

# ADAMAS: HADAMARD SPARSE ATTENTION FOR EFFICIENT LONG-CONTEXT INFERENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) now support context windows of hundreds of thousands to millions of tokens, enabling applications such as long-document summarization, large-scale code synthesis, multi-document question answering and persistent multi-turn dialogue. However, such extended contexts exacerbate the quadratic cost of self-attention, leading to severe latency in autoregressive decoding. Existing sparse attention methods alleviate these costs but rely on heuristic patterns that struggle to recall critical key-value (KV) pairs for each query, resulting in accuracy degradation. We introduce **Adamas**, a lightweight yet highly accurate sparse attention mechanism designed for long-context inference. Adamas applies the Hadamard transform, bucketization and 2-bit compression to produce compact representations, and leverages Manhattan-distance estimation for efficient top- $k$  selections. Experiments show that Adamas matches the accuracy of full attention with only a 64-token budget, achieves near-lossless performance at 128, and supports up to  $8\times$  higher sparsity than prior state-of-the-art (SOTA) methods while delivering up to  $4.4\times$  self-attention and  $1.5\times$  end-to-end speedups on 32K-length sequences. Remarkably, Adamas attains comparable or even lower perplexity than full attention, underscoring its effectiveness in maintaining accuracy under aggressive sparsity. Code is publicly available at <https://anonymous.4open.science/r/Adamas-36EA>.

## 1 INTRODUCTION

The rapid progress of LLMs has dramatically extended their affordable context windows. Contemporary systems such as Anthropic’s Claude Sonnet 4 support up to 1M tokens (Anthropic, 2025), while OpenAI’s GPT-5 effectively handles contexts of 128K–256K tokens (OpenAI, 2025). These expanded capacities enable advanced applications, including long-document summarization (Chang et al., 2024), large-scale code synthesis (Nijkamp et al., 2023; Dainese et al., 2024), multi-document question answering (Wu et al., 2025a; Wang et al., 2024b), and persistent multi-turn dialogue (Wu et al., 2025b). By processing extensive information without manual segmentation, LLMs are increasingly able to tackle tasks that demand global coherence and long-term memory. However, these impressive expansions come with a cost. The quadratic complexity of self-attention (Zaheer et al., 2020; Vaswani et al., 2017), together with the scaling of the KV cache, leads to significant per-token latency and memory overhead.

Sparse attention offers a promising solution by restricting each query to attend only a carefully selected subset of tokens (Yuan et al., 2025). This reduces the number of KV pairs involved in attention computation, lowering both computational complexity and memory access costs while largely preserving modeling capacity. Nevertheless, existing approaches often fall into two categories with notable limitations. Static sparsity patterns (e.g., StreamingLLM (Xiao et al., 2023)) are often hand-designed from empirical attention heatmaps, yielding fixed patterns, such as fixed local windows or vertical stripes. However, these patterns fail to capture the dynamic nature of query–key interactions, leading to low recall and degraded accuracy (See Table 2). In contrast, dynamic methods like Quest (Tang et al., 2024) adapt token selection during inference, but their page-level granularity remains overly coarse. This coarse selection introduces token redundancy and limits the achievable sparsity ratio, since higher sparsity levels lead to accuracy degradation.

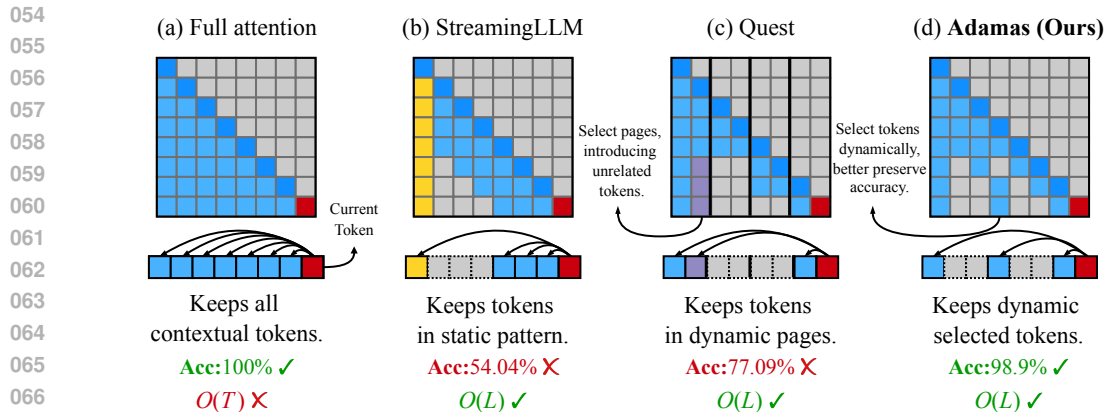


Figure 1: Illustration of Adamas compared with existing methods. While StreamingLLM employs a fixed sparse pattern and Quest inherently selects pages, Adamas dynamically selects KV pairs at the token level, thereby achieving better preservation of model accuracy and inference efficiency.<sup>1</sup>

In this paper, we propose Adamas, a lightweight yet highly accurate token-level sparse attention mechanism for long-context inference. Adamas achieves high sparsity while maintaining the **attention-quality model accuracy**, as illustrated in Figure 1. ~~Inspired by QuaRot (Ashkboos et al., 2024), which uses the Hadamard transform to suppress outlier features, we extend this idea by combining the Hadamard transform with bucketization to efficiently approximate query-key similarity at token-level.~~ Inspired by the general idea of using orthogonal transforms to stabilize feature distributions, we incorporate the classical Hadamard transform into Adamas to uniformly redistribute information across dimensions, thereby improving the robustness of our bucketization-based similarity estimation at the token level. In Adamas, queries and keys are first transformed by the Hadamard transform and then compressed into 2-bit codes via bucketization, which are stored in the KV cache with negligible overhead. During decoding, candidate keys are rapidly pre-selected using a lightweight Manhattan-distance estimator on the compressed codes, followed by top- $k$  filtering and sparse attention over the reduced candidate set. The proposed Adamas delivers substantial efficiency gains while preserving attention quality comparable to dense methods. We evaluate both the accuracy and efficiency of Adamas. Due to the dynamical selection of the most relevant KV pairs for each query, Adamas achieves up to  $8\times$  higher sparsity than previous SOTA methods while preserving comparable accuracy with full attention under a constrained budget of 128 tokens. Comprehensive evaluations demonstrate that Adamas delivers up to  $4.4\times$  self-attention and  $1.5\times$  end-to-end speedups on 32K-length sequences, accompanied by perplexity that is even lower than that of full attention, outperforming prior SOTA methods by more than  $2\times$ . The main contributions of this work are as follows:

- We propose Adamas, a novel sparse-attention mechanism that integrates Hadamard transform, bucketization, 2-bit compression and Manhattan-distance estimator, achieving  $8\times$  higher sparsity than prior SOTA methods while preserving comparable accuracy with full attention under constrained budget of 128 tokens.
- We develop high-performance GPU kernels for Adamas, featuring fused bucketization and 2-bit compression together with a lightweight Manhattan-distance estimator, enabling up to  $4.4\times$  self-attention and  $1.5\times$  end-to-end speedups over full attention in long-context decoding.
- Through extensive ablation studies, we validate the effectiveness of each component in Adamas, demonstrating that every step, including Hadamard transform, bucketization, compression, and estimation, contributes to its overall performance and efficiency.

<sup>1</sup>Accuracy (Acc) results are derived from the normalized performance of our evaluations on LongBench. Detailed results can be found in Figure 3 and Table 1.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

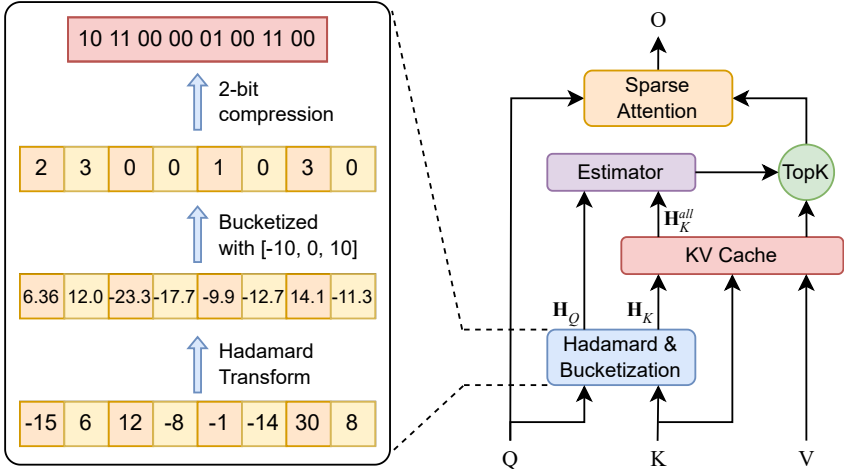


Figure 2: Overview of Adamas. Queries  $Q$  and keys  $K$  are processed through Hadamard transform, bucketization, and 2-bit compression. The transformed keys  $\mathbf{H}_K$  are then compared against the transformed query  $\mathbf{H}_Q$  in Manhattan-distance estimator, based on which the top- $k$  KV pairs are selected. Finally, Adamas performs sparse attention using  $Q$  and the selected KV pairs.

## 2 PRELIMINARIES

In this section, we briefly review sparse attention and the Hadamard transform to provide background and facilitate a clearer understanding of our proposed method.

**Sparse attention** is a variant of the self-attention mechanism. The core process is defined as

$$O = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V, \tag{1}$$

where  $Q$ ,  $K$ ,  $V$ , and  $O$  denote the query, key, value, and attention output, respectively, and  $d$  is the head dimension. Since  $QK^\top$  computes the query–key similarity, achieving high recall requires selecting the most relevant keys for each query. However, evaluating  $QK^\top$  against all keys incurs quadratic complexity of sequence length, making it impractical for long contexts. Sparse attention addresses this challenge by efficiently approximating  $QK^\top$  with only a carefully chosen subset of keys, thereby reducing both computational and memory costs while preserving the model’s accuracy.

**Hadamard transform**, also known as the Walsh–Hadamard transform, is an orthogonal linear transform that projects the vector to a new basis defined by Walsh functions. The Hadamard transform corresponds to multiplication by a Hadamard matrix  $\mathbf{H}_d$ , which is a square matrix of order  $d = 2^n$  with entries restricted to  $\{+1, -1\}$ . The smallest Hadamard matrix is defined as:

$$\mathbf{H}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \tag{2}$$

Higher-order matrices are constructed recursively via the Kronecker product  $\mathbf{H}_{2^n} = \mathbf{H}_2 \otimes \mathbf{H}_{2^{n-1}}$ . This recursive structure enables an efficient Hadamard transform algorithm that computes the matrix–vector product  $\mathbf{H}x$  in  $\mathcal{O}(d \log_2 d)$  rather than the naive  $\mathcal{O}(d^2)$ . For dimensions  $d$  that are not exact powers of two, the existence of a Hadamard matrix is not guaranteed. In such cases, one may exploit factorizations  $d = 2^nm$ , where  $m$  is the order of a known Hadamard matrix, and apply the Kronecker construction  $\mathbf{H}_d = \mathbf{H}_{2^n} \otimes \mathbf{H}_m$ , yielding a transform with complexity  $\mathcal{O}(d(m+n))$ .

