

BEYOND HOMOGENEOUS ATTENTION: MEMORY-EFFICIENT LLMs VIA FOURIER-APPROXIMATED KV CACHE

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models struggle with memory demands from the growing Key-Value (KV) cache as context lengths increase. Existing compression methods homogenize head dimensions or rely on attention-guided token pruning, often sacrificing accuracy or introducing computational overhead. We propose **FourierAttention**, a training-free framework that exploits the heterogeneous roles of transformer head dimensions: lower dimensions prioritize local context, while upper ones capture long-range dependencies. By projecting the long-context-insensitive dimensions onto orthogonal Fourier bases, FourierAttention approximates their temporal evolution with fixed-length spectral coefficients. Evaluations on LLaMA models show FourierAttention achieves the best long-context accuracy on LongBench and Needle-In-A-Haystack (NIAH). Besides, a custom Triton kernel, **FlashFourierAttention**, is designed to optimize memory via streamlined read-write operations, enabling efficient deployment without performance compromise.

1 INTRODUCTION

Large Language Models (LLMs) have transformed natural language processing with breakthroughs in text generation, comprehension, and reasoning (OpenAI, 2023; Sun et al., 2024; OpenAI, 2024; Guo et al., 2025). However, their autoregressive decoding relies heavily on a memory-intensive Key-Value (KV) cache, leading to significant memory allocation as context lengths scale (Vaswani et al., 2017; Fu, 2024; Liu et al., 2025). This overhead limits LLM deployment in resource-constrained environments. While approaches like sparse attention and cache compression have been explored to reduce memory needs, they often compromise accuracy or add complexity (Cai et al., 2024; Yuan et al., 2025). Developing memory-efficient methods that preserve performance remains crucial for the broader applicability of LLMs.

Existing training-free KV cache compression methods, like token eviction strategies (Xiao et al., 2024; Zhang et al., 2023; Li et al., 2024b), prune sequence subsets but overlook the heterogeneous roles of head dimensions, leaving dimension-aware allocation largely unexplored. Similarly, hidden dimension compression (Chang et al., 2024; Saxena et al., 2024) methods apply uniform ratios, both neglect their distinct contribution across dimensions (Liu et al., 2024b; Peng et al., 2024). These approaches treat head dimensions as homogeneous, static units rather than dynamically allocating resources based on their importance.

Another critical limitation of existing methods lies in their reliance on attention-guided strategies (Zhang et al., 2023; Li et al., 2024b). While these approaches enable selective token pruning with minimal accuracy degradation, they impose prohibitive memory and latency overheads due to attention score recalculation. We address this challenge by adapting the HiPPO framework (Gu et al., 2020), a mathematically grounded approach for long-sequence modeling. HiPPO approximates infinite-length sequences as compact finite states by projecting inputs onto finite-order orthogonal basis functions, such as polynomial bases or Fourier bases (Gu et al., 2020; He et al., 2023). This retains global critical and contextually vital patterns while filtering out redundant signals. By leveraging HiPPO’s theoretical foundations, we can bypass attention recomputation entirely, achieving both memory efficiency and computational efficiency.

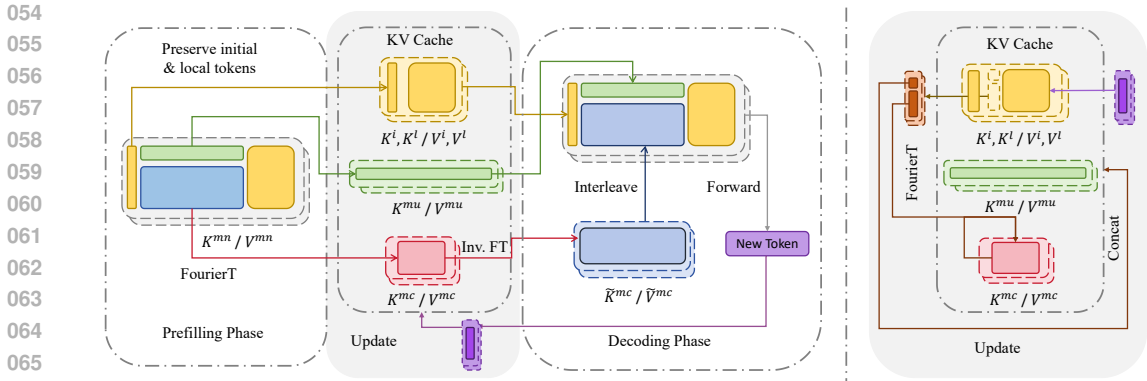


Figure 1: Overview of FourierAttention.

