



SHADOWKV: KV CACHE IN SHADOWS FOR HIGH-THROUGHPUT LONG-CONTEXT LLM INFERENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

With the widespread deployment of long-context large language models (LLMs), there has been a growing demand for efficient support of high-throughput inference. However, as the key-value (KV) cache expands with the sequence length, the increasing memory footprint and the need to access it for each token generation both result in low throughput when serving long-context LLMs. While various dynamic sparse attention methods have been proposed to speed up inference while maintaining generation quality, they either fail to sufficiently reduce GPU memory consumption or introduce significant decoding latency by offloading the KV cache to the CPU. We present SHADOWKV, a high-throughput long-context LLM inference system that stores the low-rank key cache and offloads the value cache to reduce the memory footprint for larger batch sizes and longer sequences. To minimize decoding latency, SHADOWKV employs an accurate KV selection strategy that reconstructs minimal sparse KV pairs on-the-fly. By evaluating SHADOWKV on a broad range of benchmarks, including RULER, LongBench, and Needle In A Haystack, and models like Llama-3.1-8B, Llama-3-8B-1M, GLM-4-9B-1M, Yi-9B-200K, Phi-3-Mini-128K, and Qwen2-7B-128K, we demonstrate that it can support up to $6\times$ larger batch sizes and boost throughput by up to $3.04\times$ on an A100 GPU without sacrificing accuracy, even surpassing the performance achievable with infinite batch size under the assumption of infinite GPU memory.

1 INTRODUCTION

Large language models (LLMs) have increasingly demonstrated their ability to scale and handle long contexts (Microsoft, 2024; Liu et al., 2024a; Achiam et al., 2023; Team et al., 2023), enabling them to tackle complex tasks like multi-document question answering and information retrieval from extensive contexts of up to 1M tokens (Achiam et al., 2023; Wang et al., 2024b). However, efficiently serving these long-context LLMs presents challenges related to the key-value (KV) cache (Liu et al., 2024b; Ge et al., 2023), which stores previous key-value activations to avoid re-computation. As the KV cache scales with sequence length, its growing memory footprint and the need to access it for each token generation lead to low throughput during long-context LLM inference. To address these issues, KV cache eviction or sparse attention methods have been widely explored.

However, existing methods face three primary limitations: accuracy degradation, inadequate memory reduction, and significant decoding latency overhead. KV cache eviction strategies (Zhang et al., 2024d;c) aim to reduce the memory footprint by discarding KV pairs based on specific policies, but they often result in information loss and accuracy degradation in tasks such as multi-turn conversations (Yang et al., 2024b; Tang et al., 2024a). Dynamic sparse attention methods (Tang et al., 2024b) preserve all KV pairs on the GPU and accelerate inference by computing attention with selected KV pairs. However, this line of work does not mitigate the memory footprint, thereby limiting the batch size and preventing accommodation of extremely long contexts (e.g., 1M tokens). A naive solution based on sparse attention involves offloading the KV cache to the CPU to reduce memory usage (Lee et al., 2024a; He & Zhai, 2024). Nonetheless, this approach incurs significant overhead due to the latency of fetching the selected sparse KV pairs from the CPU during decoding.

Consequently, an ideal effective system for long-context LLM inference with sparse attention should: (i) reduce GPU memory usage, (ii) minimize inference latency, and (iii) maintain accuracy

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

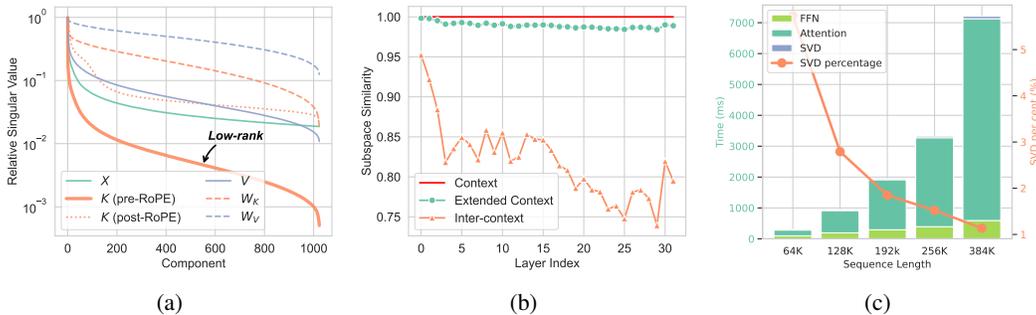


Figure 1: (a) For a sample from PG-19 (Rae et al., 2019; Gao et al., 2020) fed into Llama-3.1-8B, the pre-RoPE keys are the most low-rank, as indicated by the sharpest decay in singular values. (b) Average similarities, defined in Section 3.1, between rank-256 truncated SVD projections of pre-RoPE keys from PG-19 sequences using Llama-3.1-8B. Similarity is measured between a length 16K “Context” and either a 16K+2K continuation on “Context” (“Extended context”) or a new length 16K sequence (“Inter-context”). Pre-RoPE keys within sequences exhibit similar low-rank subspaces, while those between sequences show different patterns. (c) The relative overhead of singular value decomposition (SVD) decreases as sequence length scales for the pre-filling stage.

within limited sparse KV cache budgets. Fortunately, we can potentially overcome these challenges by leveraging our discovery that pre-Rotary Position Embedding (Su et al., 2024) (RoPE) keys are exceptionally low-rank compared to the layer inputs, post-RoPE keys, values, key weight matrix, and value weight matrix, as indicated in Figure 1a. Furthermore, our analysis in Figure 1b reveals that pre-RoPE keys lack significant similarities in low-rank subspaces across different sequences, while a sequence and its continuation tend to strongly share low-rank subspaces, enabling high compression rates within each sequence. Motivated by these findings, we have developed two key insights that pave the way for the design of an applicable system, detailed in Section 3.

Low-rank Keys and Offloaded Values for Storage: In long-context LLM inference, the quadratic scaling of attention computation with sequence length makes the linear cost of low-rank decomposition during pre-filling negligible, as illustrated in Figure 1c¹. To reduce memory footprint, we retain the low-rank pre-RoPE key cache on the GPU and offload the value cache to the CPU since the value cache does not exhibit low-rank properties, minimizing memory footprint without sacrificing accuracy. During decoding with sparse attention, we employ CUDA multi-streams to overlap the recovery of the selected key cache with the fetching of the corresponding value cache. This approach conceals key cache reconstruction and reduces data fetching overhead by $2\times$ compared to the naive offloading strategy, thereby decreasing the latency of sparse attention during decoding.

Accurate KV Selection for Fast Decoding: To further reduce decoding latency in sparse attention, we propose an accurate KV selection method that maintains accuracy with minimal number of selected tokens (i.e. the K of TopK), which we refer to as sparse budgets (1.56%). Our analysis reveals that most post-RoPE keys exhibit high cosine similarity with adjacent tokens, enabling chunk-level approximations for selecting important tokens. A minimal number of outlier chunks (0.3%), which are more challenging to approximate (Figure 3b), are stored as static cache on the GPU to preserve accuracy. As shown in Figure 2, our method outperforms the naive sparse attention approach (Tang et al., 2024b) and achieves higher sparsity, accelerating decoding.

Building on these insights, we present SHADOWKV in Section 4, depicted in Figure 2, a high-throughput system for long-context LLM inference. Specifically, during pre-filling, we offload the value cache to the CPU, retaining only the low-rank pre-RoPE keys, along with compressed landmarks of the key cache and detected outliers for larger batch sizes. During decoding, landmarks are used to select chunk indices for key cache recovery and value cache fetching. We perform accurate sparse attention computation with selected KV pairs and static outliers to achieve high throughput.

Empirically, we conduct extensive experiments and ablation studies to demonstrate the effectiveness and efficiency of SHADOWKV. In Section 5.1, we evaluate across various long-context LLMs, such

¹In practical scenarios, the key cache can be offloaded to the CPU to perform SVD asynchronously or precomputed and stored as part of the prefix cache (Juravsky et al., 2024).

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

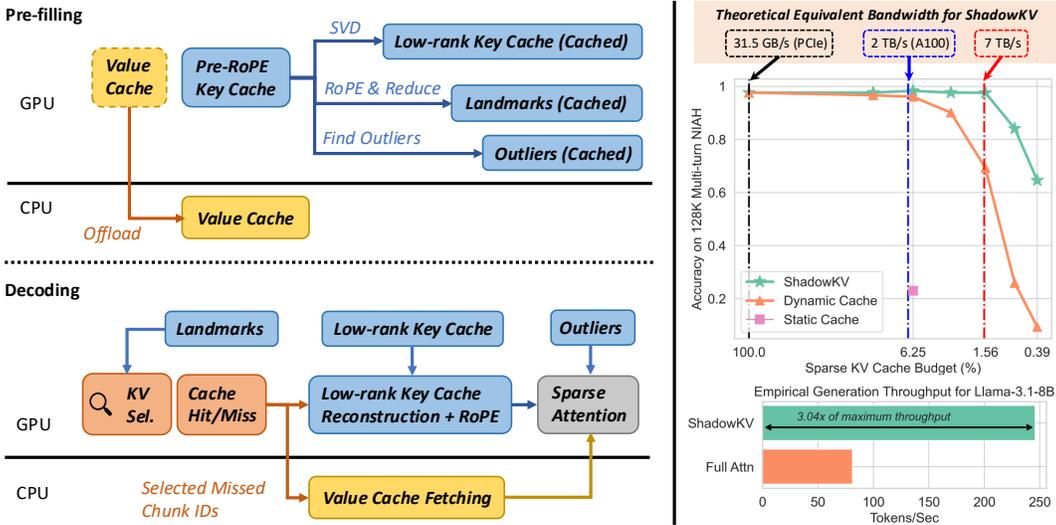


Figure 2: **Left:** SHADOWKV enhances long-context LLM inference throughput by offloading the value cache to the CPU while maintaining a low-rank key cache, landmarks, and outliers on the GPU. During decoding, it employs landmarks for efficient sparse attention, reducing computation and data movement. **Right:** SHADOWKV effectively utilizes a limited KV budget to achieve high accuracy, theoretically reaching over 7 TB/s equivalent bandwidth on an A100, and empirically boosts generation throughput by 3.04× for Llama-3.1-8B with on a batch of 122K contexts.

as Llama-3-8B-1M (Gradient., 2024), Llama-3.1-8B (Meta AI, 2024), GLM-4-9B-1M (GLM et al., 2024), Yi-9B-200K (AI et al., 2024), Phi-3-Mini-128K (Abdin et al., 2024) and Qwen2-7B-128K (Yang et al., 2024a) using benchmarks including RULER (Hsieh et al., 2024), LongBench (Bai et al., 2023), and Needle In A Haystack (Kamradt, 2023) with contexts up to 1M.