## 3 METHOD

In this section, we present Adamas, an efficient variant of the Transformer sparse attention mechanism that could substantially reduce computation overhead while preserving model accuracy. We show the workflow of Adamas in Figure 2.

**Algorithm 1** Workflow of Adamas

---

**Input:** Query  $\mathbf{Q}$ , Key  $\mathbf{K}$ , and Value  $\mathbf{V} \in \mathbb{R}^d$ , Hadamard matrix  $\mathbf{H} \in \mathbb{R}^{d \times d}$   
**Output:** Attention output matrix  $\mathbf{O}$

- 1:  $\mathbf{H}_Q \leftarrow \mathbf{Q}\mathbf{H}$ ;     $\mathbf{H}_K \leftarrow \mathbf{K}\mathbf{H}$  ▷ Apply Hadamard transform
- 2:  $\widehat{\mathbf{H}}_Q, \widehat{\mathbf{H}}_K \leftarrow \text{Bucketize}(\mathbf{H}_Q, \mathbf{H}_K)$  ▷ Quantize into  $\{0, 1, 2, 3\}^d$
- 3:  $\widehat{\mathbf{H}}_Q, \widehat{\mathbf{H}}_K \leftarrow \text{Compress}(\widehat{\mathbf{H}}_Q, \widehat{\mathbf{H}}_K)$  ▷ Pack into 2-bit codes
- 4:  $\widehat{\mathbf{H}}_K^{all}, \mathbf{K}^{all}, \mathbf{V}^{all} \leftarrow \text{KVCache.update}(\widehat{\mathbf{H}}_K, \mathbf{K}, \mathbf{V})$
- 5:  $D \leftarrow \text{ManhattanDistance}(\widehat{\mathbf{H}}_Q, \widehat{\mathbf{H}}_K^{all})$  ▷ Estimate query–key similarity
- 6:  $I \leftarrow \text{Top-}k(D)$  ▷ Select top- $k$  candidate indices
- 7:  $\mathbf{K}^s, \mathbf{V}^s \leftarrow \mathbf{K}^{all}[I], \mathbf{V}^{all}[I]$
- 8:  $\mathbf{O} \leftarrow \text{SparseAttention}(\mathbf{Q}, \mathbf{K}^s, \mathbf{V}^s)$
- 9: **return**  $\mathbf{O}$

---

## 3.1 MOTIVATIONS

The design of Adamas is motivated by the dual goals of theoretical equivalence and practical efficiency: leveraging the mathematical equivalence of Hadamard-transformed similarities with the original attention formulation, while exploiting the smoothing property of the transform to enable effective low-bit quantization and lightweight similarity estimation.

**Theoretical perspective.** The Hadamard transform is an orthogonal transformation, which ensures that the similarity computation in the Hadamard domain is mathematically equivalent to that in the original space. Specifically, for queries  $Q$  and keys  $K$ , we have

$$(Q\mathbf{H})(K\mathbf{H})^\top = Q(\mathbf{H}\mathbf{H}^\top)K^\top = QK^\top \quad (3)$$

where  $\mathbf{H}$  denotes a Hadamard matrix with  $\mathbf{H}\mathbf{H}^\top = \mathbf{I}$ . This equivalence shows that applying the Hadamard transform to both  $Q$  and  $K$  does not incur any information loss.

**Practical perspective.** In practice, directly quantizing query and key vectors can incur severe information loss due to the presence of large-magnitude outlier values, which dominate the distribution. The Hadamard transform mitigates this issue by redistributing variance more evenly across dimensions and suppressing extreme outliers, thereby producing smoother value distributions (Elhage et al., 2023; Ashkboos et al., 2024). As a result, bucketization after the Hadamard transform can approximate the original similarity structure with minimal degradation, enabling compact 2-bit representations that significantly reduce memory overhead while retaining sufficient accuracy. Moreover, since the Hadamard matrix contains only  $\pm 1$  entries, its application can be implemented via fast Hadamard transform (Dao, 2024; IBM & Meta, 2024) rather than costly dense matrix multiplications, further reducing computational overhead.

## 3.2 ADAMAS

As shown in Algorithm 1, Adamas modifies the standard Transformer attention pipeline with three innovations: 1) Hadamard transform applied to queries and keys (yielding what we call Hadamard vectors), 2) bucketization and 2-bit compression for Hadamard vectors, and 3) a Manhattan-distance estimator for candidate token selection. These components work jointly to enable faster attention computation with modest memory cost and near lossless accuracy degradation.

**Hadamard transform applied to queries and keys.** Formally, given queries and keys  $Q, K \in \mathbb{R}^d$ , we compute their Hadamard-transformed representations as

$$H_Q = Q\mathbf{H}, \quad H_K = K\mathbf{H}, \quad (4)$$

where  $\mathbf{H} \in \mathbb{R}^{d \times d}$  denotes a Hadamard matrix.

**Bucketization and 2-bit compression.** For efficient computation and storage, each element in  $H_Q, H_K \in \mathbb{R}^d$  is bucketized into one of four levels using predefined thresholds  $\{B_1, B_2, B_3\}$ , mapped to the discrete set  $\{0, 1, 2, 3\}$ , which can be encoded into 2-bit integer. The bucketization operator  $B(\cdot)$  is defined as

$$B(x) = \sum_i \mathbb{I}(x > B_i), \quad (5)$$

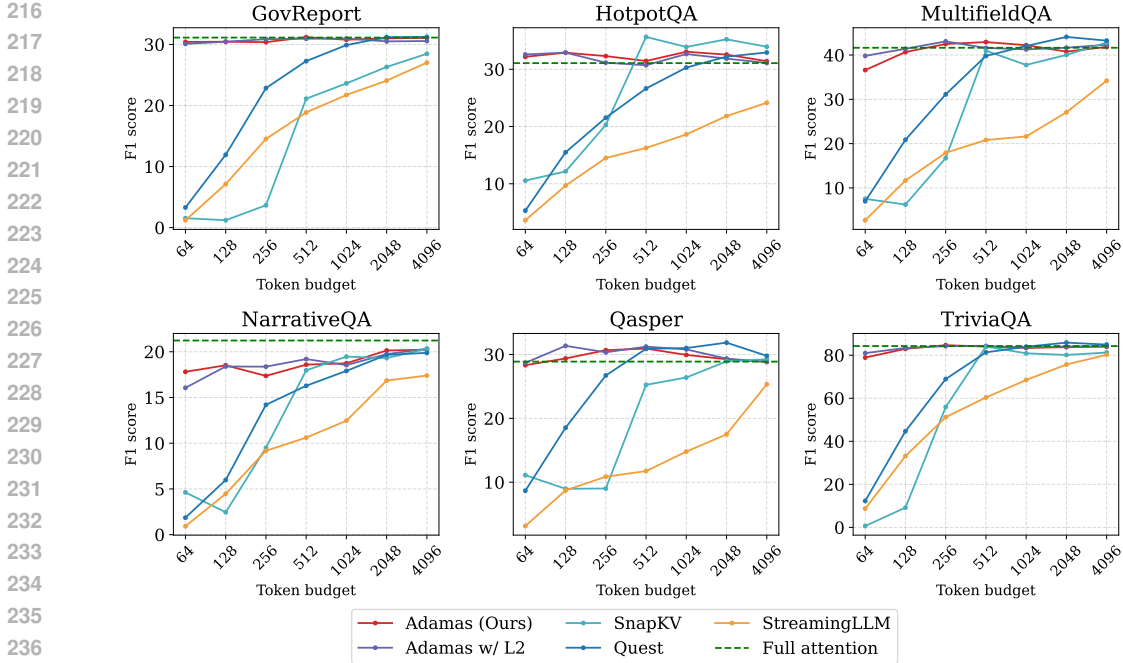


Figure 3: Evaluation results on LongBench. Adamas exhibits the smallest performance drop compared to full attention while maintaining high sparsity.

where  $\mathbb{I}(\cdot)$  denotes the indicator function, which evaluates to 1 if the condition inside holds, and 0 otherwise. We further analyze the criterion of bucketization threshold selection in Appendix C. The bucketized Hadamard vectors  $\hat{\mathbf{H}}_Q$  and  $\hat{\mathbf{H}}_K$  are obtained by applying  $B(\cdot)$  element-wise to  $\mathbf{H}_Q$  and  $\mathbf{H}_K$ , respectively.

We further pack every eight elements into a single 16-bit value. The compressed  $\hat{\mathbf{H}}_K$  is stored directly in the KV cache, which increases cache size by only 1/16, and is later used for lightweight similarity estimation. This strategy significantly reduces the memory footprint of the Hadamard-transformed vectors while retaining sufficient information for effective candidate token selection.

**Manhattan distance estimation.** The third part in Adamas is a similarity estimator based on Manhattan distance, operating directly on the 2-bit compressed representations. Given a compressed query  $\hat{\mathbf{H}}_Q \in \{0, 1, 2, 3\}^d$  and a compressed key  $\hat{\mathbf{H}}_K \in \{0, 1, 2, 3\}^d$ , we approximate similarity by the negative Manhattan distance:

$$\text{sim}(\hat{\mathbf{H}}_Q, \hat{\mathbf{H}}_K) \approx -\|\hat{\mathbf{H}}_Q - \hat{\mathbf{H}}_K\|_1. \tag{6}$$

Since  $\hat{\mathbf{H}}_Q$  and  $\hat{\mathbf{H}}_K$  are encoded as 2-bit integers, the similarity computation can be carried out with bit-wise integer operations instead of relatively expensive floating-point arithmetic. To fully exploit this compression, we further design custom kernels that efficiently process groups of eight elements in each computation step, achieving both a low memory footprint and low computational overhead.

## 4 EXPERIMENTS

In this section, we evaluate Adamas in terms of both efficacy and efficiency. We also conduct ablation studies on the Hadamard transform, bucketization, and distance metrics to better understand the contribution of each component.

### 4.1 SETUP

**Models.** We use LongChat-v1.5-7b-32k (Li et al., 2023) and Yarn-Llama-2-7b-128k (Peng et al., 2024) for evaluation, following the experimental settings of Quest (Tang et al., 2024). LongChat-

Table 1: Average accuracy (%) of all methods evaluated on LongBench, normalized by full attention results.