Building on these insights, we introduce **FourierAttention**, a training-free KV cache compression framework using a translated Fourier transform, as shown in Figure 1. Departing from prior methods that uniformly process all head dimensions, FourierAttention identifies localized, context-insensitive dimensions in KV states and approximates their temporal evolution via a fixed set of orthogonal Fourier basis functions. By retaining only the dominant Fourier coefficients ($N \ll L$, where L is the sequence length), our method projects sequences into a compact spectral representation. Unlike polynomial bases that are widely used in HiPPO, FourierAttention exploits the shift-invariance and temporal parallelism of Fourier transforms, allowing for efficient computation and higher performance. During decoding, a customized Triton kernel **FlashFourierAttention** is used to decompose KV cache states during attention calculation, minimizing memory overhead via streamlined read-write operations. Our contributions can be summarized as follows:

- We reveal a bifurcation in Transformer head dimensions: lower dimensions prioritize local context, while upper ones capture long-range dependencies. This inspires us to compress long-context-insensitive dimensions without sacrificing contextual awareness.
- We introduce FourierAttention, which optimizes KV cache by projecting its temporal evolution onto a fixed set of orthogonal Fourier bases. This method efficiently eliminates redundant components while preserving contextual fidelity, achieving a balance between memory and computational efficiency.
- We evaluate FourierAttention’s performance on the LLaMA Series using LongBench and NIAH. Our FourierAttention achieves the consistent superiority of long-context performance over other cache optimization methods while maintaining lower memory consumption.

2 RELATED WORK

KV cache optimization is a crucial technique for enhancing efficiency in attention-based LLMs (Fu, 2024; Liu et al., 2025). As context length increases, the KV cache in LLMs grows linearly, creating substantial memory overhead that becomes a bottleneck for long-context applications. Beyond architectural modifications during pretraining (Ainslie et al., 2023; Liu et al., 2024a), existing training-free optimization approaches mainly involve token eviction or compression. The former discards tokens based on positional or attention patterns, including StreamingLLM (Xiao et al., 2024), H2O (Zhang et al., 2023), SnapKV (Li et al., 2024b), and PyramidKV (Cai et al., 2024), while the latter compresses KV cache through quantization or low-rank projection, such as KIVI (Liu et al., 2024c), KVQuant (Hooper et al., 2024), and Palu (Chang et al., 2024). However, these methods lack fine-grained consideration of different head dimensions in the KV cache, applying uniform optimization across all dimensions. In contrast, our FourierAttention compresses most dimensions to a fixed length while preserving long-context-sensitive dimensions, effectively reducing KV cache size while maintaining the original long-context capabilities.