In Section 5.2, we demonstrate that SHADOWKV can support 6× larger batch sizes and boost throughput by 3.04× compared to small batches on an A100 using Llama-3.1-8B, with each sample having a context length of 122K. We also present results across different models and context lengths, increasing throughput up to 2.97× for Llama-3-8B-1M, 2.56× for GLM-4-9B-1M, and 2.66× for Yi-9B-200K, even surpassing infinite batch size under the assumption of infinite GPU memory.

2 RELATED WORKS

Token Eviction. To reduce memory footprint, eviction-based strategies keep a fixed size of KV cache to store the critical token KV pairs and discard unnecessary tokens. StreamingLLM (Xiao et al., 2023b) addresses the limitations of window attention by retaining attention sinks and recent KV pairs. H₂O (Zhang et al., 2024d) introduces a low-cost eviction policy, updating the KV cache based on cumulative attention scores. LESS (Dong et al., 2024b) accumulates evicted token information by a constant-sized low-rank cache, which allows partial access to previously evicted information, along with tokens maintained by a sparse policy. SnapKV (Li et al., 2024) uses the local window of prompts to select important tokens for future generations. However, they suffer from performance degradation and information loss since the evicted tokens will never be recovered.

Dynamic Sparse Attention. This line of work retains all KV cache but performs dynamic sparse attention within selected KV pairs to reduce inference latency. SparQ (Ribar et al., 2023) uses the norm of the query to decide an important subset of the key cache’s channels to calculate a metric to select relevant tokens. Quest (Tang et al., 2024b) segments tokens into pages and selects pages by approximating the highest attention within each page. Loki (Singhania et al., 2024) performs principal component analysis on key caches using a calibration dataset, selecting tokens based on attention scores computed in low-dimensional space. TriForce (Sun et al., 2024) combines sparse attention with speculative decoding (Leviathan et al., 2023) for lossless acceleration. InfiniGen (Lee et al., 2024a) offloads the entire KV cache to the CPU and prefetches essential entries using an predefined projection matrix via SVD for KV selection. In contrast, SHADOWKV employs an online, prompt-dependent SVD for key cache compression rather than for KV selection.

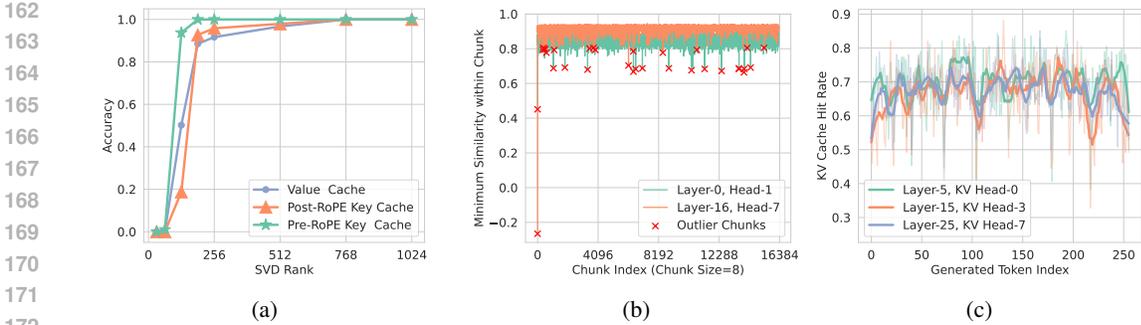


Figure 3: (a) Accuracy on the needle retrieval task across various ranks shows that the pre-RoPE key cache can be compressed by over 6 times without a drop in accuracy. (b) The number of notable outlier chunks is small, taking only 0.2-0.3%. (c) The KV cache has a high hit rate, reducing computations and data movements by over 60% for each decoding step.

Quantization. Several methods have been introduced to optimize KV cache quantization (Hooper et al., 2024; Yue et al., 2024; Xiao et al., 2023a), reducing memory consumption while retaining accuracy. KIVI (Liu et al., 2024c) applies different quantization strategies for keys and values, quantizing the keys per-channel and the values per-token to 2-bit. Palu (Chang et al., 2024) decomposes KV weight matrices offline, caching low-rank KV projections to achieve a higher compression rate. Quantization methods reduce the KV cache bit width, which is orthogonal to our approach.

3 OBSERVATIONS

We present two key insights of long-context LLMs that inspire SHADOWKV’s design, as follows.

3.1 LOW-RANK KEYS AND OFFLOADED VALUES FOR STORAGE

To reduce memory footprint, the low-rank nature of the KV cache has been explored by recent studies (DeepSeek-AI, 2024; Xu et al., 2024; Chang et al., 2024). However, these methods focus on data-independent decomposition, either requiring training or achieving limited compression rates.

Observation. In our study, by conducting SVD on the model weights W_k , W_v , the input X , the pre-/post-RoPE key cache, and the value cache of Llama-3.1-8B, we visualize the relative singular value distributions in Figure 1a together with the accuracy in Figure 3a. As we observed, pre-RoPE keys have the lowest rank and can be compressed by $6\times$ without performance degradation.

We also identify striking dynamic and static behaviors in low-rank keys between and within sequences, inspired by a related investigation in FFN layers (Dong et al., 2024a). Analogous to cosine similarity, we define $\mathcal{D}(\mathbf{H}_1, \mathbf{H}_2) = \langle \mathbf{H}_1, \mathbf{H}_2 \rangle / r$ to be the similarity metric between low-rank subspaces of two rank- r projection matrices, \mathbf{H}_1 and \mathbf{H}_2 , where $\langle \cdot, \cdot \rangle$ is the Frobenius inner product². In our case with truncated SVDs of pre-RoPE keys, let $\mathbf{K}_1, \mathbf{K}_2 \in \mathbb{R}^{n \times d}$ have rank- r truncated SVDs, $\Phi_1 \Sigma_1 \Psi_1^\top$ and $\Phi_2 \Sigma_2 \Psi_2^\top$, respectively, where $\Phi_1 \in \mathbb{R}^{n \times r}$, $\Sigma_1 \in \mathbb{R}^{r \times r}$, $\Psi_1 \in \mathbb{R}^{d \times r}$, and similarly for Φ_2, Σ_2 , and Ψ_2 . Then, $\mathcal{D}(\Psi_1 \Psi_1^\top, \Psi_2 \Psi_2^\top)$ can measure the similarity between the low-rank subspaces of the two right singular matrices. Depicted in Figure 1b, pre-RoPE keys between sequences do not strongly share similar low-rank subspaces, but extensions of the same sequence do.

Insights. Our observation of the low-rank nature in the pre-RoPE keys indicates that storing the low-rank projections is sufficient for each sequence. By keeping the low-rank key cache on the GPU and offloading the value cache to the CPU since it is not low-rank, we can largely reduce the memory footprint. During decoding, selected KV pairs can be reconstructed on-the-fly for computation.

²Since \mathbf{H}_1 and \mathbf{H}_2 are projection matrices, their squared Frobenius norms are the sum of their singular values which consist of r 1’s and $d - r$ 0’s, i.e., $\|\mathbf{H}_1\|_F^2 = r$. Thus, by Cauchy-Schwarz, $|\mathcal{D}(\mathbf{H}_1, \mathbf{H}_2)| \leq 1$. Additionally, $\mathcal{D}(\mathbf{H}_1, \mathbf{H}_2) \geq 0$ by the cyclic property of trace and positive semidefiniteness of projection matrices. Together, this shows $\mathcal{D}(\mathbf{H}_1, \mathbf{H}_2) \in [0, 1]$, maximized or minimized when the projection matrices project onto identical or orthogonal subspaces, respectively.

Algorithm 1: SHADOWKV Pre-filling

Input: $K, K^{\text{RoPE}}, V \in \mathbb{R}^{b \times h_{kv} \times s \times d}$, SVD rank r , chunk size c , number of outlier chunks o

\triangleright Store low-rank projection of pre-RoPE key cache
 $A \in \mathbb{R}^{b \times s \times r}, B \in \mathbb{R}^{b \times h_{kv} \times r \times d} \leftarrow \text{SVD}(K)$

\triangleright Segment post-RoPE key cache into chunks and compute the mean of each chunk
 $C \in \mathbb{R}^{b \times h_{kv} \times s / c \times d} \leftarrow \text{Reduce}(K^{\text{RoPE}})$

\triangleright Compute cosine similarity within each chunk
 $S \in \mathbb{R}^{b \times h_{kv} \times s / c \times c} \leftarrow \text{CosineSimilarity}(C, K^{\text{RoPE}})$

\triangleright Find lowest cosine similarity as outliers
 $I \in \mathbb{R}^{b \times h_{kv} \times o} \leftarrow \text{ArgTopK}(-\text{Min}(S, \text{dim} = -1), o)$

$K^{\text{outlier}}, V^{\text{outlier}} \leftarrow \text{Gather}(K^{\text{RoPE}}, V, I)$

\triangleright Offload the rest of values to the CPU and store the non-outlier chunks' mean as landmarks

$V^{\text{CPU}} \leftarrow V \setminus V^{\text{outlier}}, L \leftarrow C \setminus \text{Gather}(C, I)$

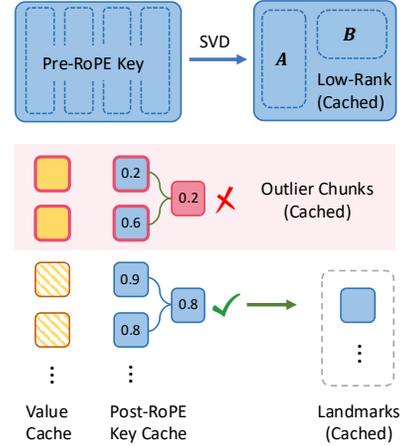


Figure 4: SHADOWKV pre-filling.

3.2 ACCURATE KV SELECTION FOR FAST DECODING

To further reduce the latency overhead in sparse attention, including fetching the selected value cache from the CPU and reconstructing the corresponding key cache, an accurate KV selection method is needed to minimize the sparse KV cache budget while maintaining the accuracy.

Observation. We found most post-RoPE key cache exhibits spatial locality, with high cosine similarity to adjacent tokens, except for a few outliers. To quantify this, we conducted inference experiments on 128K contexts. We divided the post-RoPE keys into chunks of eight tokens and visualized the minimum cosine similarity between the chunk’s mean and its key cache, as shown in Figure 3b. The results indicate that, apart from a few outliers, there is generally high cosine similarity, suggesting the mean values can serve as landmarks to approximate attention well within normal chunks.

Analysis. This finding suggests that for the majority of chunks, we can maintain the mean value as compressed landmarks to select minimal important KV pairs (1.56%) accurately during decoding. Outlier chunks, which may contain dense or critical information and are difficult to approximate, are retained to ensure accuracy. Given their relatively small number (0.2–0.3%), storing them on the GPU is feasible without affecting memory capacity. Furthermore, as shown in Figure 3c, considering the temporal locality of the KV cache—meaning that the KV pairs selected by the queries of two adjacent decoding steps have a high repetition rate, a cache policy (Zhang et al., 2024a) can be leveraged to further reduce the latency overhead by 60% during decoding with optimized kernels.