Methods	Adamas(Ours)	Adamas w/ L2	Quest	SnapKV	StreamingLLM
Average Accuracy	98.89	98.22	77.09	66.23	54.04

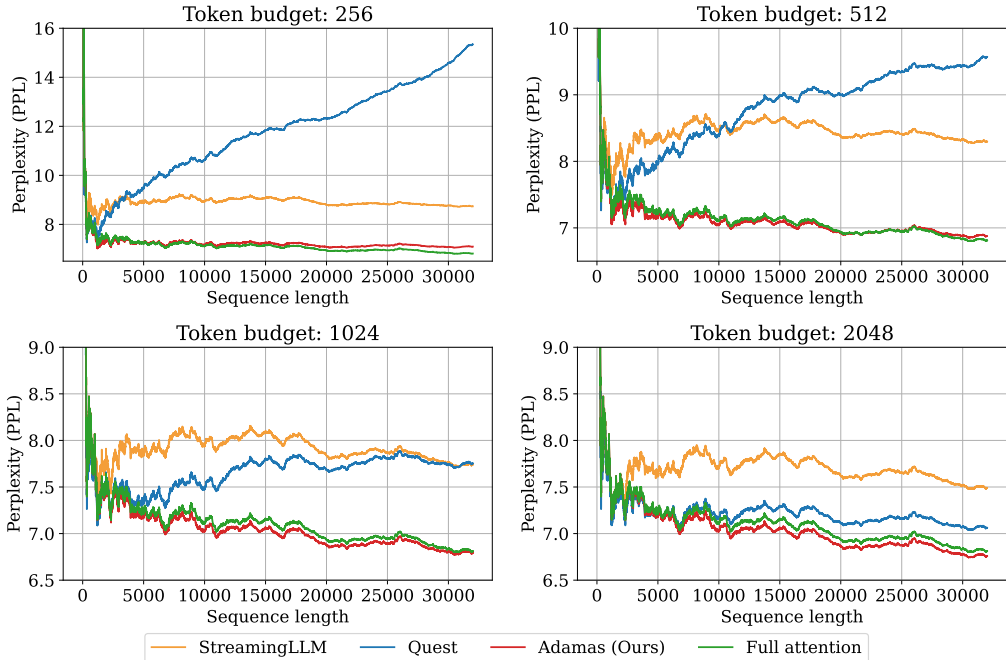


Figure 4: Perplexity results of StreamingLLM, Quest, Adamas, and full attention evaluated on PG19 with LongChat-7b-v1.5-32k under varying token budgets. Adamas consistently matches full attention and even shows lower perplexity at larger token budget.

v1.5-7b-32k is used in most experiments, while Yarn-Llama-2-7b-128k is employed for extremely long-context scenarios (up to 100K tokens).

**Tasks.** We consider three categories of evaluation tasks: (1) PG19 (Rae et al., 2019), a long-form language modeling dataset for assessing prediction confidence; (2) the passkey retrieval task (Peng et al., 2024), which measures retrieval capability in long contexts; and (3) six datasets from LongBench (Bai et al., 2024), a benchmark designed for long-context understanding. The LongBench tasks cover multiple settings, including single-document QA (NarrativeQA (Kočíšký et al., 2018), Qasper (Dasigi et al., 2021), MultiFieldQA (Bai et al., 2024)), multi-document QA (HotpotQA (Yang et al., 2018)), summarization (GovReport (Huang et al., 2021)), and few-shot learning (TriviaQA (Joshi et al., 2017)).

**Baselines.** For fair comparison, we benchmark against training-free sparse attention methods. Specifically, we include StreamingLLM (Xiao et al., 2023) as a representative static-sparse-mask method, SnapKV (Li et al., 2024) and Quest (Tang et al., 2024) as representative dynamic-KV-selection methods.

#### 4.2 EFFICACY EVALUATION

**LongBench evaluation.** We evaluate Adamas and baselines on six datasets from LongBench. As shown in Figure 3, Adamas consistently surpasses these prior SOTA methods across all datasets, with particularly strong advantages under low token budgets, demonstrating its ability to preserve critical KV pairs. In comparison, StreamingLLM consistently underperforms full attention due to its KV cache eviction strategy. SnapKV exhibits mediocre overall accuracy and only achieves barely

Table 2: Passkey retrieval accuracy (%) of StreamingLLM, Quest, and Adamas under different token budgets. Results highlight the limitations of StreamingLLM’s sliding window design, the budget sensitivity of Quest, and the robustness of Adamas across both short and long contexts.

(a) Results of 10K length passkey retrieval test on LongChat-7b-v1.5-32k

Methods / Budget	16	32	64	128	256	512
StreamingLLM	1%	1%	1%	1%	3%	5%
Quest	52%	67%	<b>99%</b>	98%	100%	100%
Adamas(Ours)	<b>68%</b>	<b>85%</b>	93%	<b>98%</b>	<b>100%</b>	<b>100%</b>

(b) Results of 100K length passkey retrieval test on Yarn-Llama-2-7b-128k

Methods / Budget	64	128	256	512	1024	2048	4096
StreamingLLM	1%	1%	1%	1%	1%	2%	4%
Quest	25%	58%	84%	95%	<b>99%</b>	99%	99%
Adamas(Ours)	<b>54%</b>	<b>71%</b>	<b>87%</b>	<b>95%</b>	98%	<b>100%</b>	<b>100%</b>

acceptable performance when the token budget exceeds 512, yet it still falls short of full attention. Quest achieves competitive performance when the token budget exceeds 1024, but its accuracy drops sharply below 512. Specifically, as shown in Table 1, Adamas achieves over a 20% accuracy improvement over Quest, more than a 30% improvement over SnapKV, and accuracy comparable to Adamas with L2 on average accuracy values normalized by full attention results. For more detailed results, please refer to Table 4. These results confirm that Adamas sustains reliable accuracy across diverse long-context scenarios, even under constrained token budgets.

**Perplexity results.** We evaluate perplexity on the PG19 dataset using LongChat-7b-v1.5-32k with a 32K sequence length. As shown in Figure 4, Adamas closely tracks full attention across all settings and even achieves lower perplexity at larger token budgets, demonstrating its accuracy and robustness. In comparison, Quest suffers from high perplexity under high sparsity conditions, as its page-wise strategy introduces many irrelevant KV pairs, obscuring truly relevant ones when the token budget is limited. StreamingLLM maintains stable but consistently higher perplexity due to its eviction strategy, which permanently discards past information.

**Passkey results.** We evaluate passkey retrieval on LongChat-7b-v1.5-32k with 10K-length inputs and on Yarn-Llama-2-7b-128k with 100K-length inputs. Results are listed in Table 2. Adamas maintains consistently high accuracy across all budgets and significantly outperforms prior SOTA under constrained settings. In comparison, StreamingLLM performs poorly across all budgets due to its sliding window design, which discards keys outside the retained span. Quest achieves strong accuracy with large budgets but degrades sharply when the budget is small ( $\leq 32$  for 10K and  $\leq 128$  for 100K), as its page-wise partitioning fails to recall scattered keys. Overall, these results highlight that Adamas robustly retrieves critical KV pairs even under severely limited token budgets.

### 4.3 EFFICIENCY EVALUATION

We conduct efficiency evaluations with LongChat-7b-v1.5-32k on NVIDIA RTX A6000.

**Kernels.** We implement high-performance CUDA kernels for the Hadamard transform, fused bucketization-compression, Manhattan-distance estimation, top- $k$  selection, and sparse attention. Manhattan-distance estimation, the most computationally intensive component (shown in Appendix Table 8), is significantly accelerated through lightweight bit-wise integer operations enabled by our bucketization and compression design. The Hadamard transform and fused bucketization-compression incur negligible overhead and are thus omitted from further analysis. As shown in Figure 5, Adamas achieves consistent acceleration across various token budgets and sequence lengths. In particular, it delivers up to  $4.4\times$  speedup over full attention implemented by FlashInfer (Ye et al., 2025) under a 256-token budget, with nearly lossless accuracy.

**End-to-end.** Table 3 shows that as the token budget increases from 256 to 4096, Adamas maintains stable end-to-end latency and achieves up to  $1.5\times$  speedup over FlashInfer on 32K-length

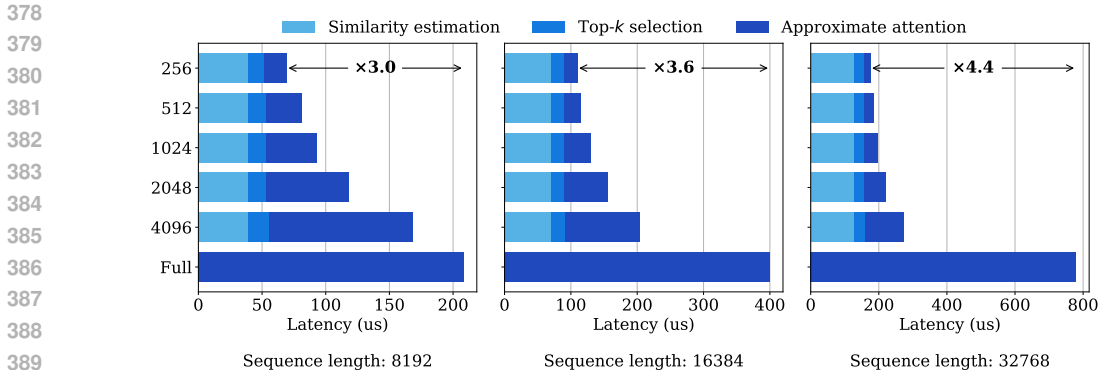


Figure 5: Kernel-level breakdown of Adamas self-attention compared to full attention (FlashInfer).

Table 3: End-to-end decoding latency (ms). Adamas achieves up to  $1.5\times$  speedup.

Sequence / Budget	256	512	1024	2048	4096	Full Attn
8192	26.0	26.3	26.7	27.5	28.8	30.2
16384	29.3	29.5	30.0	30.9	32.3	38.0
32768	34.9	35.2	35.8	36.9	38.6	53.7

sequences. Since the attention module accounts for only 30% – 40% of the total decoding time, the overall acceleration is somewhat diluted compared to the kernel-level results. Nevertheless, Adamas exhibits only marginal increases in latency with larger token budgets, highlighting its scalability and efficiency in long-context scenarios. We further conduct end-to-end efficiency evaluations on an A800 80G PCIe GPU, which better reflects real-world server environments. Notably, Adamas achieves up to a  $1.4\times$  speedup on 32K-length sequences, demonstrating substantial and practically meaningful efficiency gains. Please refer to Appendix B for more detailed results.

#### 4.4 ABLATION STUDY

In this part, we conduct ablation studies on the Hadamard transform, bucketization, and distance metrics. To better understand the contribution of each component, we benchmark three variants of Adamas on LongBench using LongChat-7b-v1.5-32k: (1) Adamas without the Hadamard transform, (2) Adamas with different bucketization granularities, and (3) Adamas with L2 distance (replacing Manhattan distance). The results are shown in Figure 6, with detailed data provided in Appendix H.

**Hadamard transform.** As illustrated in Figure 6, Adamas without the Hadamard transform (Adamas w/o Hadamard) directly bucketizes queries and keys before Manhattan-distance estimation, achieving near-zero scores across multiple datasets under small token budgets. Its performance improves only gradually with larger budgets and remains below the baseline even at a 4096-token budget. This behavior indicates that the model fails to generate accurate answers without the Hadamard transform, underscoring its critical role in Adamas: Hadamard transform smooths the value distribution of queries and keys, suppresses extreme outliers, and thereby mitigates the information loss introduced by bucketization.

