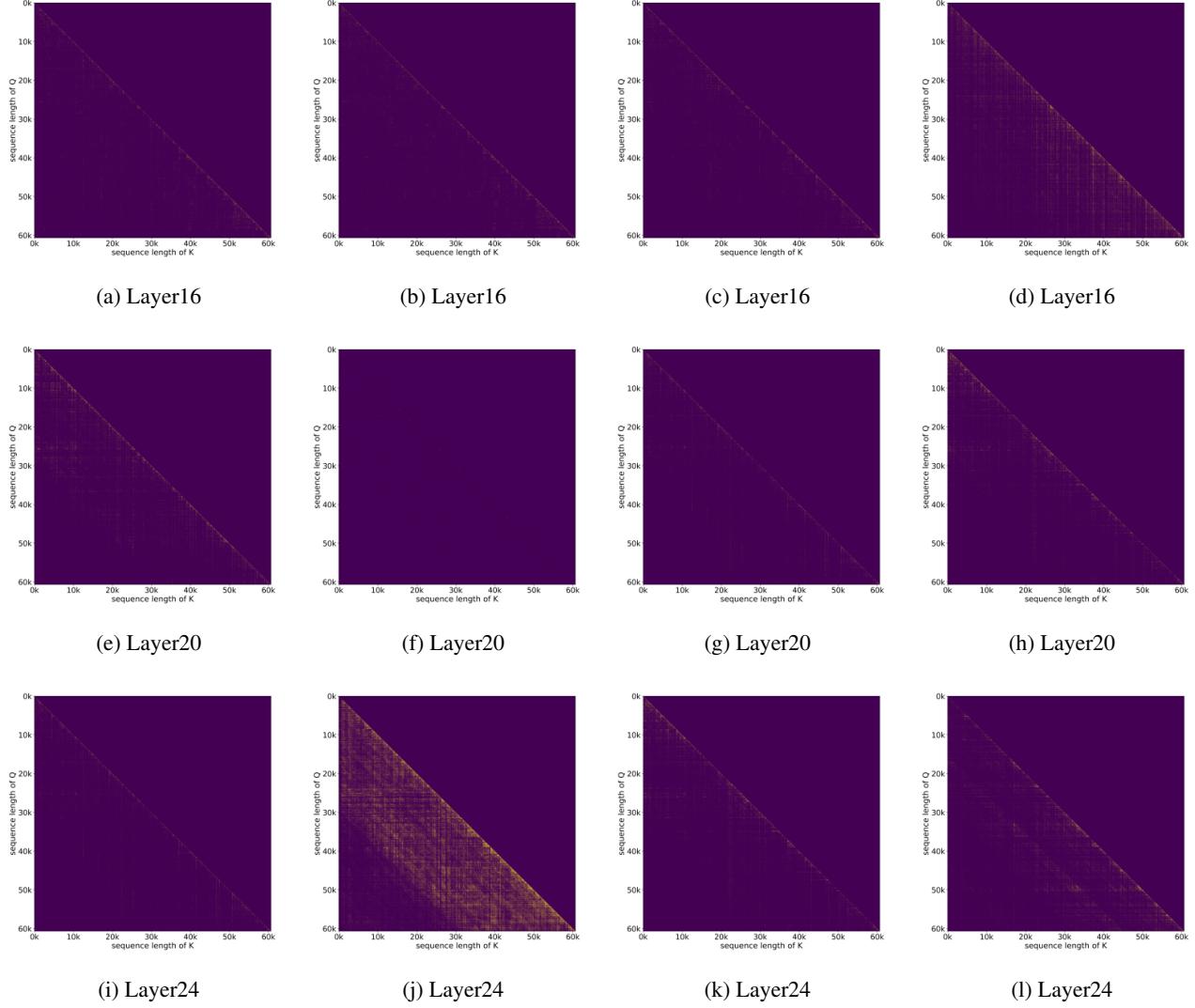
*Figure 13.* The visualization attention based on a content length of 61K, displays the sparse patterns for randomly chosen heads from layers 16, 20 and 24.

*Table 5.* The **CRA** percentages can be achieved by selecting different ratios of *top-k* stripes under varying sampling rates. The tests were conducted on the RULER task with a sequence length of 64K using the GLM4-9B model.

| ratio of `top-k` | 2.5% | | 10% | | 20% | | 40% | | 80% | |
|---|---|---|---|---|---|---|---|---|---|---|
| sampling ratio | 100% | 0.4% | 100% | 0.4% | 100% | 0.4% | 100% | 0.4% | 100% | 0.4% |
| HEAD-1 | 16.35% | 12.74% | 26.91% | 23.33% | 45.99% | 42.14% | 58.21% | 55.34% | 96.30% | 93.65% |
| HEAD-2 | 55.43% | 48.40% | 63.89% | 58.63% | 85.92% | 81.98% | 89.07% | 84.21% | 99.15% | 98.08% |
| HEAD-3 | 93.20% | 90.44% | 98.32% | 97.62% | 99.14% | 98.43% | 99.41% | 99.12% | 99.98% | 99.66% |

*Table 6.* Hyperparameters selected for GLM-4B and InterLM2-7B after tuning across sequence length ranges. Acc | Spd represents the optimal configuration that maintains accuracy without loss compared to the full-attention, as well as the configuration achieving the highest speedups while retaining at least 90% of the accuracy.

| range of sequence length | < 16K | | [16K,48K) | | [48K,80K) | | [80K,112K) | | >=112K | |
|---|---|---|---|---|---|---|---|---|---|---|
| models Acc \| Spd | GLM4 | InternLM2 | GLM4 | InternLM2 | GLM4 | InternLM2 | GLM4 | InternLM2 | GLM4 | InternLM2 |
| CRA Column | 0.98 \| 0.85 | 0.98 \| 0.85 | 0.95 \| 0.85 | 0.95 \| 0.80 | 0.95 \| 0.80 | 0.92 \| 0.80 | 0.90 \| 0.80 | 0.92 \| 0.80 | 0.95 \| 0.80 | 0.95 \| 0.80 |
| CRA Slash | 0.90 \| 0.85 | 0.95 \| 0.85 | 0.95 \| 0.80 | 0.90 \| 0.80 | 0.95 \| 0.80 | 0.92 \| 0.80 | 0.90 \| 0.80 | 0.90 \| 0.80 | 0.85 \| 0.80 | 0.90 \| 0.80 |
| Num of Chunks | 1 \| 1 | 1 \| 1 | 1 \| 1 | 1 \| 1 | 1 \| 1 | 1 \| 1 | 2 \| 1 | 2 \| 1 | 1 \| 1 | 2 \| 1 |