4 SHADOWKV

In this section, we introduce SHADOWKV, a high-throughput long-context LLM inference system. We first elaborate our algorithm in Section 4.1, covering both the pre-filling and decoding phases. Subsequently, in Section 4.2, we discuss the concept of theoretical equivalent bandwidth to illustrate the benefits of our approach.

4.1 ALGORITHM

The algorithm of SHADOWKV is divided into two main phases: pre-filling and decoding. The pre-filling phase involves low-rank decomposition of the post-RoPE key cache, offloading the value cache, and constructing landmarks to facilitate subsequent high-throughput decoding. The decoding phase includes accurate KV selection and efficient sparse KV cache reconstruction.

Pre-filling. During the pre-filling phase, we optimize GPU memory usage by performing low-rank compression on the key cache of each layer and offloading values to the CPU. Specifically, as demonstrated in Algorithm 1 and Figure 4, we apply SVD on the pre-RoPE key cache and store only

Algorithm 2: SHADOWKV Decoding

Input: $A, B, L, V^{\text{CPU}}, Q \in \mathbb{R}^{b \times h_q \times s_q \times d}$, $K^{\text{outlier}}, V^{\text{outlier}}, K, V \in \mathbb{R}^{b \times h_{kv} \times s_q \times d}$, number of chunks n_c , number of selected chunk budget k

▷ *Compute chunk attention score*

$$P \in \mathbb{R}^{b \times h_q \times s_q \times n_c} \leftarrow \text{MatMul}(Q, L^T)$$

$$S \in \mathbb{R}^{b \times h_q \times s_q \times n_c} \leftarrow \text{Softmax}(P/\sqrt{d})$$

$$S_1 \in \mathbb{R}^{b \times h_q \times n_c} \leftarrow \text{sum}(S, \text{dim} = -2)$$

$$S_2 \in \mathbb{R}^{b \times h_{kv} \times n_c} \leftarrow \max_{\text{kv_group}}(S_1)$$

▷ *Select top-k chunks for each KV head*

$$I \in \mathbb{R}^{b \times h_{kv} \times k} \leftarrow \text{ArgTopK}(S_2, k)$$

▷ *Gather values from CPU*

$$V^{\text{sparse}} \leftarrow \text{Gather}(V^{\text{CPU}}, I)$$

$$V \leftarrow [V^{\text{outlier}}, V^{\text{sparse}}, V]$$

▷ *Recover keys from low-rank projection*

$$K^{\text{sparse}} \leftarrow \text{MatMul}(\text{Gather}(A, I), B)$$

$$K \leftarrow [K^{\text{outlier}}, \text{RoPE}(K^{\text{sparse}}); K]$$

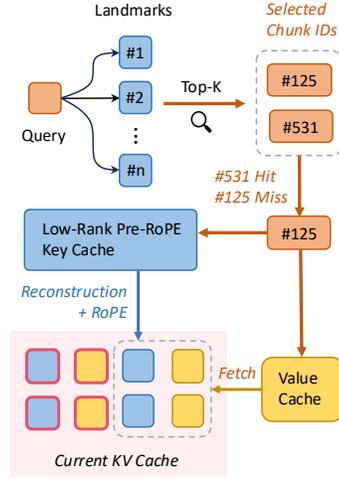


Figure 5: SHADOWKV decoding phase.

the low-rank representations for each layer. Post-RoPE key cache is segmented into chunks, with the mean of each chunk computed as landmarks. By computing the cosine similarity within these chunks, we identify poorly approximated tokens as outliers. This small set of outliers is gathered and stored on the GPU as the static cache, while the remaining key cache is maintained as compact landmarks, with the corresponding values offloaded to the CPU memory.

High-throughput Decoding. For incoming queries, we first compute the approximate attention scores using the landmarks. As detailed in Algorithm 2, by identifying the top-k scoring chunk indices, the corresponding values are retrieved from the CPU, and the key cache is simultaneously reconstructed from low-rank projections, effectively concealing the construction of the key cache. Based on the insight that the KV cache has temporal locality, we build cache-aware CUDA kernels, reducing computation and value fetching by 60%. As shown in Figure 5, we conduct an index scan to detect the missed chunks and only rebuild the necessary KV pairs on-the-fly.

Based on our observations in Section 3.1, future pre-RoPE keys within a sequence reside in a shared low-rank subspace with the context. As a result, an extension of our algorithm would be to store generated tokens as low-rank states using the same projections obtained from pre-filling to reduce the memory usage for future generations³. We evaluate it and include the results in Appendix A.1.

4.2 THEORETICAL EQUIVALENT BANDWIDTH

The benefit of SHADOWKV in terms of increasing throughput can be analyzed through the concept of equivalent bandwidth. Consider each K or V vector as being M bytes in size, with a sequence length of S , a chunk size of C , a selected chunk budget of K , O outliers, and hit rate α . During KV selection, SHADOWKV loads $M \times S/C$ bytes using the GPU memory bandwidth B_{GPU} . For value cache fetching, it loads $M \times K \times C$ bytes using the PCIe bandwidth B_{PCIe} (Sheng et al., 2023). Since value movement and key cache reconstruction can be overlapped, we do not need to count key cache reconstruction here. Following this, SHADOWKV performs standard attention computation for the top-k chunks and predefined outliers, requiring $2M \times (K + O) \times C$ bytes. The equivalent bandwidth of SHADOWKV is defined as below and the GPU memory savings is detailed in Appendix A.2.

$$B_{\text{equivalent}} = \frac{2SB_{\text{GPU}}}{S/C + 2(K + O)C + (1 - \alpha)KCB_{\text{GPU}}/B_{\text{PCIe}}}$$

For example, assuming $C=8$, $S=128\text{K}$, $K=256$, $O=48$, $B_{\text{PCIe}}=31.5$ GB/s, and $B_{\text{GPU}}=2$ TB/s for A100, the equivalent bandwidth of SHADOWKV is calculated as 7.2 TB/s, which is $3.6 \times$ higher than A100 memory bandwidth. This result indicates that SHADOWKV theoretically achieves a high equivalent bandwidth to accelerate attention computation. System implementation is detailed in Appendix B.1.

³If $\Psi \in \mathbb{R}^{d \times r}$ is the right singular matrix calculated from the SVD of pre-RoPE context keys $K \in \mathbb{R}^{s \times d}$, new pre-RoPE keys $K' \in \mathbb{R}^{s_q \times d}$ can be stored as $K'\Psi$ and projected back up with Ψ^T when needed.

Table 1: Performance of different models and different methods on RULER (Hsieh et al., 2024) evaluated at length of 128K. SHADOWKV outperforms other methods with a 1.56% sparse budget.

Methods	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
<i>Llama-3-8B-1M</i>	100.00	100.00	98.96	98.96	98.96	95.57	75.00	48.96	78.54	71.85	86.68
Loki	18.75	1.04	2.08	0.00	1.56	0.78	4.17	13.54	26.04	25.35	9.33
Loki (V only)	41.67	6.25	37.50	1.04	8.07	30.73	10.42	19.79	51.67	37.50	24.46
InfiniGen	100.00	98.96	84.38	53.13	63.28	54.95	65.63	48.96	81.67	50.35	70.13
InfiniGen (V only)	100.00	98.96	96.88	76.04	81.25	77.08	67.71	50.00	81.67	53.47	78.31
Quest	100.00	100.00	98.96	77.08	97.65	93.49	60.42	50.00	77.08	65.63	82.03
Quest (V only)	100.00	100.00	98.96	85.42	97.92	95.49	70.83	46.88	78.75	65.63	83.99
SHADOWKV	100.00	100.00	97.92	98.96	96.88	95.83	72.92	52.08	81.67	72.57	86.88
<i>GLM-4-9B-1M</i>	100.00	100.00	94.79	87.50	99.74	93.75	67.71	55.21	97.29	72.22	86.82
Loki	71.88	27.08	22.92	2.08	9.90	11.46	28.13	27.08	31.04	54.17	28.57
Loki (V only)	96.88	55.21	56.25	18.75	51.04	50.52	45.83	39.58	72.71	59.72	54.65
InfiniGen	100.00	93.75	82.29	0.00	79.43	60.16	57.29	53.13	92.71	57.29	67.60
InfiniGen (V only)	100.00	96.88	87.50	7.29	95.31	75.26	56.25	54.17	95.63	60.76	72.91
Quest	100.00	95.83	90.62	54.17	94.01	76.30	55.21	52.08	95.83	64.58	77.86
Quest (V only)	100.00	96.88	93.75	72.92	95.83	83.07	56.25	53.13	96.88	65.97	81.47
SHADOWKV	100.00	100.00	95.83	83.33	98.70	87.76	69.79	55.21	97.50	68.06	85.62
<i>Llama-3.1-8B</i>	100.00	100.00	98.96	91.67	98.96	95.31	82.29	47.92	68.96	71.18	85.53
Loki	68.75	32.29	32.29	20.83	42.71	28.65	41.67	33.33	24.79	29.86	35.52
Loki (V only)	95.83	36.46	57.29	62.50	77.86	70.83	69.79	39.58	35.21	37.50	58.29
InfiniGen	100.00	77.08	78.13	13.54	58.07	47.40	65.63	41.67	60.83	50.35	59.27
InfiniGen (V only)	100.00	88.54	87.50	26.04	79.43	77.08	72.92	43.75	57.08	55.21	68.76
Quest	100.00	98.96	97.92	34.38	93.49	88.54	70.83	44.79	65.63	68.40	76.29
Quest (V only)	100.00	98.96	98.96	56.25	95.83	90.63	76.04	46.88	66.25	67.36	79.72
SHADOWKV	100.00	100.00	100.00	83.33	97.92	92.19	81.25	48.96	67.08	64.93	83.57

5 EMPIRICAL EVALUATION

In this section, we showcase the effectiveness and efficiency of SHADOWKV. Specifically,

- In Section 5.1, we show that SHADOWKV can reduce the GPU memory footprint of the KV cache by over $6\times$ without accuracy degradation on a wide range of models and evaluation benchmarks.
- In Section 5.2, we demonstrate SHADOWKV can support up to $6\times$ larger batch sizes and increase the inference throughput by up to $3.04\times$ without compromising model quality.
- In Section 5.3, we present extensive ablation studies that validate the effectiveness of each component of SHADOWKV in optimizing GPU memory usage and enhancing performance.

5.1 ACCURACY EVALUATION

We demonstrate that SHADOWKV can reduce the GPU memory usage of the KV cache by $6\times$ while maintaining accuracy on a range of long-context tasks with a minimal sparse KV cache budget.