**Bucketization.** To investigate how bucketization granularity affects accuracy, we compare three variants: 1-bit (Adamas-1bit, with threshold), 2-bit (Adamas), and 3-bit (Adamas-3bit). As shown in Figure 6, Adamas and Adamas-3bit achieve nearly identical performance, with Adamas-3bit showing only a slight advantage. Adamas-1bit performs comparably to Adamas when the token budget exceeds 256, but its accuracy drops significantly under smaller budgets, suggesting that retaining more bits helps preserve information. However, the marginal improvement of Adamas-3bit over Adamas highlights diminishing returns when increasing bit width. Since higher bit widths also incur greater storage costs, we adopt 2-bit bucketization as the most memory-efficient choice without sacrificing accuracy.

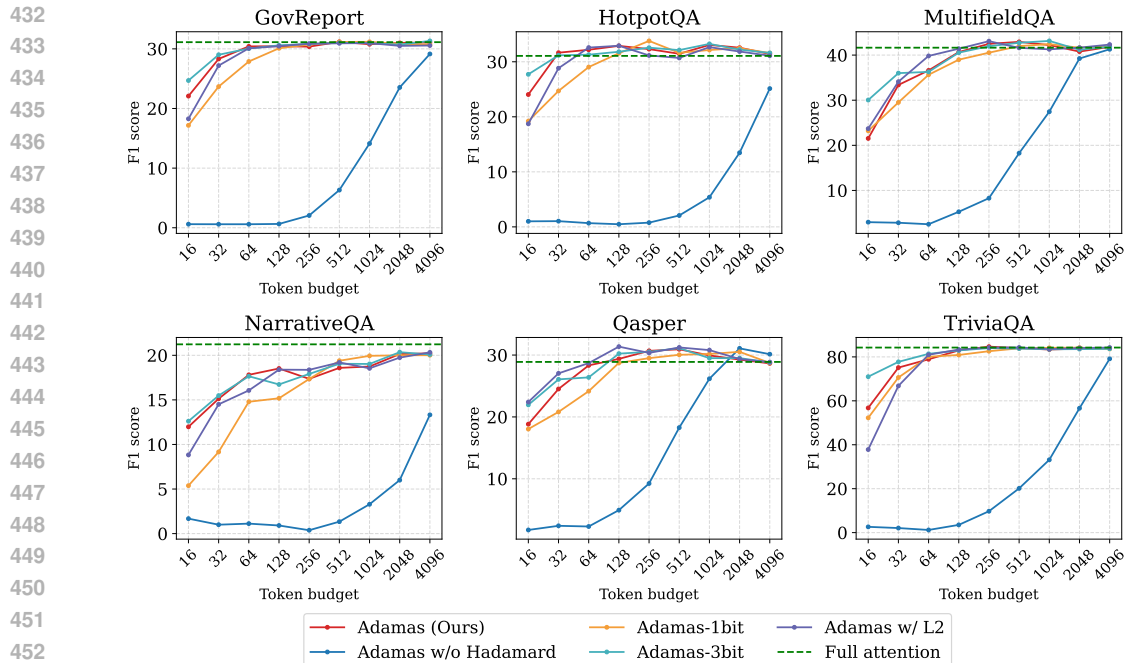


Figure 6: Ablation studies on the Hadamard transform, bucketization, and distance metrics.

**Distance metrics.** We replace Manhattan distance (L1) estimator with L2 distance estimator to assess the impact of distance metrics. As shown in Figure 6, Adamas with L2 distance achieves comparable performance with vanilla Adamas, showing similar performance curves. Specifically, Adamas with L2 performs slightly better on single-document QA such as MultifieldQA and Qasper, while vanilla Adamas shows margin advantages on other tasks. This discrepancy reflects the different sensitivities of the two distance metrics: L2 distance emphasizes global geometric consistency and is more effective for reasoning within a single coherent context, whereas L1 distance is more robust to noise and sparsity, making it better suited for integrating dispersed or partially relevant information across multiple sources.

## 5 RELATED WORK

A large amount of work has explored sparsity in attention mechanisms to improve computational efficiency, with a primary focus on how to effectively select critical KV pairs. Existing methods can be broadly categorized into three groups:

**Static sparse methods.** These methods adopt predefined sparsity patterns, such as sliding windows or attention sinks. For example, H2O (Zhang et al., 2023) and infLLM (Xiao et al., 2024) rely on the sliding window pattern, while StreamingLLM (Xiao et al., 2023) and DUOAttention (Xiao et al., 2025) combine sliding window with attention sink to preserve accuracy. Moreover, DUOAttention requires an extra training phase to determine streaming/retrieval heads, which introduces substantial training overhead. MOA (Wang et al., 2024a) further refines this idea by profiling each attention head across layers to adjust the window span and then incorporating attention sinks to build the sparse mask. Although these static designs yield high throughput, their fixed structures risk discarding tokens that are critical to a given query (Tang et al., 2024), often resulting in accuracy degradation when the token budget is constrained, therefore significantly suffering from retrieval tasks.

**Dynamic KV selection methods.** To better preserve accuracy, these approaches select KV pairs adaptively based on the input query. Quest (Tang et al., 2024) estimates page importance using an upper-bound approximation of attention weights and selects the top- $k$  pages for sparse attention, though its page-level granularity may miss critical tokens. MInference (Jiang et al., 2024) defines three dynamic patterns and performs online estimation per attention head to select the proper pattern,

486 while FlexPrefill (Lai et al., 2025) adaptively determines sparse patterns and selection ratios per head  
487 to accelerate long-sequence prefilling. Both methods, however, primarily target the prefilling stage  
488 rather than full-sequence decoding, therefore are orthogonal and can, in principle, be combined with  
489 [Adamas](#). SeerAttention (Gao et al., 2024) introduces a learnable gating mechanism that activates  
490 important blocks within the attention map, but at the cost of training additional parameters, which  
491 complicates deployment by preventing the model from being used out-of-the-box at inference.

492 **Training-based sparse methods.** Beyond static and dynamic masking, an alternative approach inte-  
493 integrates sparsity into the attention mechanism through specialized training procedures. Reformer (Ki-  
494 taev et al., 2020) replaces standard dot-product attention with locality-sensitive hashing (LSH), re-  
495 ducing complexity from  $O(L^2)$  to  $O(L \log L)$ . NSA (Yuan et al., 2025) adopts a dynamic hierarchi-  
496 cal sparse strategy that combines coarse-grained token compression, fine-grained token selection,  
497 and sliding window via gating to balance global context and local precision. MoBA (Lu et al., 2025)  
498 incorporates a Mixture-of-Experts (MoE) design (Shazeer et al., 2017) into sparse attention, using  
499 a gating mechanism to select historical blocks. While these methods often achieve better accuracy  
500 and efficiency than standard full attention, they require pretraining the model, which is computationally  
501 costly and limits their applicability. In contrast, our method is training-free and can be directly  
502 applied to pretrained LLMs, offering broad compatibility without retraining overhead.

503 In parallel to sparse attention, another line of research focuses on **KV cache compression**, which  
504 aims to reduce the memory footprint of stored key/value tensors. While decoding with a KV cache  
505 has a baseline time and space complexity of  $O(T)$  with respect to sequence length  $T$ , these two  
506 families of techniques address different bottlenecks. Sparse attention focuses on selecting the most  
507 relevant query-key pairs at the attention level, with the overall goal of approximating full attention  
508 computation using a subset of tokens to achieve acceleration. This reduces I/O during attention  
509 computation, especially under high sparsity, and thus primarily improves latency (i.e., reduces time  
510 complexity), rather than memory usage. In contrast, KV cache compression reduces the memory  
511 footprint of stored key/value tensors (space complexity). However, because compression or pruning  
512 of tokens is generally irreversible, such methods often incur accuracy loss and may introduce ad-  
513 ditional decoding latency. Moreover, simply compressing the KV cache alone may not necessarily  
514 reduce I/O during attention computation, limiting its effectiveness for speedup. Overall, sparse at-  
515 tention and KV cache compression are largely complementary: the former focuses on computational  
516 acceleration, while the latter focuses on memory efficiency. Each addresses a different bottleneck in  
517 efficient long-context decoding.

## 518 6 CONCLUSION

519  
520 In this paper, we introduce [Adamas](#), a lightweight yet efficient attention mechanism for long-  
521 context inference. By integrating the Hadamard transform, bucketization, 2-bit compression, and  
522 a Manhattan-distance estimator, [Adamas](#) dynamically selects the most relevant KV pairs for each  
523 query. This design enables [Adamas](#) to achieve accuracy comparable to full attention with only a 64-  
524 token budget, and to become nearly lossless at 128 tokens, attaining up to  $8\times$  higher sparsity than  
525 prior SOTA methods. Comprehensive evaluations demonstrate that [Adamas](#) delivers up to  $4.4\times$   
526 self-attention and  $1.5\times$  end-to-end speedups on 32K-length sequences, while achieving perplexity  
527 even lower than full attention and surpassing prior SOTA by more than  $2\times$ .

540 REPRODUCIBILITY STATEMENT

541  
542 To ensure reproducibility, we provide the complete implementation of Adamas in the anonymized  
543 repository, including algorithmic details, source code and scripts for running experiments. All  
544 datasets and models used in this work are publicly available. We are confident that with the pro-  
545 vided resources, readers can reproduce the entirety of our presented results.