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

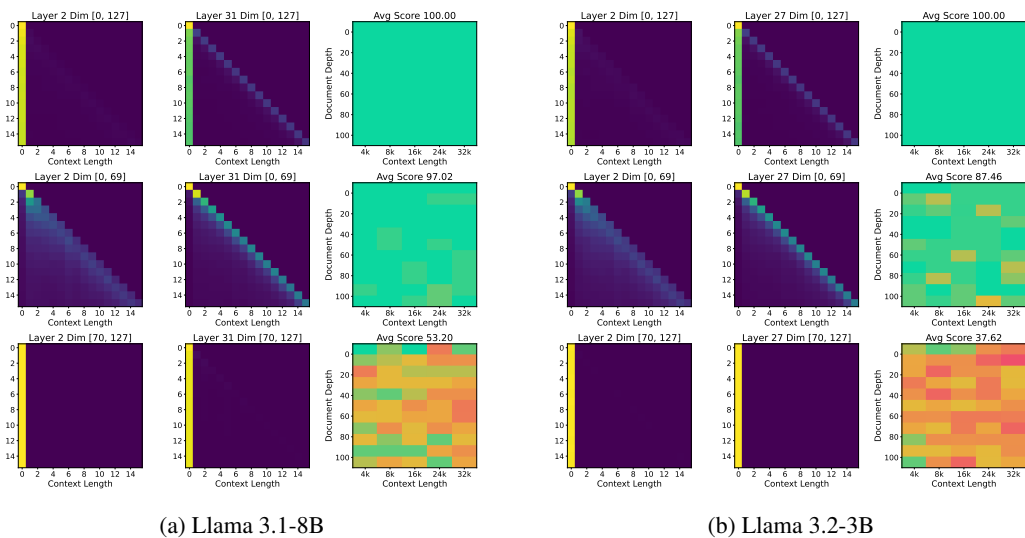


Figure 2: Visualization of the average attention score and its components in Llama 3.1-8B and Llama 3.2-3B over 32 sentences with 16 tokens. The component of the lower dimensions corresponds to the local branch in Xiao et al. (2024), while that of the upper dimensions corresponds to the global branch. This reveals the different functions of different dimensions in the attention mechanism. It can be further validated that adding Gaussian noise to the lower dimensions has little effect on NIAH performance, but adding noise to the upper dimensions will harm the performance remarkably.

3 METHODOLOGY

3.1 HEAD DIMENSION HETEROGENEITY

We analyze the heterogeneous sensitivity of transformer head dimensions to varying context lengths. By visualizing attention scores across 128 dimensions in LLaMA architecture (Figure.2), we identify a bifurcation in attention patterns: the first 70 dimensions (0–69) exhibit sharp focus on short-range context, with score distributions concentrated on recent tokens, while the latter 58 dimensions (70–127) maintain a persistent bias toward initial "sink tokens"—positional embeddings that serve as static reference points. This divergence suggests distinct contextual roles encoded within head dimensions, where specialized subsets prioritize local versus global signal retention.

To further validate this hypothesis, we evaluate the model on a Needle-In-A-Haystack retrieval task across sequences of up to 32,000 tokens. As shown in Figure 2a, the baseline model achieves perfect retrieval accuracy (100.0). Introducing Gaussian noise to the first 70 dimensions, confirming their limited role in long-range dependency resolution. Conversely, perturbing the latter 58 dimensions catastrophically reduces accuracy to 53.20 on average, with failures consistent across all tested depths and context lengths (Figure 2b mirrors this trend). This stark contrast empirically demonstrates that upper dimensions in transformers are indispensable for retaining long-range information, while lower dimensions specialize in local context encoding. These findings provide critical insights for optimizing memory-efficient architectures, as strategically prioritizing dimensions specialized in long-range retention enhances contextual awareness within memory limits. For more details on dimension selection, please refer to Section 3.4. Due to the heterogeneity further illustrated in Appendix B, we propose **FourierAttention**. In FourierAttention, most dimensions of the KV cache are compressed to a fixed length via translated Fourier transform.

3.2 PRELIMINARY: HIPPO FRAMEWORK

Inspired by HiPPO (Gu et al., 2020), we compress these less context-sensitive dimensions into fixed-length states to reduce KV cache storage. Under the HiPPO framework, an infinitely long sequence, $f_{1...L}$, can be approximated by finite-length states, $c \in \mathbb{R}^N$, as the combining coefficients of finite-order basis functions. HiPPO designs different state update equations for various basis

functions under different measure functions, such as LegT based on Legendre Polynomials in a translated fixed window size. Among these methods, FourierT measure based on Translated Fourier Transform is most suitable for token-wise parallelism in transformers, because it can be expressed in matrix form and performed independently in different order states. Therefore, we adopt FourierT to compress cache, $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{L \times d}$, which also achieves better downstream performance in Section 5.1.

3.3 ONLINE COMPRESSION VIA HIPPO-FOURIERT

We set the translated window length in FourierT to the maximum context length, ensuring effective compression within valid input-output ranges. In the prefilling phase, we preserve all dimensions of the initial L_{init} and the local L_{local} tokens,

$$\begin{aligned} \mathbf{K}^i, \mathbf{K}^l &= \mathbf{K}[:, L_{\text{init}}], \mathbf{K}[-L_{\text{local}} :] \\ \mathbf{V}^i, \mathbf{V}^l &= \mathbf{V}[:, L_{\text{init}}], \mathbf{V}[-L_{\text{local}} :] \end{aligned} \quad (1)$$

and distinguish the dimension indices $\mathcal{D}^{ku}, \mathcal{D}^{kc}, \mathcal{D}^{vu}, \mathcal{D}^{vc}$ in KV cache for uncompressing and compressing, to enable training-free integration.