Setup. We choose four widely used long-context models for our evaluation: Llama-3-8B-1M (Gradient., 2024), GLM-4-9B-1M (GLM et al., 2024), Llama-3.1-8B (Meta AI, 2024), and Yi-9B-200K (AI et al., 2024). We evaluate our approach on three challenging long-context benchmarks: RULER (Hsieh et al., 2024), LongBench (Bai et al., 2023), and Needle In A Haystack (Kamradt, 2023), covering QA, multi-hop, reasoning, summarization, code completion⁴. We set the chunk size to 8, the rank to 160, and the number of outliers to 48 for SHADOWKV.

Baselines. We include three dynamic sparse attention methods as baselines: Quest (Tang et al., 2024b), Loki (Singhania et al., 2024), and InfiniGen (Lee et al., 2024b). For all methods, we retain exact pre-filling and perform dynamic sparse attention during decoding, where the computation cost is set to 1/16 of full attention for selecting sparse KV pairs. We include two variants for each baseline: one where all the KV cache is offloaded, and another where only the value cache is offloaded.

⁴We include results for Yi-9B-200K and other models (e.g., Llama-3-70B-1M) in Appendix A. Needle In A Haystack is also tested on Phi-3-Mini-128K (Abdin et al., 2024) and Qwen2-7B-128K (Yang et al., 2024a).

Table 2: Performance of various methods on different models with LongBench (Bai et al., 2023) samples exceeding 4K tokens. SHADOWKV outperforms other methods and maintains the accuracy.

Methods	NarrQA	MultiFQA	HotpotQA	MuSiQue	DuRead	GovRep	SAMSum	PassRetr	LCC	Avg.
<i>Llama-3-8B-1M</i>	18.98	41.84	36.79	21.47	31.93	34.18	35.96	81.50	56.07	39.86
Loki	2.26	10.19	5.48	3.16	12.17	28.97	7.84	40.52	31.44	15.78
Loki (V only)	3.20	21.01	12.41	3.86	17.07	31.24	16.23	52.57	38.10	21.74
InfiniGen	14.39	31.46	33.63	17.94	26.65	27.38	21.97	74.30	38.58	31.81
InfiniGen (V only)	17.83	36.08	35.28	19.64	28.39	29.28	28.12	74.85	45.53	35.00
Quest	20.13	36.63	35.00	18.14	24.55	27.11	35.63	79.00	53.64	36.65
Quest (V only)	17.26	39.51	36.78	18.71	26.41	29.49	35.80	79.50	60.05	38.17
SHADOWKV	17.17	39.73	38.29	21.08	31.77	31.62	35.87	80.00	63.93	39.94
<i>GLM-4-9B-1M</i>	25.44	51.09	58.67	39.61	32.04	29.97	40.31	99.00	58.02	48.24
Loki	5.82	30.60	22.73	9.20	30.09	30.35	22.70	98.92	40.77	32.35
Loki (V only)	10.89	44.97	45.44	23.51	32.07	30.56	35.34	99.50	50.27	41.39
InfiniGen	23.67	46.31	55.69	33.91	27.49	25.44	33.48	91.83	36.96	41.64
InfiniGen (V only)	25.63	48.44	57.23	36.94	29.77	26.67	36.64	93.58	46.69	44.62
Quest	23.81	44.53	56.41	35.49	23.54	21.73	37.39	87.00	43.80	41.52
Quest (V only)	26.00	46.32	57.54	36.42	24.58	24.52	37.71	93.50	46.52	43.68
SHADOWKV	26.50	51.31	59.09	38.87	32.92	28.54	38.70	96.50	58.55	47.89
<i>Llama-3.1-8B</i>	31.56	55.10	57.65	29.46	35.26	34.45	29.84	100.00	67.31	48.96
Loki	2.31	18.89	10.64	5.47	19.30	31.16	15.91	94.88	44.60	27.02
Loki (V only)	3.93	38.59	22.85	12.96	27.43	32.22	26.43	98.25	56.11	35.42
InfiniGen	27.23	52.72	53.89	26.81	27.72	29.61	24.42	98.93	56.89	44.25
InfiniGen (V only)	29.73	53.47	55.11	28.72	28.55	31.42	26.76	99.17	62.66	46.18
Quest	29.70	49.04	53.96	27.18	27.16	30.43	29.85	98.50	57.35	44.80
Quest (V only)	30.02	53.97	56.39	27.06	29.06	31.65	30.23	99.00	63.89	46.81
SHADOWKV	30.93	55.20	57.32	29.13	31.85	32.79	30.40	99.50	66.03	48.13

The former has similar latency to SHADOWKV but a smaller sparse budget since SHADOWKV only needs to fetch the value cache from the CPU. The latter aligns with the same sparse KV cache budget but significantly increases GPU memory usage. The latter one is marked as “V only” in the table.

RULER. As shown in Table 1, SHADOWKV demonstrates excellent performance on 128K contexts. With a fixed sparse budget of 1.56%, other methods experience performance degradation. In contrast, SHADOWKV is more robust and even outperforms original full attention on certain tasks, such as variable tracking. For complex tasks like multi-document QA or multi-key needle retrieval, other methods suffer from significant performance degradation while SHADOWKV does not.

LongBench. On LongBench, we evaluate our method with a range of realistic scenarios, including single-/multi-document question-answering, document summarization, code completion, information retrieval, etc. We only test on samples longer than 4K and set the sparse KV cache budget to 256 for this benchmark since it has shorter inputs compared to RULER. As shown in Table 2, SHADOWKV outperforms other methods consistently and maintains the performance.

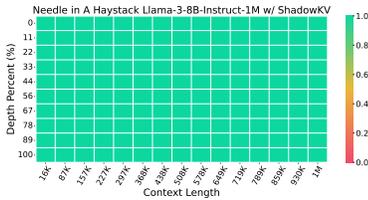


Figure 6: Needle In A Haystack.

Needle In A Haystack. On the Needle In A Haystack dataset, as shown in Figure 6, SHADOWKV shows the ability to process information at different positions across various context windows, ranging from 16K to 1M tokens. More experiments on a range of models can be found in Appendix B.3.

Integrate with Efficient Pre-filling Methods. We also combined SHADOWKV with a state-of-the-art efficient pre-filling method MInference (Jiang et al., 2024). As shown in Table 3, following the setting of MInference, we tested it on RULER with contexts scaling from 8K to 256K. This demonstrates that our method is compatible with pre-filling acceleration techniques. For some certain context length settings, we even see a slight performance improvement.

Table 3: Performance of different methods on RULER (Hsieh et al., 2024) using MInference (Jiang et al., 2024) in the pre-filling stage. SHADOWKV is compatible with MInference.

Methods	8K	16K	32K	64K	128K	256K	Avg.
Llama-3-8B-1M w/ MInference	89.92	88.02	82.81	78.45	78.12	74.57	81.98
SHADOWKV w/ MInference	90.47	88.12	83.28	77.71	78.32	74.31	82.04

Multi-turn Conversation Capability. To simulate multi-turn conversations, we challenged SHADOWKV with a multi-turn needle retrieval task (Multi-turn NIAH). We also test two eviction-based methods in Figure 7, including SnapKV (Li et al., 2024) and StreamingLLM (Xiao et al., 2023b). The performance of SnapKV drops significantly from the second round due to the required context information being different from the first round. Since SnapKV inevitably evicted tokens based on the first-turn conversation, it cannot successfully retrieve related information for future queries. In contrast, SHADOWKV can maintain accuracy in the multi-turn conversation setting.

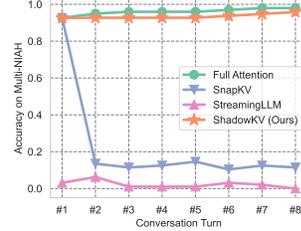


Figure 7: Multi-turn NIAH.

5.2 EFFICIENCY EVALUATION

To demonstrate the efficiency of SHADOWKV, we deploy it into real-world large batch serving scenarios. By measuring the throughput during decoding across different models on A100, we show that SHADOWKV can support up to $6\times$ larger batch sizes and boost throughput by up to $3.04\times$. [The detailed latency breakdown can be found in Appendix A.6.](#)

Baselines. The baseline selects the largest batch size that can fit entirely on the GPU with full attention. We also include results for the same batch size of SHADOWKV and the infinite batch size, assuming infinite GPU memory capabilities⁵. We set the sparse budget to 1.56% for SHADOWKV.

Results. As shown in Table 4, SHADOWKV demonstrates significant throughput improvements for various models on an A100, surpassing even those with infinite GPU memory. Notably, SHADOWKV supports batch sizes up to $6\times$ larger and enhances throughput by up to $3.04\times$ compared to full attention, even surpassing infinite batch size assuming infinite GPU memory. While the gains for GLM-4-9B-1M and Yi-9B-200K are slightly lower, the improvements still reach up to $2.56\times$ and $2.66\times$ respectively, highlighting SHADOWKV’s adaptability even with fewer KV heads.

Table 4: Generation throughput (tokens/s) on an A100. The gray text in brackets denotes batch size.

Model	Context	Full Attention	SHADOWKV	Gain	Full Attention (Inf)
Llama-3-8B-1M (8 KV heads)	60K	160.62 (8)	455.14 (48)	2.83 \times	168.72 (48) / 273.07 (Inf)
	122K	80.77 (4)	239.51 (24)	2.97 \times	83.05 (24) / 134.30 (Inf)
	244K	40.37 (2)	119.01 (12)	2.95 \times	52.00 (12) / 67.15 (Inf)
Llama-3.1-8B (8 KV heads)	60K	160.93 (8)	472.77 (48)	2.94 \times	168.72 (48) / 273.07 (Inf)
	122K	80.78 (4)	245.90 (24)	3.04 \times	83.05 (24) / 134.30 (Inf)
GLM-4-9B-1M (4 KV heads)	60K	241.05 (12)	615.89 (50)	2.56 \times	266.24 (50) / 436.91 (Inf)
	122K	122.67 (6)	293.40 (25)	2.39 \times	158.83 (25) / 214.87 (Inf)
	244K	61.13 (3)	136.51 (12)	2.23 \times	78.84 (12) / 107.44 (Inf)
Yi-9B-200K (4 KV heads)	60K	204.81 (10)	544.36 (42)	2.66 \times	271.21 (42) / 364.09 (Inf)
	122K	101.44 (5)	260.03 (21)	2.56 \times	133.53 (21) / 179.06 (Inf)
	244K	46.74 (2)	118.55 (10)	2.54 \times	65.79 (10) / 89.53 (Inf)

⁵For the equivalent SHADOWKV batch size, we evaluate a single Transformer block with FlashAttention and then project the number to the entire model. For the infinite batch size, we leverage A100’s theoretical memory bandwidth (2 TB/s) for attention computations.