547 REFERENCES

- 548  
549 Anthropic. Claude sonnet 4 now supports 1m tokens of context, 2025. [https://www.](https://www.anthropic.com/news/lm-context)  
550 [anthropic.com/news/lm-context](https://www.anthropic.com/news/lm-context), Accessed: 2025-09-04.
- 551 Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Pashmina Cameron, Martin  
552 Jaggi, Dan Alistarh, Torsten Hoeffler, and James Hensman. Quarot: Outlier-free 4-bit inference in  
553 rotated llms. *Advances in Neural Information Processing Systems*, 37:100213–100240, 2024.
- 554  
555 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,  
556 Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench: A bilingual,  
557 multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meet-*  
558 *ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3119–3137,  
559 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/  
560 2024.acl-long.172. URL <https://aclanthology.org/2024.acl-long.172>.
- 561 Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer. Boookscore: A systematic explo-  
562 ration of book-length summarization in the era of LLMs. In *The Twelfth International Confer-*  
563 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=7Ttk3RzDeu)  
564 [7Ttk3RzDeu](https://openreview.net/forum?id=7Ttk3RzDeu).
- 565 Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. Generating code  
566 world models with large language models guided by monte carlo tree search. In *The Thirty-*  
567 *eighth Annual Conference on Neural Information Processing Systems*, 2024. URL [https://](https://openreview.net/forum?id=9SpWvX9ykp)  
568 [openreview.net/forum?id=9SpWvX9ykp](https://openreview.net/forum?id=9SpWvX9ykp).
- 569  
570 Tri Dao. Fast hadamard transform in cuda, with a pytorch interface, 2024. [https://github.](https://github.com/Dao-AI-Lab/fast-hadamard-transform)  
571 [com/Dao-AI-Lab/fast-hadamard-transform](https://github.com/Dao-AI-Lab/fast-hadamard-transform), Accessed: 2025-09-23.
- 572 Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. A dataset of  
573 information-seeking questions and answers anchored in research papers. In Kristina Toutanova,  
574 Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cot-  
575 terrell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of*  
576 *the North American Chapter of the Association for Computational Linguistics: Human Lan-*  
577 *guage Technologies*, pp. 4599–4610, Online, June 2021. Association for Computational Linguis-  
578 tics. doi: 10.18653/v1/2021.naacl-main.365. URL [https://aclanthology.org/2021.](https://aclanthology.org/2021.naacl-main.365/)  
579 [naacl-main.365/](https://aclanthology.org/2021.naacl-main.365/).
- 580 Nelson Elhage, Robert Lasenby, and Christopher Olah. Privileged bases in the transformer residual  
581 stream. *Transformer Circuits Thread*, pp. 24, 2023.
- 582  
583 Yizhao Gao, Zhichen Zeng, Dayou Du, Shijie Cao, Peiyuan Zhou, Jiaying Qi, Junjie Lai, Hayden  
584 Kwok-Hay So, Ting Cao, Fan Yang, et al. Seerattention: Learning intrinsic sparse attention in  
585 your llms. *arXiv preprint arXiv:2410.13276*, 2024.
- 586 Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. Efficient attentions for  
587 long document summarization. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek  
588 Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou  
589 (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association*  
590 *for Computational Linguistics: Human Language Technologies*, pp. 1419–1436, Online, June  
591 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.112. URL  
592 <https://aclanthology.org/2021.naacl-main.112/>.
- 593 IBM and Meta. Hadacore: Tensor core accelerated hadamard transform kernel, 2024. [https://](https://pytorch.org/blog/hadacore/)  
[pytorch.org/blog/hadacore/](https://pytorch.org/blog/hadacore/), Accessed: 2025-09-23.

- 594 Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua  
595 Han, Amir H Abdi, Dongsheng Li, Chin-Yew Lin, et al. Minference 1.0: Accelerating pre-filling  
596 for long-context llms via dynamic sparse attention. *Advances in Neural Information Processing*  
597 *Systems*, 37:52481–52515, 2024.
- 598  
599 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly  
600 supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan  
601 (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*  
602 *(Volume 1: Long Papers)*, pp. 1601–1611, Vancouver, Canada, July 2017. Association for Com-  
603 putational Linguistics. doi: 10.18653/v1/P17-1147. URL [https://aclanthology.org/  
604 P17-1147/](https://aclanthology.org/P17-1147/).
- 605 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
606 *International Conference on Learning Representations*, 2020. URL [https://openreview.  
607 net/forum?id=rkgNKkHtvB](https://openreview.net/forum?id=rkgNKkHtvB).
- 608 Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,  
609 and Edward Grefenstette. The narrativeqa reading comprehension challenge. *Transactions of the*  
610 *Association for Computational Linguistics*, 6:317–328, 2018.
- 611  
612 Xunhao Lai, Jianqiao Lu, Yao Luo, Yiyuan Ma, and Xun Zhou. Flexprefill: A context-aware sparse  
613 attention mechanism for efficient long-sequence inference. In *The Thirteenth International Con-  
614 ference on Learning Representations*, 2025. URL [https://openreview.net/forum?  
615 id=OfjI1belrT](https://openreview.net/forum?id=OfjI1belrT).
- 616 Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph E. Gonzalez, Ion Stoica,  
617 Xuezhe Ma, and Hao Zhang. How long can open-source llms truly promise on context length?,  
618 June 2023. URL <https://lmsys.org/blog/2023-06-29-longchat>.
- 619 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle  
620 Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before  
621 generation. *Advances in Neural Information Processing Systems*, 37:22947–22970, 2024.
- 622  
623 Enzhe Lu, Zhejun Jiang, Jingyuan Liu, Yulun Du, Tao Jiang, Chao Hong, Shaowei Liu, Weiran He,  
624 Enming Yuan, Yuzhi Wang, et al. Moba: Mixture of block attention for long-context llms. *arXiv*  
625 *preprint arXiv:2502.13189*, 2025.
- 626 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,  
627 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program  
628 synthesis. In *The Eleventh International Conference on Learning Representations*, 2023. URL  
629 [https://openreview.net/forum?id=iaYcJKpY2B\\_](https://openreview.net/forum?id=iaYcJKpY2B_).
- 630  
631 OpenAI. Introducing gpt-5 for developers, 2025. [https://openai.com/index/  
632 introducing-gpt-5-for-developers/](https://openai.com/index/introducing-gpt-5-for-developers/), Accessed: 2025-09-04.
- 633 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. YaRN: Efficient context win-  
634 dow extension of large language models. In *The Twelfth International Conference on Learning*  
635 *Representations*, 2024. URL <https://openreview.net/forum?id=wHBfxhZulu>.
- 636  
637 Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive  
638 transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
- 639 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
640 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.  
641 *arXiv preprint arXiv:1701.06538*, 2017.
- 642  
643 Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest: query-  
644 aware sparsity for efficient long-context llm inference. In *Proceedings of the 41st International*  
645 *Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- 646 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
647 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
tion processing systems*, 30, 2017.

- 648 Kuan-Chieh Wang, Daniil Ostashev, Yuwei Fang, Sergey Tulyakov, and Kfir Aberman. Moa:  
649 Mixture-of-attention for subject-context disentanglement in personalized image generation. In  
650 *SIGGRAPH Asia 2024 Conference Papers*, pp. 1–12, 2024a.
- 651  
652 Minzheng Wang, Longze Chen, Fu Cheng, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu,  
653 Nan Xu, Lei Zhang, Run Luo, Yunshui Li, Min Yang, Fei Huang, and Yongbin Li. Leave no  
654 document behind: Benchmarking long-context LLMs with extended multi-doc QA. In Yaser Al-  
655 Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Em-  
656 pirical Methods in Natural Language Processing*, pp. 5627–5646, Miami, Florida, USA, Novem-  
657 ber 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.322.  
658 URL <https://aclanthology.org/2024.emnlp-main.322/>.
- 659 Jian Wu, Linyi Yang, Dongyuan Li, Yuliang Ji, Manabu Okumura, and Yue Zhang. MMQA: Evalu-  
660 ating LLMs with multi-table multi-hop complex questions. In *The Thirteenth International Con-  
661 ference on Learning Representations*, 2025a. URL [https://openreview.net/forum?  
662 id=GG1pykXDCa](https://openreview.net/forum?id=GG1pykXDCa).
- 663 Shirley Wu, Michel Galley, Baolin Peng, Hao Cheng, Gavin Li, Yao Dou, Weixin Cai, James Zou,  
664 Jure Leskovec, and Jianfeng Gao. Collabllm: From passive responders to active collaborators. In  
665 *International Conference on Machine Learning (ICML)*, 2025b.
- 666  
667 Chaojun Xiao, Pengl Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan  
668 Liu, and Maosong Sun. Inllm: Training-free long-context extrapolation for llms with an efficient  
669 context memory. *Advances in Neural Information Processing Systems*, 37:119638–119661, 2024.
- 670 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming  
671 language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- 672  
673 Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, junxian guo, Shang Yang, Haotian Tang, Yao Fu,  
674 and Song Han. Duoattention: Efficient long-context LLM inference with retrieval and streaming  
675 heads. In *The Thirteenth International Conference on Learning Representations*, 2025. URL  
676 <https://openreview.net/forum?id=cFu7ze7xUm>.
- 677 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,  
678 and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question  
679 answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceed-  
680 ings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–  
681 2380, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.  
682 doi: 10.18653/v1/D18-1259. URL <https://aclanthology.org/D18-1259/>.
- 683 Zihao Ye, Lequn Chen, Ruihang Lai, Wuwei Lin, Yineng Zhang, Stephanie Wang, Tianqi Chen,  
684 Baris Kasikci, Vinod Grover, Arvind Krishnamurthy, and Luis Ceze. Flashinfer: Efficient and  
685 customizable attention engine for llm inference serving. *arXiv preprint arXiv:2501.01005*, 2025.  
686 URL <https://arxiv.org/abs/2501.01005>.
- 687  
688 Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie,  
689 Yuxing Wei, Lean Wang, Zhiping Xiao, Yuqing Wang, Chong Ruan, Ming Zhang, Wenfeng  
690 Liang, and Wangding Zeng. Native sparse attention: Hardware-aligned and natively trainable  
691 sparse attention. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher  
692 Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational  
693 Linguistics (Volume 1: Long Papers)*, pp. 23078–23097, Vienna, Austria, July 2025. Association  
694 for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.1126.  
695 URL <https://aclanthology.org/2025.acl-long.1126/>.
- 696  
697 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago  
698 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for  
699 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- 700  
701 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,  
Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient gen-  
erative inference of large language models. *Advances in Neural Information Processing Systems*,  
36:34661–34710, 2023.

## A LONGBENCH EVALUATIONS

We present detailed data of Figure 3 in Table 4 for reference. As shown in Table 4, both variants of Adamas (vanilla Adamas and Adamas with L2 distance) show consistent better performance over other sparse methods, especially under limited token budget (less than 512 tokens)

Table 4: Detailed LongBench evaluation results of StreamingLLM, SnapKV, Quest, Adamas with L2 and Adamas(Ours) on LongChat-7b-v1.5-32k.