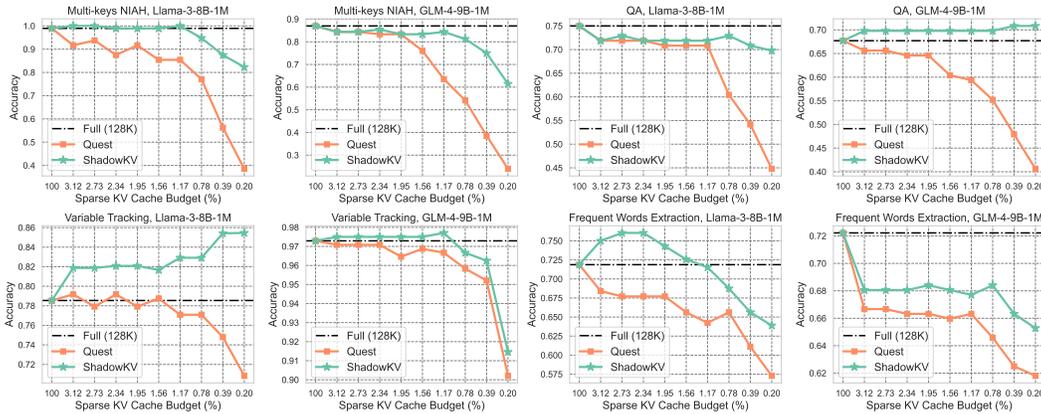


Figure 8: Comparison results between the models with full cache, our SHADOWKV, and Quest.

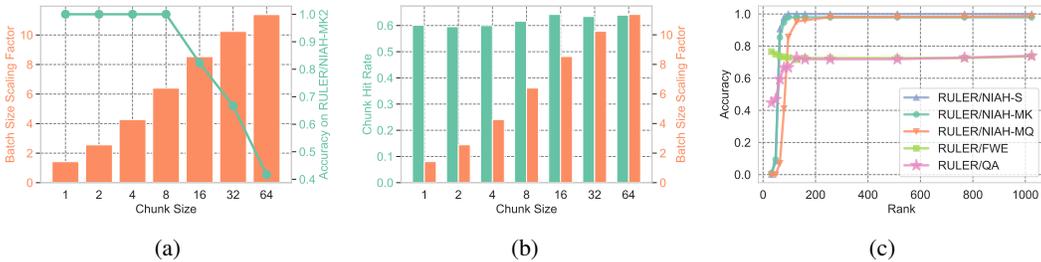


Figure 9: (a) Impact of chunk size on batch size and accuracy. (b) Minimal effect of chunk size on hit rate. (c) Accuracy trends across different ranks with Llama-3-8B-1M on different tasks.

5.3 ABLATION RESULTS

We present extensive ablation studies of SHADOWKV, focusing on three key points: (1) sparse KV cache budget variations, (2) chunk size selections, and (3) pre-RoPE key cache rank choices. [Additional ablations, including precision sensitivity analysis and the effectiveness of outliers, are provided in Appendix A.](#)

Sparse KV Cache Budget. We examine SHADOWKV’s performance across various tasks with different sparse budgets, as illustrated in Figure 8. SHADOWKV consistently surpasses Quest under the same sparse budgets and achieves higher throughput. On most tasks, it maintains accuracy with just a 1.56% sparse budget compared to full attention and even improves slightly on some tasks.

Chunk Size. As shown in Figure 9a, increasing the chunk size allows for larger batch sizes. However, accuracy declines when the chunk size exceeds eight. Meanwhile, the chunk size choice has minimal impact on the chunk hit rate, which remains around 60%, as illustrated in Figure 9b.

Rank of Pre-RoPE Keys. We assess SHADOWKV’s performance across various tasks using different ranks for pre-RoPE keys. As illustrated in Figure 9c, accuracy increases with the rank up to approximately 160, after which it stabilizes near full-rank performance. Interestingly, the trends vary across tasks, and in some cases, low-rank approximations achieve better performance.

6 CONCLUSION

We present SHADOWKV, a high-throughput inference system for long-context LLM inference. SHADOWKV optimizes GPU memory usage through the low-rank key cache and offloaded value cache, allowing for larger batch sizes. It reduces decoding overhead by accurate sparse attention, boosting throughput while maintaining accuracy. Our empirical experiments demonstrate SHADOWKV can support up to $6\times$ larger batch sizes and enhance throughput by up to $3.04\times$ on an A100 across various long-context models, including Llama-3.1-8B, Llama-3-8B-1M, GLM-4-9B-1M, and Yi-9B-200K. SHADOWKV holds great promise for improving long-context LLM inference.

REFERENCES

- 540
541
542 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany
543 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical re-
544 port: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*,
545 2024.
- 546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
547 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
548 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 549
550 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng
551 Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue,
552 Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao
553 Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu,
554 and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.
- 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. *arXiv preprint arXiv:2308.14508*, 2023.
- 558
559 Chi-Chih Chang, Wei-Cheng Lin, Chien-Yu Lin, Chong-Yan Chen, Yu-Fang Hu, Pei-Shuo Wang,
560 Ning-Chi Huang, Luis Ceze, and Kai-Chiang Wu. Palu: Compressing kv-cache with low-rank
561 projection. *arXiv preprint arXiv:2407.21118*, 2024.
- 562 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv*
563 *preprint arXiv:2307.08691*, 2023.
- 564
565 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
566 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,
567 35:16344–16359, 2022.
- 568 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language
569 model, 2024.
- 570
571 Harry Dong, Beidi Chen, and Yuejie Chi. Prompt-prompted adaptive structured pruning for efficient
572 llm generation. In *First Conference on Language Modeling*, 2024a.
- 573 Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. Get more
574 with less: Synthesizing recurrence with kv cache compression for efficient llm inference. *arXiv*
575 *preprint arXiv:2402.09398*, 2024b.
- 576
577 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
578 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text
579 for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 580
581 Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. Model tells
582 you what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*,
583 2023.
- 584
585 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
586 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng,
587 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu,
588 Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao,
589 Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu,
590 Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan
591 Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang,
592 Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language
593 models from glm-130b to glm-4 all tools, 2024.
- 594
595 Gradient. Llama-3-8b-instruct gradient 4194k (v0.1), 2024. URL <https://huggingface.co/gradientai/Llama-3-8B-Instruct-Gradient-1048k>.

- 594 Jiaao He and Jidong Zhai. Fastdecode: High-throughput gpu-efficient llm serving using heteroge-
595 neous pipelines. *arXiv preprint arXiv:2403.11421*, 2024.
- 596
- 597 Ke Hong, Guohao Dai, Jiaming Xu, Qiuli Mao, Xiuhong Li, Jun Liu, Kangdi Chen, Hanyu Dong,
598 and Yu Wang. Flashdecoding++: Faster large language model inference on gpus. *arXiv preprint*
599 *arXiv:2311.01282*, 2023.
- 600 Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao,
601 Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with
602 kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024.
- 603
- 604 Cheng-Ping Hsieh, Simeng Sun, Samuel Krirman, Shantanu Acharya, Dima Rekish, Fei Jia, Yang
605 Zhang, and Boris Ginsburg. Ruler: What’s the real context size of your long-context language
606 models? *arXiv preprint arXiv:2404.06654*, 2024.
- 607 Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua
608 Han, Amir H Abdi, Dongsheng Li, Chin-Yew Lin, et al. Minference 1.0: Accelerating pre-filling
609 for long-context llms via dynamic sparse attention. *arXiv preprint arXiv:2407.02490*, 2024.
- 610 Jordan Juravsky, Bradley Brown, Ryan Ehrlich, Daniel Y Fu, Christopher Ré, and Azalia
611 Mirhoseini. Hydragen: High-throughput llm inference with shared prefixes. *arXiv preprint*
612 *arXiv:2402.05099*, 2024.
- 613
- 614 Greg Kamradt. Needle in a haystack - pressure testing llms. 2023.
- 615 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
616 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
617 serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Prin-*
618 *ciples*, pp. 611–626, 2023.
- 619
- 620 Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. InfiniGen: Efficient generative infer-
621 ence of large language models with dynamic kv cache management. In *18th USENIX Symposium*
622 *on Operating Systems Design and Implementation (OSDI 24)*, 2024a.
- 623 Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. {InfiniGen}: Efficient generative
624 inference of large language models with dynamic {KV} cache management. In *18th USENIX*
625 *Symposium on Operating Systems Design and Implementation (OSDI 24)*, pp. 155–172, 2024b.
- 626
- 627 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative
628 decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
- 629 Yucheng Li, Bo Dong, Chenghua Lin, and Frank Guerin. Compressing context to enhance inference
630 efficiency of large language models. *arXiv preprint arXiv:2310.06201*, 2023.
- 631
- 632 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle
633 Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before
634 generation. *arXiv preprint arXiv:2404.14469*, 2024.
- 635 Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and
636 language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024a.
- 637
- 638 Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios
639 Kyriillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance
640 hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing*
641 *Systems*, 36, 2024b.
- 642 Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi
643 Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint*
644 *arXiv:2402.02750*, 2024c.
- 645
- 646 Meta AI. Introducing Llama 3.1, 2024. URL [https://ai.meta.com/blog/
meta-llama-3-1/](https://ai.meta.com/blog/meta-llama-3-1/). Accessed: 2024-08-21.
- 647
- Microsoft. Microsoft bingchat, 2024. URL <https://www.bing.com/chat>.

- 648 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
649 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
650 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
651
- 652 QwenTeam. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.](https://qwenlm.github.io/blog/qwen2.5/)
653 [github.io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 654 Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive
655 transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
656
- 657 Luka Ribar, Ivan Chelombiev, Luke Hudlass-Galley, Charlie Blake, Carlo Luschi, and Douglas Orr.
658 Sparq attention: Bandwidth-efficient llm inference. *arXiv preprint arXiv:2312.04985*, 2023.
- 659 Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang,
660 Christopher Ré, Ion Stoica, and Ce Zhang. Flexgen: High-throughput generative inference of
661 large language models with a single gpu. In *International Conference on Machine Learning*, pp.
662 31094–31116. PMLR, 2023.
- 663
- 664 Prajwal Singhanian, Siddharth Singh, Shwai He, Soheil Feizi, and Abhinav Bhatel. Loki: Low-rank
665 keys for efficient sparse attention. *arXiv preprint arXiv:2406.02542*, 2024.
- 666
- 667 Zezheng Song, Jiaxin Yuan, and Haizhao Yang. Fmint: Bridging human designed and data pre-
668 trained models for differential equation foundation model. *arXiv preprint arXiv:2404.14688*,
669 2024.
- 670 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
671 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 672
- 673 Hanshi Sun, Zhuoming Chen, Xinyu Yang, Yuandong Tian, and Beidi Chen. Triforce: Lossless
674 acceleration of long sequence generation with hierarchical speculative decoding. *arXiv preprint*
675 *arXiv:2404.11912*, 2024.
- 676
- 677 Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Shikuan Hong, Yiwu Yao, and Gongyi Wang.
678 Razorattention: Efficient kv cache compression through retrieval heads. *arXiv preprint*
arXiv:2407.15891, 2024a.
- 679
- 680 Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest:
681 Query-aware sparsity for efficient long-context llm inference. *arXiv preprint arXiv:2406.10774*,
682 2024b.
- 683
- 684 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
685 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
686 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 687
- 688 Vijay Thakkar, Pradeep Ramani, Cris Cecka, Aniket Shivam, Honghao Lu, Ethan Yan, Jack Kosa-
689 an, Mark Hoemmen, Haicheng Wu, Andrew Kerr, Matt Nicely, Duane Merrill, Dustyn Blasig,
690 Fengqi Qiao, Piotr Majcher, Paul Springer, Markus Hohnerbach, Jin Wang, and Manish Gupta.
691 CUTLASS, January 2023. URL <https://github.com/NVIDIA/cutlass>.
- 692
- 693 Haixin Wang, Xinlong Yang, Jianlong Chang, Dian Jin, Jinan Sun, Shikun Zhang, Xiao Luo, and
694 Qi Tian. Parameter-efficient tuning of large-scale multimodal foundation model. *Advances in*
695 *Neural Information Processing Systems*, 36, 2024a.
- 696
- 697 Minzheng Wang, Longze Chen, Cheng Fu, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu,
698 Nan Xu, Lei Zhang, Run Luo, et al. Leave no document behind: Benchmarking long-context llms
699 with extended multi-doc qa. *arXiv preprint arXiv:2406.17419*, 2024b.
- 700
- 701 T Wolf. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint*
arXiv:1910.03771, 2019.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant:
Accurate and efficient post-training quantization for large language models. In *International
Conference on Machine Learning*, pp. 38087–38099. PMLR, 2023a.