Datasets	Methods / Budget	64	128	256	512	1024	2048	4096
GovReport (31.12)	StreamingLLM	1.21	7.13	14.52	18.86	21.72	24.08	27.01
	SnapKV	1.55	1.21	3.66	21.11	23.62	26.31	28.47
	Quest	3.3	11.93	22.84	27.27	29.89	<b>31.19</b>	<b>31.23</b>
	Adamas w/ L2	<b>30.09</b>	<b>30.45</b>	<b>30.86</b>	<b>30.93</b>	<b>30.96</b>	30.50	30.56
	Adamas(Ours)	<b>30.41</b>	<b>30.44</b>	<b>30.37</b>	<b>31.18</b>	<b>30.77</b>	<b>31.00</b>	<b>31.07</b>
HotpotQA (31.07)	StreamingLLM	3.63	9.68	14.51	16.25	18.62	21.83	24.14
	SnapKV	10.55	12.16	20.28	<b>35.65</b>	<b>33.90</b>	<b>35.23</b>	<b>33.94</b>
	Quest	5.31	15.49	21.54	26.64	30.29	32.21	<b>32.93</b>
	Adamas w/ L2	<b>32.58</b>	<b>32.91</b>	<b>31.15</b>	30.70	32.66	31.86	31.10
	Adamas(Ours)	<b>32.18</b>	<b>32.89</b>	<b>32.30</b>	<b>31.45</b>	<b>33.06</b>	<b>32.56</b>	31.41
MultifieldQA (41.64)	StreamingLLM	2.70	11.67	17.93	20.01	21.62	27.07	34.21
	SnapKV	7.53	6.24	16.74	40.98	37.76	40.03	<b>42.82</b>
	Quest	6.99	20.87	31.12	39.8	<b>42.03</b>	<b>44.09</b>	<b>43.25</b>
	Adamas w/ L2	<b>39.81</b>	<b>41.41</b>	<b>43.08</b>	<b>41.60</b>	41.28	<b>41.65</b>	42.33
	Adamas(Ours)	<b>36.60</b>	<b>40.67</b>	<b>42.49</b>	<b>42.93</b>	<b>42.22</b>	40.77	41.83
NarrativeQA (21.23)	StreamingLLM	0.92	4.47	9.17	10.61	12.46	16.85	17.40
	SnapKV	4.62	2.45	9.5	17.97	<b>19.47</b>	19.31	<b>20.32</b>
	Quest	1.86	5.98	14.2	16.28	17.91	19.66	19.88
	Adamas w/ L2	<b>16.06</b>	<b>18.39</b>	<b>18.37</b>	<b>19.20</b>	18.54	<b>19.74</b>	<b>20.33</b>
	Adamas(Ours)	<b>17.81</b>	<b>18.52</b>	<b>17.36</b>	<b>18.59</b>	<b>18.74</b>	<b>20.14</b>	20.22
Qasper (28.89)	StreamingLLM	3.15	8.72	10.87	11.74	14.79	17.50	25.33
	SnapKV	11.11	8.97	9.02	25.25	26.40	28.92	<b>29.25</b>
	Quest	8.67	18.54	26.72	30.89	<b>31.01</b>	<b>31.86</b>	<b>29.79</b>
	Adamas w/ L2	<b>28.70</b>	<b>31.36</b>	<b>30.35</b>	<b>31.23</b>	<b>30.80</b>	<b>29.36</b>	28.86
	Adamas(Ours)	<b>28.33</b>	<b>29.38</b>	<b>30.68</b>	<b>30.93</b>	29.94	29.26	28.65
TriviaQA (84.25)	StreamingLLM	8.73	33.14	51.20	60.37	68.48	75.67	80.21
	SnapKV	0.67	9.19	55.92	<b>84.01</b>	80.88	80.12	81.25
	Quest	12.29	44.68	68.92	81.32	<b>83.95</b>	<b>85.84</b>	<b>84.94</b>
	Adamas w/ L2	<b>80.97</b>	<b>83.35</b>	<b>84.21</b>	<b>84.29</b>	<b>83.60</b>	<b>84.03</b>	<b>84.25</b>
	Adamas(Ours)	<b>78.91</b>	<b>82.95</b>	<b>84.67</b>	83.99	83.36	83.75	83.95

## B END-TO-END EFFICIENCY EVALUATIONS

We further conduct end-to-end efficiency evaluations on an A800 80G PCIe GPU, which better reflects real-world server environments. As shown in Table 5, Adamas achieves up to a 1.4 $\times$  speedup on 32K-length sequences, demonstrating substantial and practically meaningful efficiency gains.

It is worth noting that Quest employs page-wise sparsity, trading accuracy for speed, especially under high-sparsity regimes, whereas Adamas consistently maintains high accuracy, achieving over 20% accuracy gains over Quest on average (Table 1). Despite operating at the token-wise level, Adamas

Table 5: End-to-end efficiency(s) evaluated on an A800 80G PCIe GPU, encompassing both prefill and 64-token decoding stages. Adamas achieves up to  $1.4\times$  speedup with high accuracy.

Sequence	Methods / Budget	256	512	1024	2048	4096	Full Attn
8192	Adamas(Ours)	1.833	1.863	1.844	1.751	1.785	
	Quest	1.626	1.606	1.573	1.596	1.621	2.343
	SnapKV	2.373	2.339	2.351	2.315	2.331	
16384	Adamas(Ours)	2.828	2.834	2.861	2.894	2.924	
	Quest	2.473	2.536	2.503	2.492	2.512	3.930
	SnapKV	3.851	3.934	3.901	3.938	3.914	
32768	Adamas(Ours)	5.680	5.711	5.703	5.730	5.804	
	Quest	4.901	4.908	4.918	4.945	4.984	7.978
	SnapKV	7.966	7.945	7.953	7.976	7.967	

is only about 10% slower than Quest, underscoring the effectiveness of our carefully optimized high-performance kernel.

### C BUCKETIZATION THRESHOLD ANALYSIS

We analyzed the distributions of Queries and Keys using the LongBench dataset, with the results presented in Figure 10. The distributions indicate that both Queries and Keys closely follow a zero-mean normal distribution.

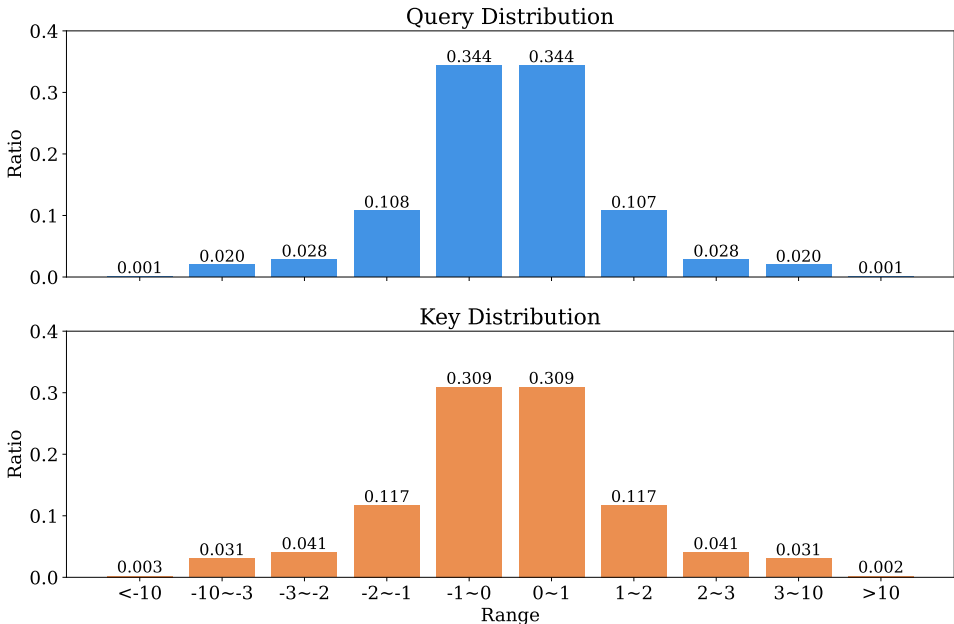


Figure 7: Distribution of Queries and Keys derived from LongBench statistics. Both Queries and Keys exhibit an approximately normal distribution.

To study the sensitivity of bucketization thresholds, we evaluate a range of values  $b \in \{1, 3, 5, 7, 10, 15\}$  under a fixed token budget of 128 using the bucketization rule  $[-b, 0, b]$ . The results in Table 6, evaluated on six LongBench datasets (with top-2 scores bolded), show that  $b = 1$  and  $b = 10$  consistently achieve the best accuracy across tasks.

Table 6: Sensitivity study of bucketization threshold, evaluated on six datasets of LongBench. Top-2 scores are bolded. Both threshold=1 and threshold=10 show consistent better accuracy performance.

Dataset / Threshold	1	3	5	7	10	15
GovReport	<b>31.06</b>	25.29	29.99	29.62	<b>30.44</b>	29.20
HotpotQA	<b>33.38</b>	26.12	30.81	29.00	<b>32.89</b>	31.54
MultifieldQA	40.28	33.54	37.75	<b>41.65</b>	<b>40.67</b>	39.35
NarrativeQA	<b>17.96</b>	16.47	15.45	16.02	<b>18.52</b>	14.90
Qasper	<b>28.95</b>	23.47	27.15	28.16	<b>29.38</b>	<b>29.38</b>
TriviaQA	<b>83.39</b>	74.7	80.46	81.04	<b>82.95</b>	80.65

### C.1 EMPIRICAL SENSITIVITY RESULTS

The experimental results reveal two threshold regimes that yield notably strong performance:

- **Fine-grained threshold:**  $b = 1$  (in implicit standardized units by LayerNorm)
- **Coarse-grained threshold:**  $b = 10$  (in raw magnitude)

Intermediate thresholds (e.g.,  $b = 3, 5, 7$ ) and those exceeding  $b = 10$  exhibit moderate degradation but remain functional, indicating robustness to suboptimal threshold selection.

### C.2 COMPREHENSIVE LONGBENCH COMPARISON FOR $b = 1$ AND $b = 10$

To further validate the effectiveness of the two best-performing thresholds, we conduct comprehensive evaluations of  $b = 1$  and  $b = 10$  across the LongBench benchmark. The results, presented in Figure 8, show that the performance curves of  $b = 1$  and  $b = 10$  almost completely overlap on every task. Moreover, both thresholds significantly outperform all other baselines, demonstrating that these two bucketization schemes consistently preserve the essential features required for accurate QK ranking under sparse attention.

### C.3 WHY $b = 1$ AND $b = 10$ PERFORM WELL

The effectiveness of these two thresholds can be explained by how they isolate dominant contributors to the QK dot product:

- **Large-magnitude components** ( $|x| > 10$ ) disproportionately influence the QK dot product and often encode strong semantic similarity.
- **Near-zero components** ( $|x| < 1$ ) contribute minimally and behave largely as noise.

Under this interpretation:

- $b = 1$  provides a principled separation of “strong positive,” “strong negative,” and “near-zero” signals in standardized units, enabling a clean distinction between informative and uninformative dimensions.
- $b = 10$  emphasizes dimensions with high raw magnitude that dominate the attention computation.

Both thresholds thus preserve the structural information necessary for reliably approximating QK rankings.

### C.4 WHY INTERMEDIATE OR LARGER THRESHOLDS UNDERPERFORM

Thresholds in the mid-range or in overly coarse regimes may impair bucketization by:

- **Over-partitioning** weak-signal regions, thereby amplifying noise, or
- **Over-merging** strong-signal regions, diminishing contrast among informative dimensions.

These effects lead to mild degradation in similarity estimation and consequently reduced accuracy.