- 702 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
703 language models with attention sinks. In *The Twelfth International Conference on Learning Rep-*
704 *resentations*, 2023b.
- 705
706 Yuhui Xu, Zhanming Jie, Hanze Dong, Lei Wang, Xudong Lu, Aojun Zhou, Amrita Saha, Caiming
707 Xiong, and Doyen Sahoo. Think: Thinner key cache by query-driven pruning. *arXiv preprint*
708 *arXiv:2407.21018*, 2024.
- 709 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
710 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
711 *arXiv:2407.10671*, 2024a.
- 712
713 June Yong Yang, Byeongwook Kim, Jeongin Bae, Beomseok Kwon, Gunho Park, Eunho Yang,
714 Se Jung Kwon, and Dongsoo Lee. No token left behind: Reliable kv cache compression via
715 importance-aware mixed precision quantization. *arXiv preprint arXiv:2402.18096*, 2024b.
- 716 Zihao Ye, Ruihang Lai, Bo-Ru Lu, Lin Chien-Yu, Size Zheng, Lequn Chen, Tianqi Chen, and
717 Luis Ceze. Cascade inference: Memory bandwidth efficient shared prefix batch decoding,
718 2024. URL <https://flashinfer.ai/2024/02/02/cascade-inference.html>.
719 Accessed: 2024-09-25.
- 720 Yixiao Yuan, Yangchen Huang, Yu Ma, Xinjin Li, Zhenglin Li, Yiming Shi, and Huapeng Zhou.
721 Rhyme-aware chinese lyric generator based on gpt. *arXiv preprint arXiv:2408.10130*, 2024.
- 722
723 Yuxuan Yue, Zhihang Yuan, Haojie Duanmu, Sifan Zhou, Jianlong Wu, and Liqiang Nie. Wkvquant:
724 Quantizing weight and key/value cache for large language models gains more. *arXiv preprint*
725 *arXiv:2402.12065*, 2024.
- 726 Hailin Zhang, Xiaodong Ji, Yilin Chen, Fangcheng Fu, Xupeng Miao, Xiaonan Nie, Weipeng Chen,
727 and Bin Cui. Pqcache: Product quantization-based kvcache for long context llm inference. *arXiv*
728 *preprint arXiv:2407.12820*, 2024a.
- 729
730 Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han,
731 Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. ∞ bench: Extending long context
732 evaluation beyond 100k tokens, 2024b.
- 733 Yichi Zhang, Bofei Gao, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie
734 Hu, Wen Xiao, et al. Pyramidkv: Dynamic kv cache compression based on pyramidal information
735 funneling. *arXiv preprint arXiv:2406.02069*, 2024c.
- 736
737 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,
738 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient gener-
739 ative inference of large language models. *Advances in Neural Information Processing Systems*,
740 36, 2024d.
- 741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A ADDITIONAL EXPERIMENT RESULTS

In this section, we present additional experiments not covered in the main text, including the handling of newly generated tokens (as discussed in Section 4.1), scalability analysis for larger models and longer sequences (mentioned in Section 5.1), latency breakdown (mentioned in Section 5.2), additional ablation studies (referenced in Section 5.3), and etc.

A.1 HANDLING OF NEWLY GENERATED TOKENS

To address the handling of newly generated tokens, we project these tokens’ key cache into a low-rank space using the same projections applied during the prefilling phase. This approach preserves the benefits of reduced GPU memory usage, particularly for long output sequences.

As shown in Table 5 and Table 6, we refer to this extension as SHADOWKV +. Our evaluation across various models demonstrates that SHADOWKV + effectively maintains accuracy while optimizing memory usage.

Table 5: Performance of SHADOWKV and SHADOWKV + across different models on RULER (Hsieh et al., 2024) evaluated at length of 128K.

Methods	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
<i>Llama-3-8B-1M</i>	100.00	100.00	98.96	98.96	98.96	95.57	75.00	48.96	78.54	71.85	86.68
SHADOWKV	100.00	100.00	97.92	98.96	96.88	95.83	72.92	52.08	81.67	72.57	86.88
SHADOWKV +	100.00	100.00	98.96	100.00	95.83	93.49	71.88	50.00	80.21	71.88	86.23
<i>GLM-4-9B-1M</i>	100.00	100.00	94.79	87.50	99.74	93.75	67.71	55.21	97.29	72.22	86.82
SHADOWKV	100.00	100.00	95.83	83.33	98.70	87.76	69.79	55.21	97.50	68.06	85.62
SHADOWKV +	100.00	100.00	95.83	85.42	98.17	85.16	69.79	56.25	97.92	67.71	85.63
<i>Llama-3.1-8B</i>	100.00	100.00	98.96	91.67	98.96	95.31	82.29	47.92	68.96	71.18	85.53
SHADOWKV	100.00	100.00	100.00	83.33	97.92	92.19	81.25	48.96	67.08	64.93	83.57
SHADOWKV +	100.00	100.00	100.00	84.38	96.88	91.67	81.25	52.08	65.63	62.85	83.47
<i>Yi-9B-200K</i>	100.00	100.00	86.46	62.50	64.58	32.55	44.79	39.58	36.87	89.93	65.73
SHADOWKV	100.00	100.00	82.29	67.71	63.28	31.51	43.75	38.54	56.04	72.22	65.53
SHADOWKV +	100.00	100.00	81.25	67.71	61.72	31.51	46.88	38.54	53.54	72.92	65.41

Table 6: Performance of SHADOWKV and SHADOWKV + on different models with LongBench (Bai et al., 2023) samples exceeding 4K tokens.

Methods	NarratQA	MultiFQA	HotpotQA	MuSiQue	DuRead	GovRep	SAMSum	PassRetr	LCC	Avg.
<i>Llama-3-8B-1M</i>	18.98	41.84	36.79	21.47	31.93	34.18	35.96	81.50	56.07	39.86
SHADOWKV	17.17	39.73	38.29	21.08	31.77	31.62	35.87	80.00	63.93	39.94
SHADOWKV +	20.42	41.16	37.22	21.03	31.77	31.98	35.80	80.00	63.89	40.36
<i>GLM-4-9B-1M</i>	25.44	51.09	58.67	39.61	32.04	29.97	40.31	99.00	58.02	48.24
SHADOWKV	26.50	51.31	59.09	38.87	32.92	28.54	38.70	96.50	58.55	47.89
SHADOWKV +	27.59	51.31	59.17	38.34	33.55	31.25	39.46	96.50	55.86	48.11
<i>Llama-3.1-8B</i>	31.56	55.10	57.65	29.46	35.26	34.45	29.84	100.00	67.31	48.96
SHADOWKV	30.93	55.20	57.32	29.13	31.85	32.79	30.40	99.50	66.03	48.13
SHADOWKV +	32.25	54.29	57.75	28.37	31.07	32.89	28.73	98.75	67.59	47.97
<i>Yi-9B-200K</i>	13.88	30.02	52.46	28.20	22.29	30.25	19.08	67.00	73.50	37.41
SHADOWKV	12.44	30.82	52.43	27.73	20.79	29.83	20.73	64.00	72.89	36.85
SHADOWKV +	14.08	30.94	51.16	27.00	19.50	29.34	21.16	66.00	73.47	36.96

A.2 QUANTITATIVE ANALYSIS OF GPU MEMORY SAVINGS

The GPU memory savings provided by SHADOWKV can be quantitatively analyzed as follows. Let each K or V vector have a size of M bytes, with a sequence length S , a chunk size C , a selected chunk budget K , O outliers, and a pre-RoPE key cache rank r . The GPU memory savings of SHADOWKV can then be expressed as:

$$\text{Memory Savings} = \frac{2SM}{SM/C + 2(K + O)C + Sr + rM}$$

For example, assuming $M = 1024$, $C = 8$, $S = 128\text{K}$, $K = 256$, $O = 48$, $r = 160$, the memory savings of SHADOWKV is calculated as $7.08\times$. This result demonstrates that SHADOWKV can theoretically reduce the KV cache memory footprint on the GPU by $7.08\times$ for longer sequences and larger batch sizes.

A.3 ACCURACY RESULTS FOR YI-9B-200K

We present accuracy results for Yi-9B-200K (AI et al., 2024) on RULER (Hsieh et al., 2024) and LongBench (Bai et al., 2023), highlighting SHADOWKV’s superior performance across diverse tasks compared to other methods.

Table 7: Performance of Yi-9B-200K with different methods on RULER (Hsieh et al., 2024) evaluated at length of 128K. SHADOWKV outperforms other methods with a 1.56% sparse budget.

Methods	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
<i>Yi-9B-200K</i>	100.00	100.00	86.46	62.50	64.58	32.55	44.79	39.58	36.87	89.93	65.73
Loki	34.38	2.08	2.08	0.00	0.00	0.52	22.92	21.88	0.00	25.00	10.89
Loki (V only)	59.38	11.46	18.75	5.21	4.43	2.08	22.92	31.25	0.00	35.07	19.06
InfiniGen	100.00	94.79	77.08	1.04	40.10	20.57	37.50	34.38	41.46	46.18	49.31
InfiniGen (V only)	100.00	98.96	78.13	2.08	58.33	24.48	40.63	35.42	52.92	55.90	54.69
Quest	100.00	98.96	79.17	26.04	56.51	31.77	32.29	31.25	51.04	71.88	57.89
Quest (V only)	100.00	100.00	80.21	45.83	59.37	31.90	36.45	34.37	53.54	71.88	61.36
SHADOWKV	100.00	100.00	82.29	67.71	63.28	31.51	43.75	38.54	56.04	72.22	65.53

Table 8: Performance of Yi-9B-200K with LongBench (Bai et al., 2023) samples exceeding 4K tokens. SHADOWKV outperforms other methods and maintains the accuracy.