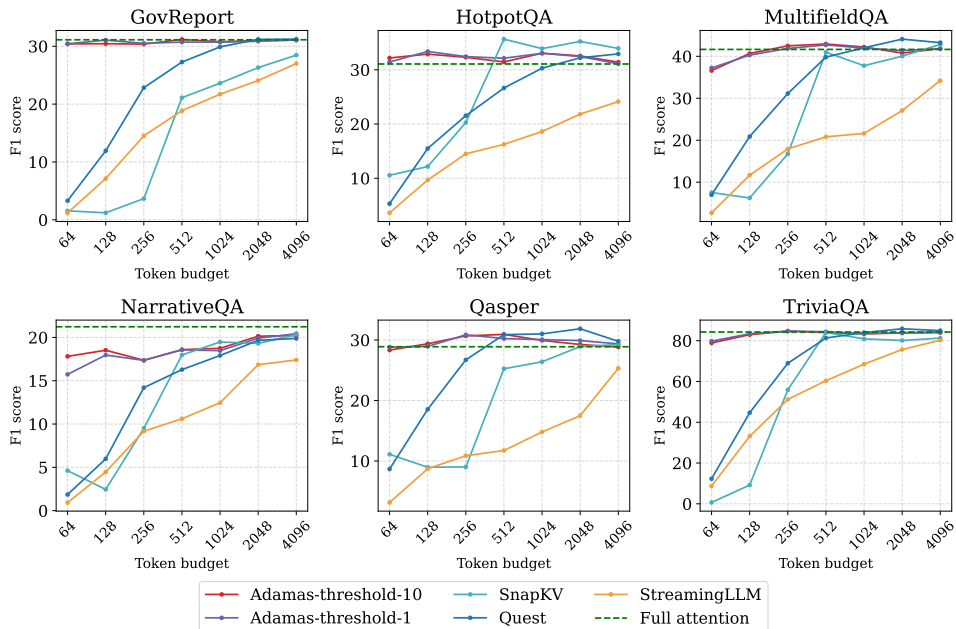


Figure 8: Full LongBench comparison of bucketization thresholds  $b = 1$  and  $b = 10$ . The accuracy curves nearly coincide and significantly outperform other baselines, validating the effectiveness of both threshold choices.

## C.5 CONCLUSION

Given that the dimensions of Queries and Keys approximately follow a normal distribution, an effective threshold need only separate strong-contribution features from weak or noisy ones. Thresholds that fail to achieve this separation exhibit moderate performance loss. In contrast,  $b = 1$  and  $b = 10$  provide interpretable and complementary partitioning strategies that retain the essential information needed for robust approximation of the QK ranking. These findings establish both thresholds as sound, empirically validated choices for bucketized attention computation.

## D PERPLEXITY EVALUATION OF ADAMAS WITH 4096-TOKEN BUDGET

For completeness, we additionally include the perplexity curves under a 4096-token budget, evaluated on PG19. As shown in the Figure 9, Adamas achieves noticeably lower perplexity than full attention, indicating that it can effectively filter out noisy tokens, select the most critical ones, and perform more accurate inference.

## E NEEDLE-IN-A-HAYSTACK STRESS TESTS

We conduct the Needle-in-a-Haystack stress tests using Yarn-Llama-2-7b-128k model, to evaluate the recall rates of Adamas, Quest, and full attention at a context length of 32K. As shown in Figure 10 and 11, Adamas maintains a recall performance that is nearly identical to that of full attention even under high sparsity (with only 256-token budget), consistently outperforming the baseline. Specifically, both Adamas and full attention achieve scores around 7/10 at a sequence length of 32K, whereas Quest reaches only a little over 2. This finding robustly supports our claim of "better recall under high sparsity".

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

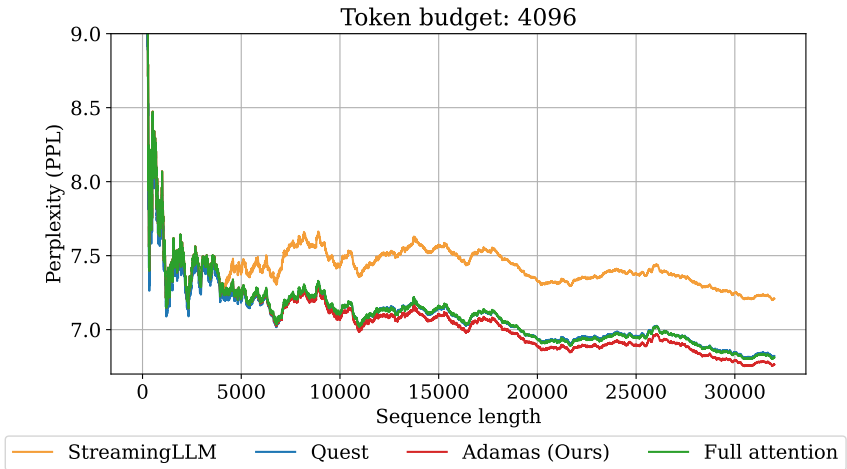


Figure 9: Perplexity comparison on PG19 under a 4096-token budget. Adamas achieves noticeably lower perplexity than full attention.

## F LLM USAGE

In this work, we used LLMs as auxiliary tools for writing and figure preparation. Specifically, LLMs were employed to refine the clarity of text, improve grammar, and correct minor wording errors. In addition, for some figures, initial plotting scripts were generated with the help of LLMs, which we subsequently modified and finalized. Importantly, all research ideas, theoretical results, algorithm designs, and experimental analyses were developed independently by the authors without reliance on LLMs.

## G COMPUTATIONAL WORKLOADS AND MEMORY ACCESS ANALYSIS

### G.1 MATHEMATICAL NOTATIONS

Table 7: Mathematical notations.

Symbol	Description.
$b$	Batch size
$s$	Sequence length
$h$	Hidden dimension
$h_{kv}$	Hidden dimension of Key and Value
$d$	Head dimension
$n$	Number of buckets
$p$	Page size
$k$	Number of selected tokens

### G.2 COMPUTATIONAL WORKLOAD

Detailed breakdown of computational workload is shown in Table 8. In summary:

- Computational workload of full attention:  $4(h + h_{kv} + s)bh$  FLOPs.
- Computational workload of Quest:  $[4h + 4h_{kv} + (\frac{4}{p} + \frac{4\log_2(\frac{k}{p})}{pd})s + 4k] \cdot bh + 2ph$  FLOPs.
- Computational workload of Adamas:  $[4h + 4h_{kv} + 2n + 2 + (\frac{4\log_2(k)}{d} + 3)s + 4k] \cdot bh$  FLOPs.

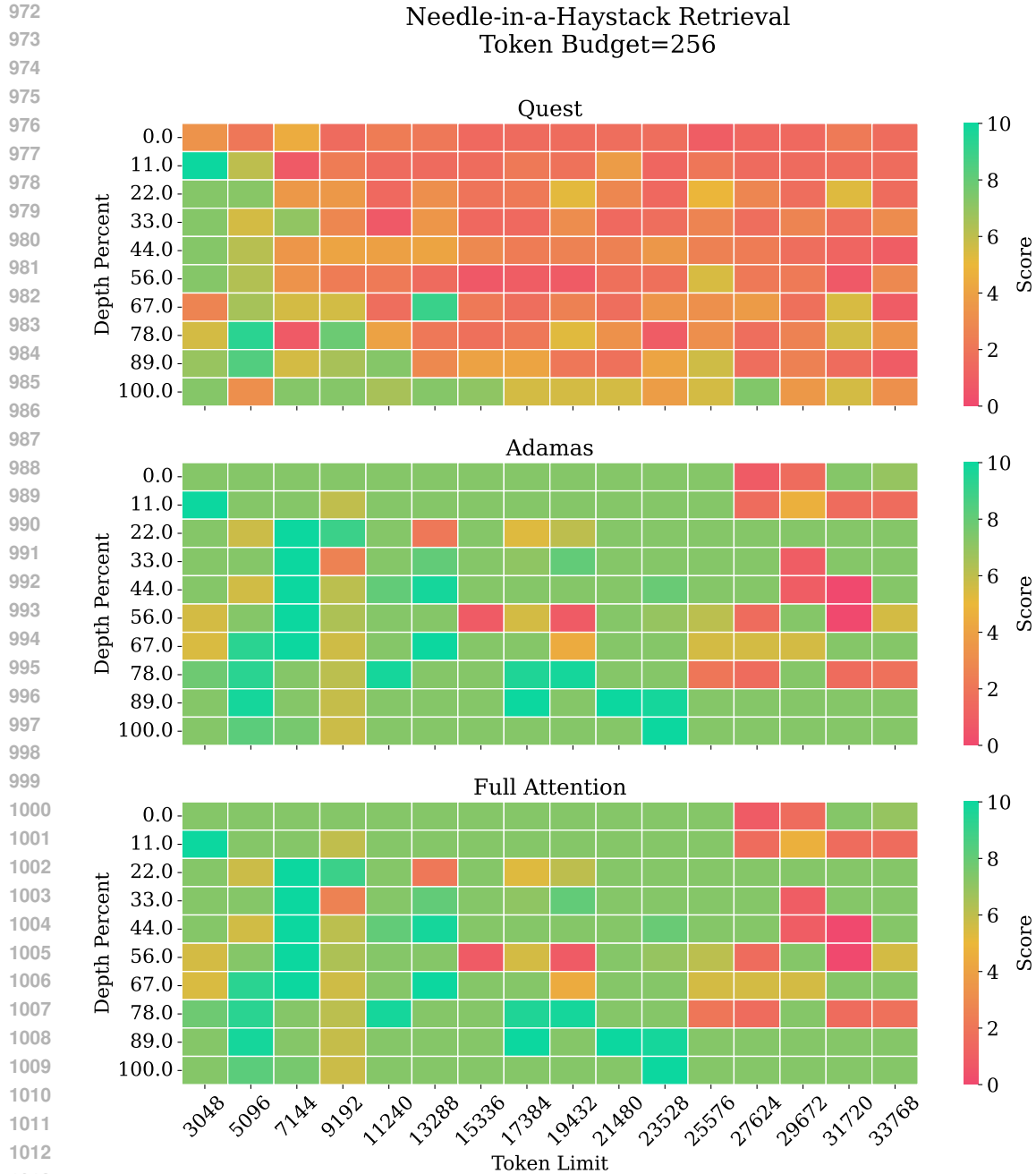


Figure 10: Needle-in-a-Haystack stress test using Yarn-Llama-2-7b-128k, with token budget=256.

### G.3 MEMORY ACCESS

Detailed breakdown of memory access is shown in Table 9. In summary:

- Memory access of full attention:
  - Read access:  $3bh + 2bsh_{kv} + 2h^2 + 2hh_{kv} + h + h_{kv}$
  - Write access:  $3bh + 2bh_{kv}$
- Memory access of Quest:

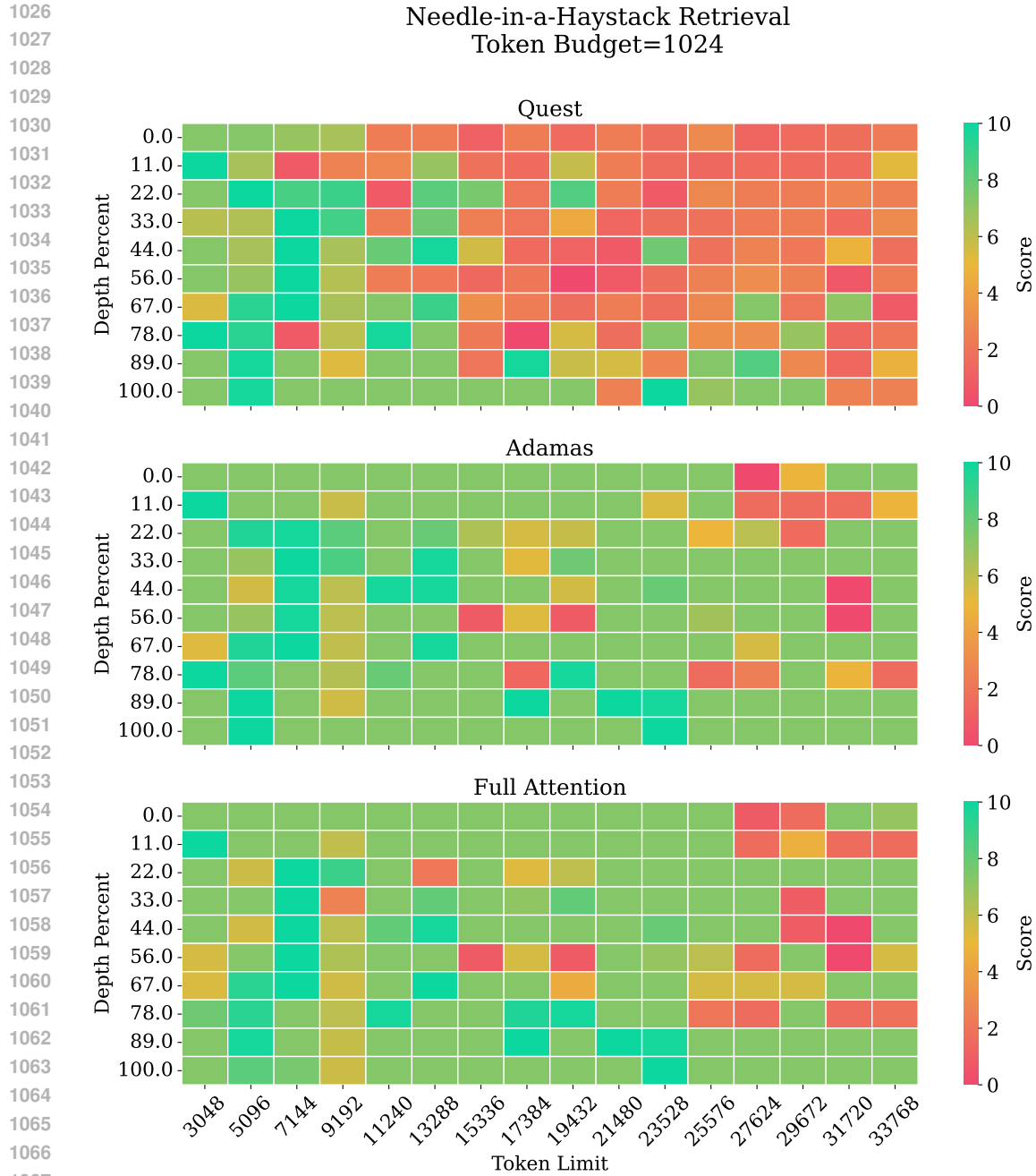


Figure 11: Needle-in-a-Haystack stress test using Yarn-Llama-2-7b-128k, with token budget=1024.

- Read access:  $7bh + \frac{2d+1}{pd} \cdot bsh + 2bkh_{kv} + 2h^2 + 2hh_{kv} + h + 2h_{kv}$
- Write access:  $(5 + \frac{s+k}{pd})bh + 2bh_{kv}$
- Memory access of Adamas:
  - Read access:  $\frac{57}{8}bh + \frac{bsh}{d} + 2bkh_{kv} + 2d^2 + 2h^2 + 2hh_{kv} + h + 2h_{kv}$
  - Write access:  $(\frac{21}{4} + \frac{s+k}{d})bh + 2bh_{kv}$

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Table 8: Computational workload

	Component	FLOPs
Full attention	QKV projection	$2bh(h + 2h_{kv})$
	$P = \text{softmax}(QK^\top)$	$2bsh$
	$PV$	$2bsh$
	Output projection	$2bh^2$
	<hr/>	
Quest	QKV projection	$2bh(h + 2h_{kv})$
	Reduce keys	$2ph$
	QK element-wise product	$\frac{2}{p} \cdot bsh$
	Per channel max	$\frac{1}{p} \cdot bsh$
	Page sum	$\frac{1}{p} \cdot bsh$
	Top- $k$	$\frac{4 \log_2(\frac{k}{p})}{pd} \cdot bsh$
	$P = \text{softmax}(QK^\top)$	$2bkh$
	$PV$	$2bkh$
	Output projection	$2bh^2$
	<hr/>	
Adamas	QKV projection	$2bh(h + 2h_{kv})$
	Hadamard transform	$2bh$
	Bucketizaion	$2n \cdot bh$
	Manhattan-distance estimation	$3bsh$
	Top- $k$	$\frac{4 \log_2(k)}{d} \cdot bsh$
	$P = \text{softmax}(QK^\top)$	$2bkh$
	$PV$	$2bkh$
Output projection	$2bh^2$	

Table 9: Memory access

	Component	Read access	Write access
Full attention	QKV projection	$bh + h(h + 2h_{kv}) + (h + 2h_{kv})$	$b(h + 2h_{kv})$
	$A = \text{softmax}(QK^\top)V$	$bh + 2bsh_{kv}$	$bh$
	Output projection	$bh + h^2$	$bh$
<hr/>			
Quest	QKV projection	$bh + h(h + 2h_{kv}) + (h + 2h_{kv})$	$b(h + 2h_{kv})$
	Reduce keys	$3bh$	$2ph$
	Criticality estimation	$\frac{2bsh}{p} + bh$	$\frac{bsh}{pd}$
	Top- $k$	$\frac{bsh}{pd}$	$\frac{bkh}{pd}$
	$A = \text{softmax}(QK^\top)V$	$bh + 2bsh_{kv}$	$bh$
	Output projection	$bh + h^2$	$bh$
<hr/>			
Adamas	QKV projection	$bh + h(h + 2h_{kv}) + (h + 2h_{kv})$	$b(h + 2h_{kv})$
	Hadamard transform	$2bh + 2hd^2$	$2bh$
	Bucketizaion	$2bh$	$\frac{bh}{4}$
	Manhattan-distance estimation	$\frac{bh+bsh}{8}$	$\frac{bsh}{d}$
	Top- $k$	$\frac{bsh}{d}$	$\frac{bkh}{d}$
	$A = \text{softmax}(QK^\top)V$	$bh + 2bsh_{kv}$	$bh$
	Output projection	$bh + h^2$	$bh + h^2$

## H ABLATION STUDIES

Below we present detailed ablation results evaluated on LongBench with LongChat-7b-v1.5-32k in Table H.

Table 10: Ablation results evaluated on LongBench with LongChat-7b-v1.5-32k.

Datasets	Methods / Budget	16	32	64	128	256	512	1024	2048	4096
GovReport (31.12)	Adamas (Ours)	22.08	28.31	<b>30.41</b>	30.44	30.37	<b>31.18</b>	30.77	<b>31.00</b>	31.07
	Adamas w/o Hadamard	0.60	0.59	0.59	0.64	2.06	6.32	14.12	23.53	29.14
	Adamas-1bit	17.17	23.67	27.88	30.09	30.72	31.16	<b>31.38</b>	30.84	30.70
	Adamas-3bit	<b>24.69</b>	<b>29.01</b>	30.05	<b>30.59</b>	30.83	31.01	30.95	30.71	<b>31.33</b>
	Adamas w/ L2 distance	18.27	27.19	30.09	30.45	<b>30.86</b>	30.93	30.96	30.50	30.56
HotpotQA (31.07)	Adamas (Ours)	24.05	<b>31.65</b>	32.18	32.89	32.30	31.45	33.06	<b>32.56</b>	31.41
	Adamas w/o Hadamard	1.01	1.04	0.68	0.48	0.77	2.07	5.38	13.47	25.13
	Adamas-1bit	19.22	24.70	29.04	31.55	<b>33.79</b>	31.63	32.14	32.47	31.58
	Adamas-3bit	<b>27.72</b>	31.17	31.29	31.82	32.56	<b>32.12</b>	<b>33.25</b>	32.21	<b>31.61</b>
	Adamas w/ L2 distance	18.73	28.83	<b>32.58</b>	<b>32.91</b>	31.15	30.70	32.66	31.86	31.10
MultifieldQA (41.64)	Adamas (Ours)	21.53	33.38	36.60	40.67	42.49	<b>42.93</b>	42.22	40.77	41.83
	Adamas w/o Hadamard	2.97	2.84	2.51	5.27	8.28	18.26	27.45	39.24	41.29
	Adamas-1bit	23.21	29.51	35.64	38.98	40.53	41.99	42.43	41.62	42.01
	Adamas-3bit	<b>30.04</b>	<b>36.00</b>	36.28	40.54	41.96	42.73	<b>43.14</b>	41.17	41.86
	Adamas w/ L2 distance	23.70	34.11	<b>39.81</b>	<b>41.41</b>	<b>43.08</b>	41.60	41.28	<b>41.65</b>	<b>42.33</b>
NarrativeQA (21.23)	Adamas (Ours)	11.98	15.13	<b>17.81</b>	<b>18.52</b>	17.36	18.59	18.74	20.14	20.22
	Adamas w/o Hadamard	1.68	1.00	1.12	0.91	0.38	1.34	3.29	6.00	13.32
	Adamas-1bit	5.38	9.16	14.79	15.18	17.34	<b>19.39</b>	<b>19.94</b>	20.00	20.01
	Adamas-3bit	<b>12.61</b>	<b>15.48</b>	17.65	16.73	17.91	19.05	19.03	<b>20.35</b>	20.09
	Adamas w/ L2 distance	8.83	14.50	16.06	18.39	<b>18.37</b>	19.20	18.54	19.74	<b>20.33</b>
Qasper (28.89)	Adamas (Ours)	18.83	24.52	28.33	29.38	<b>30.68</b>	30.93	29.94	29.26	28.65
	Adamas w/o Hadamard	1.73	2.40	2.29	4.93	9.24	18.27	26.17	<b>31.08</b>	<b>30.14</b>
	Adamas-1bit	18.04	20.81	24.16	28.73	29.50	30.04	30.18	30.52	28.73
	Adamas-3bit	21.95	26.07	26.39	30.23	30.52	31.12	29.50	29.51	28.81
	Adamas w/ L2 distance	<b>22.41</b>	<b>27.04</b>	<b>28.70</b>	<b>31.36</b>	30.35	<b>31.23</b>	<b>30.80</b>	29.36	28.86
TriviaQA (84.25)	Adamas (Ours)	56.77	75.13	78.91	82.95	<b>84.67</b>	83.99	83.36	83.75	83.95
	Adamas w/o Hadamard	2.63	2.08	1.21	3.52	9.73	20.14	33.15	56.62	79.13
	Adamas-1bit	52.26	70.62	80.08	80.88	82.59	83.90	<b>84.11</b>	<b>84.22</b>	83.63
	Adamas-3bit	<b>70.99</b>	<b>77.74</b>	<b>81.40</b>	82.95	83.99	83.79	83.46	83.55	83.68
	Adamas w/ L2 distance	37.82	66.80	80.97	<b>83.35</b>	84.21	<b>84.29</b>	83.60	84.03	<b>84.25</b>