Methods	NarrQA	MultiFQA	HotpotQA	MuSiQue	DuRead	GovRep	SAMSum	PassRetr	LCC	Avg.
<i>Yi-9B-200K</i>	13.88	30.02	52.46	28.20	22.29	30.25	19.08	67.00	73.50	37.41
Loki	1.63	2.73	16.21	4.87	4.75	2.13	4.95	0.00	38.72	8.44
Loki (V only)	1.96	10.39	21.31	7.36	6.78	9.15	10.02	4.00	58.75	14.41
InfiniGen	10.01	23.61	50.47	25.91	15.11	27.96	18.97	30.00	56.46	28.72
InfiniGen (V only)	11.31	26.46	51.13	26.77	16.09	28.67	19.33	34.00	62.07	30.65
Quest	10.57	25.83	46.06	23.04	17.09	17.11	20.59	50.50	67.70	30.94
Quest (V only)	14.56	25.73	48.73	24.73	18.44	20.83	20.08	57.50	71.13	33.53
SHADOWKV	12.44	30.82	52.43	27.73	20.79	29.83	20.73	64.00	72.89	36.85

A.4 PRECISION SENSITIVITY

In the main experiments, we used BF16 for both model weights and KV cache. To further investigate the impact of precision on SHADOWKV’s performance, we conducted additional experiments using FP8 precision (`torch.float8_e5m2`). These tests aim to determine whether SHADOWKV can retain its accuracy at this lower precision, addressing concerns about precision sensitivity, particularly in SVD computations.

As detailed in Table 9 and Table 10, SHADOWKV and baseline methods were evaluated using FP8. Results show that SHADOWKV maintains accuracy and achieves consistently high performance even with FP8 precision. This robustness, despite FP8’s reduced numerical range, confirms that SHADOWKV can continue to deliver efficiency gains without compromising accuracy.

Table 9: Performance comparison of SHADOWKV and baseline methods on the RULER (Hsieh et al., 2024) using FP8 precision, evaluated at a sequence length of 128K.

Methods	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
<i>Llama-3-8B-1M</i>	100.00	100.00	98.96	95.83	97.40	95.57	63.54	48.96	75.83	73.26	84.94
Loki	5.21	1.04	0.00	0.00	0.78	0.26	5.21	13.54	28.33	28.82	8.32
Loki (V only)	36.46	9.38	31.25	0.00	6.25	21.09	11.46	15.63	57.08	35.76	22.44
Quest	100.00	98.96	98.96	71.88	96.61	93.49	63.54	45.83	78.13	67.01	81.44
Quest (V only)	100.00	100.00	98.96	85.42	97.40	93.49	70.83	48.96	78.13	65.63	83.88
SHADOWKV	100.00	100.00	97.92	94.79	95.31	93.49	75.00	48.96	80.42	73.61	85.95

Table 10: Evaluation of SHADOWKV and baseline methods on LongBench (Bai et al., 2023) with sequence lengths exceeding 4K tokens, using FP8 precision.

Methods	NarratQA	MultiFQA	HotpotQA	MuSiQue	DuRead	GovRep	SAMSum	PassRetr	LCC	Avg.
<i>Llama-3-8B-1M</i>	18.69	41.21	35.76	21.59	31.81	33.77	35.29	80.50	56.77	39.49
Loki	2.21	11.12	5.70	1.84	15.42	28.59	11.41	41.91	33.99	16.91
Loki (V only)	2.68	22.33	12.69	3.35	21.43	30.57	16.32	47.68	36.64	21.52
Quest	19.41	38.92	34.02	19.64	23.13	26.40	28.04	78.50	49.81	35.32
Quest (V only)	16.19	36.73	36.64	19.59	25.57	29.46	27.14	79.50	60.05	36.76
SHADOWKV	18.29	39.39	36.06	21.04	30.47	31.87	35.56	78.50	62.11	39.25

A.5 SCALABILITY ANALYSIS FOR LARGER MODELS AND LONGER SEQUENCES

To demonstrate the scalability of SHADOWKV, we present experiments with Llama-3-8B-1M on 1M contexts and Llama-3-70B-1M on 512K contexts, using the RULER benchmark (Hsieh et al., 2024). Additionally, we evaluate Llama-3-70B-1M on the Needle In A Haystack dataset, testing context lengths ranging from 16K to 1M tokens.

As shown in Figure 10 and Table 11, SHADOWKV maintains robust performance across increasing context lengths and model sizes, demonstrating its scalability in handling large-scale inputs. This scalability allows SHADOWKV to process extensive contexts with high accuracy, making it a valuable solution for real-world applications requiring extensive sequences.

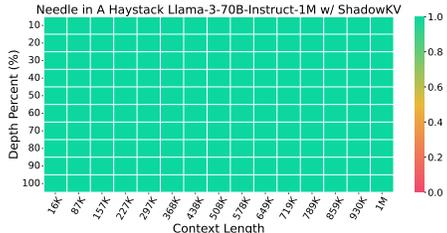


Figure 10: Needle In A Haystack.

Table 11: Performance of different methods on RULER (Hsieh et al., 2024) evaluated at length of 1M. The Llama-3-8B-1M is evaluated on 1M contexts while the Llama-3-70B-1M is evaluated on 512K contexts.

Methods	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
<i>Llama-3-70B-1M</i>	100.00	82.29	90.63	54.17	85.16	96.61	69.79	35.42	68.75	69.44	75.23
Loki	100.00	1.04	0.00	0.00	0.00	0.00	13.54	11.46	34.30	22.92	18.33
Loki (V only)	100.00	15.63	26.04	0.00	0.00	0.00	25.00	19.79	40.00	31.94	25.84
Quest	100.00	76.04	78.13	35.42	85.47	92.19	53.21	34.38	38.33	58.33	65.15
Quest (V only)	100.00	77.08	79.17	36.49	86.19	95.31	54.17	36.58	47.70	58.68	67.14
SHADOWKV	100.00	82.29	88.54	53.04	88.02	94.79	67.71	37.50	68.54	68.25	74.87
<i>Llama-3-8B-1M</i>	96.88	100.00	96.88	69.79	91.15	85.68	64.58	42.71	25.00	56.25	72.89
Loki	9.38	1.04	10.42	0.00	2.60	4.43	38.54	11.46	1.67	0.69	8.02
Loki (V only)	68.75	29.17	60.42	1.04	26.56	43.23	59.38	15.63	6.46	0.69	31.13
Quest	94.79	92.71	80.21	4.17	76.30	69.27	57.29	28.13	25.67	30.56	55.91
Quest (V only)	94.79	93.75	81.25	4.17	79.69	69.27	62.50	31.25	26.00	32.99	57.57
SHADOWKV	96.88	100.00	96.88	65.63	89.38	83.16	69.79	42.71	26.04	59.38	72.98

A.6 LATENCY BREAKDOWN

We present a detailed latency breakdown in Table 12 and Table 13 to illustrate the efficiency of each operation under various context lengths for both the prefilling and decoding stages.

Scalability for Longer Sequences. As shown in Table 12, the overhead of SVD, reduce, cosine similarity, topK, and gather computing is very low and tends to decrease as the sequence scales, proving that SHADOWKV’s scalability to longer sequences.

Table 12: Latency breakdown (ms) of a Transformer block of Llama-3-8B-1M during prefilling.

Context	Attention	FFN	SVD	Reduce	CosineSimilarity	TopK	Gather	Cost
64K	186.23	96.47	17.19	0.10	1.41	0.08	0.01	6.65%
128K	721.13	193.32	26.62	0.20	2.77	0.14	0.02	3.25%
256K	2880.21	392.77	50.56	0.42	6.11	0.11	0.03	1.75%
512K	11720.30	789.23	108.38	0.84	12.19	0.15	0.06	0.97%

Overlapping Operations for Latency Reduction. In Table 13, we demonstrate how overlapping the recomputation of the key cache with value cache fetching from the CPU significantly reduces decoding latency. This concurrent processing approach ensures that SHADOWKV minimizes overhead when handling long-context models.

Table 13: Latency breakdown (ms) of a Transformer block of Llama-3-8B-1M during decoding.

Context	GEMM+ Softmax	Max	TopK	Recompute K (Overlapped)	Fetch V	Attention	FFN	QKV
48×64K	0.56	0.07	0.14	1.25	1.84	0.23	0.33	0.05
24×128K	0.58	0.07	0.15	1.36	1.66	0.21	0.29	0.05
12×256K	0.65	0.07	0.16	1.49	1.75	0.19	0.25	0.05
6×512K	0.71	0.07	0.17	1.51	1.69	0.18	0.23	0.05

A.7 EFFICIENCY COMPARISON WITH QUEST

We present an efficiency comparison with Quest, particularly under long contexts or high batch sizes where the GPU memory alone cannot accommodate the KV cache. In such cases, both Full Attention and Quest must offload the KV cache to the CPU. As shown in Table 14, SHADOWKV significantly outperforms both Full Attention and Quest under the same sparse budget.

The efficiency advantage of SHADOWKV over Quest is due to two key factors: (1) SHADOWKV only fetches the value cache from the CPU, rather than the entire KV pair, minimizing data transfer and reducing latency, and (2) SHADOWKV integrates a cache mechanism that leverages the temporal locality of the KV cache.

Table 14: Efficiency comparison with Quest.

Context	Full Attention	Full Attention (CPU)	Quest	Quest (CPU)	SHADOWKV
3×1M	OOM	0.21 tokens/s	OOM	9.34 tokens/s	45.32 tokens/s

A.8 ACCURACY CONTRIBUTION OF OUTLIER KV CACHE

We conduct experiments using different numbers of outlier chunks for Llama-3-8B-1M on the RULER benchmark with 128K context length. As presented in Table 15, our findings indicate that outliers play a crucial role. For instance, the first chunk, a significant outlier, has previously been shown to act as an attention sink (Xiao et al., 2023b), underscoring its importance in maintaining model accuracy.

The results demonstrate that increasing the number of outlier chunks has a positive impact on accuracy, especially in complex tasks. This indicates that even a small number of outliers can effectively capture essential information, reducing the need for full attention. Remarkably, with just 8 outliers (0.049%), SHADOWKV outperforms the Quest baseline and nearly matches the accuracy achieved by full attention. However, when outliers are not adequately managed, the performance of the mean-based landmarks in SHADOWKV may fall below the min-max approach used by Quest, underscoring the importance of handling outliers properly.

Table 15: Performance across different number of outlier chunks on RULER (Hsieh et al., 2024) evaluated at length of 128K.

# Outliers	N-S1	N-S2	N-MK1	N-MK2	N-MQ	N-MV	QA-1	QA-2	VT	FWE	Avg.
0 (0.000 %)	100.00	100.00	96.88	85.42	73.18	70.83	43.75	39.58	73.54	57.29	74.05
1 (0.006 %)	100.00	100.00	97.92	98.96	95.83	94.79	70.83	51.04	70.63	70.14	85.01
2 (0.012 %)	100.00	100.00	97.92	98.96	95.57	95.57	70.83	51.04	72.08	70.49	85.25
4 (0.024 %)	100.00	100.00	97.92	98.96	95.83	95.57	71.88	51.04	74.38	71.18	85.68
8 (0.049 %)	100.00	100.00	97.92	98.96	95.57	95.05	72.92	51.04	78.13	72.57	86.22
16 (0.098 %)	100.00	100.00	97.92	98.96	96.09	95.31	72.92	51.04	80.42	71.53	86.42
32 (0.195 %)	100.00	100.00	97.92	98.96	96.35	95.57	72.92	52.08	81.25	72.22	86.73
48 (0.293 %)	100.00	100.00	97.92	98.96	96.88	95.83	72.92	52.08	81.67	72.57	86.88
<i>Quest (Ref.)</i>	100.00	100.00	98.96	77.08	97.65	93.49	60.42	50.00	77.08	65.63	82.03
<i>Full Attn (Ref.)</i>	100.00	100.00	98.96	98.96	98.96	95.57	75.00	48.96	78.54	71.85	86.68

A.9 DETAILED COMPARISON WITH INFINIGEN

We provide further clarification on the key distinctions and conduct additional experiments between SHADOWKV and InfiniGen. These experiments show that SHADOWKV significantly outperforms InfiniGen across a wide range of downstream tasks.

Differences in SVD Usage. Infinigen performs an offline SVD to get a projection matrix, which is applied to post-RoPE key and query states for KV selection, while SHADOWKV applies an online, prompt-dependent SVD directly to the pre-RoPE key cache for compression, not for KV selection.

Methodological Differences. While InfiniGen uses SVD for KV selection, it requires fetching selected, exact KV pairs from the CPU. In contrast, SHADOWKV only fetches the value cache from the CPU, reconstructing the key cache from its low-rank storage on the GPU. By overlapping these processes, SHADOWKV reduces data-fetch overhead and achieves improved efficiency in KV cache management.

Accuracy Comparison. To empirically validate SHADOWKV’s advantages, we conducted accuracy evaluations. Results confirm SHADOWKV’s effectiveness in maintaining accuracy while optimizing memory usage. Although InfiniGen performs well on simpler tasks like RULER-N-S1, it shows significant accuracy drops on more complex tasks, such as RULER-N-MK2, RULER-FWE, LongBench-LCC, and others, where SHADOWKV maintains consistently high accuracy.

B EXPERIMENT DETAILS

In this section, our goal is to provide the details of the system implementation (mentioned in Section 4.2), experiment settings, and additional experiments (mentioned in Section 5).

B.1 SYSTEM IMPLEMENTATION.

We implement the framework based on PyTorch (Paszke et al., 2019; Wolf, 2019) and dedicated kernels (Thakkar et al., 2023). FlashAttention (Dao et al., 2022; Dao, 2023; Hong et al., 2023) is used for attention computation and some efficient fused kernels in Flashinfer (Ye et al., 2024) and vLLM (Kwon et al., 2023) are used, including layer norm. To reduce memory movement and kernel launch overhead, we fuse some operations into CUDA kernels, including attention approximation, key cache low-rank reconstruction, value cache fetching, cache mechanism, etc. We leverage multi-streams to overlap the reconstruction of key cache and value cache fetching. We set the rank of pre-RoPE key cache to 160, chunk size to 8, and sparse KV cache budget to 1.56% for most cases.

B.2 DATASET DETAILS

LLMs are widely used in various fields (Li et al., 2023; Yuan et al., 2024; QwenTeam, 2024; Wang et al., 2024a; Song et al., 2024), and we select three long-context benchmarks, detailed below.

- RULER (Hsieh et al., 2024) consists of 13 complex tasks and supports adjustable context lengths, including retrieval, multi-hop tracking, aggregation, and QA tasks. For the test with MInference (Jiang et al., 2024), we set up test sets scaling from 8K to 256K for evaluation.
- LongBench (Bai et al., 2023) is a challenging long-context benchmark that assesses the performance of LLMs in extended contexts. Featuring Chinese and English languages, LongBench encompasses 6 main categories and 21 diverse tasks, evaluating LLM capabilities across crucial long-text applications like single-/multi-document QA, summarization, code completion, etc.
- Needle In A Haystack (Kamradt, 2023) is a long-context retrieval benchmark testing LLM’s performance with context window scales up to 1M tokens where information placed at various positions. We tested the retrieval capabilities of six long-context LLMs based on their context length.

B.3 NEEDLE IN A HAYSTACK

In addition to the Needle In A Haystack results for Llama-3-8B-1M shown in Figure 6, we also present results for GLM-4-9B-1M, Llama-3.1-8B, Yi-9B-200K, Phi-3-Mini-128K, and Qwen2-7B-128K, shown in Figure 11. Compared to full attention, using SHADOWKV has minimal impact on the ability to understand semantic information across different context windows and needle depths. There is even a slight performance improvement for Yi-9B-200K.

B.4 INFINITEBENCH

InfiniteBench (Zhang et al., 2024b) is a challenging long-context benchmark that consists of 10 tasks, including QA, coding, dialogue, summarization, and retrieval, with an average length of 214K.

Table 16: Accuracy of different methods on InfiniteBench (Zhang et al., 2024b).

Methods	En.Sum	En.QA	En.MC	En.Dia	Zh.QA	Code.Debug	Math.Find	Retr.PassKey	Retr.Num
<i>Llama-3-8B-1M</i>	23.05	18.14	65.06	10.50	12.47	24.36	37.14	100.00	100.00
SHADOWKV	21.50	17.73	64.63	10.50	12.45	23.86	37.43	100.00	100.00
<i>GLM-4-9B-1M</i>	28.61	9.25	68.12	39.50	11.77	30.20	40.00	100.00	100.00
SHADOWKV	23.22	8.48	68.56	32.50	11.27	30.46	40.00	100.00	100.00
<i>Llama-3.1-8B</i>	26.42	14.48	66.38	16.00	12.92	21.07	34.00	100.00	99.66
SHADOWKV	24.23	13.83	66.38	16.50	12.76	21.07	34.00	100.00	94.41
<i>Yi-9B-200K</i>	8.88	10.61	61.57	5.50	13.88	21.57	23.71	100.00	99.66
SHADOWKV	8.92	10.06	59.39	6.00	13.89	20.56	24.29	100.00	99.83

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

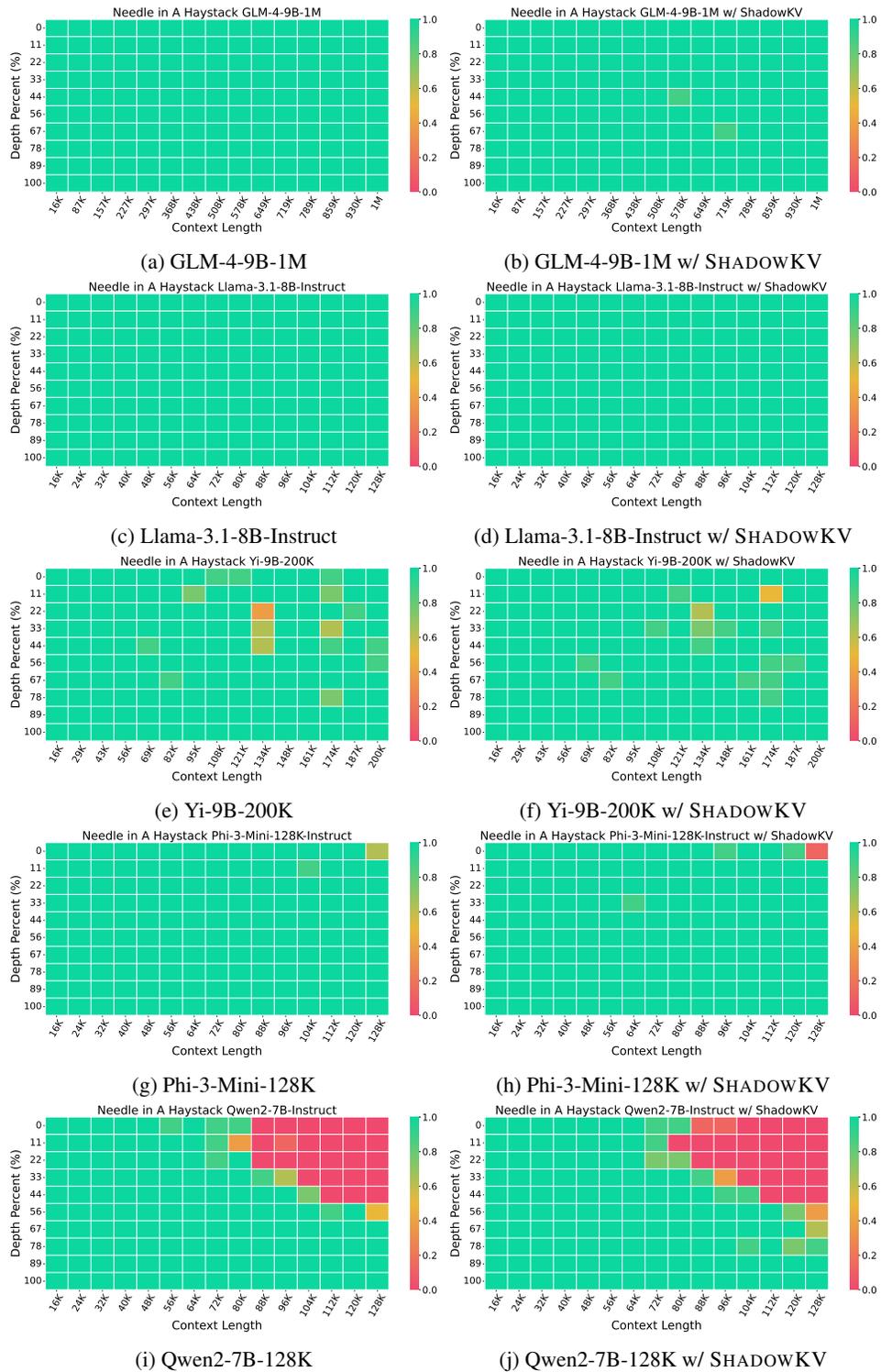


Figure 11: Needle In A Haystack (Kamradt, 2023) results using GLM-4-9B-1M (GLM et al., 2024), Llama-3.1-8B-Instruct (Meta AI, 2024), Yi-9B-200K (AI et al., 2024), Phi-3-Mini-128K (Abdin et al., 2024), and Qwen2-7B-128K (Yang et al., 2024a).