$$\begin{aligned} \mathbf{K}^{mn} &= \mathbf{K}[L_{\text{init}} : -L_{\text{local}}, \mathcal{D}^{kc}], & \mathbf{V}^{mn} &= \mathbf{V}[L_{\text{init}} : -L_{\text{local}}, \mathcal{D}^{vc}], \\ \mathbf{K}^{mu} &= \mathbf{K}[L_{\text{init}} : -L_{\text{local}}, \mathcal{D}^{ku}], & \mathbf{V}^{mu} &= \mathbf{V}[L_{\text{init}} : -L_{\text{local}}, \mathcal{D}^{vu}]. \end{aligned} \quad (2)$$

We preserve $\mathbf{K}^{mu}, \mathbf{V}^{mu}$, compress $\mathbf{K}^{mn}, \mathbf{V}^{mn}$ to fix-sized $\mathbf{K}^{mc} \in \mathbb{R}^{2N \times |\mathcal{D}^{kc}|}$, $\mathbf{V}^{mc} \in \mathbb{R}^{2N \times |\mathcal{D}^{vc}|}$ and use the original KV for forward propagation.

$$\begin{aligned} \mathbf{K}^{mc} &= \frac{1}{T} \mathcal{F} \mathbf{K}^{mn}, & \mathbf{V}^{mc} &= \frac{1}{T} \mathcal{F} \mathbf{V}^{mn}, \\ \mathbf{O} &= \text{flash_attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}). \end{aligned} \quad (3)$$

We detail the mathematical derivation in Appendix C. When it comes to the compression matrix in FourierT, originally, $\mathcal{F} \in \mathbb{C}^{N \times L_{\text{middle}}}$, where $L_{\text{middle}} = L - L_{\text{init}} - L_{\text{local}}$ and $\mathcal{F}_{nt} = e^{i \frac{2\pi nt}{T}}$. However, since caches in mainstream LLMs are real-valued, we convert complex numbers to corresponding 2D vectors, transforming N -order complex states into $2N$ -order real states. We denote $N' = 2N$, and the real compression matrix in FourierT is $\mathcal{F} \in \mathbb{R}^{N' \times L_{\text{middle}}}$ as shown in Equation 4.

$$\mathcal{F} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 0 & 0 & \cdots & 0 \\ 1 & \cos \frac{2\pi}{T} & \cdots & \cos \frac{2\pi(L_{\text{middle}}-1)}{T} \\ 0 & \sin \frac{2\pi}{T} & \cdots & \sin \frac{2\pi(L_{\text{middle}}-1)}{T} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \cos \frac{2\pi(N-1)}{T} & \cdots & \cos \frac{2\pi(N-1)(L_{\text{middle}}-1)}{T} \\ 0 & \sin \frac{2\pi(N-1)}{T} & \cdots & \sin \frac{2\pi(N-1)(L_{\text{middle}}-1)}{T} \end{bmatrix}. \quad (4)$$

In the decoding phase, FourierAttention reconstructs intermediate cache $\tilde{\mathbf{K}}^m, \tilde{\mathbf{V}}^m$ via inverse Fourier transform in attention computation with the current query vector \mathbf{q}_{t+1} ,

$$\begin{aligned} \tilde{\mathbf{K}}^m[\mathcal{D}^{ku}] &\leftarrow \mathbf{K}^{mu}, & \tilde{\mathbf{K}}^m[\mathcal{D}^{kc}] &\leftarrow \mathcal{F}^T \mathbf{K}^{mc}, & \tilde{\mathbf{K}} &= \text{cat}(\mathbf{K}^i, \tilde{\mathbf{K}}^m, \mathbf{K}^l) \\ \tilde{\mathbf{V}}^m[\mathcal{D}^{vu}] &\leftarrow \mathbf{V}^{mu}, & \tilde{\mathbf{V}}^m[\mathcal{D}^{vc}] &\leftarrow \mathcal{F}^T \mathbf{V}^{mc}, & \tilde{\mathbf{V}} &= \text{cat}(\mathbf{V}^i, \tilde{\mathbf{V}}^m, \mathbf{V}^l) \\ \mathbf{o}_{t+1} &= \text{flash_attention}(\mathbf{q}_{t+1}, \tilde{\mathbf{K}}, \tilde{\mathbf{V}}). \end{aligned} \quad (5)$$

and compresses tokens out of the local range individually,

$$\begin{aligned} \mathbf{K}^{mc} &\leftarrow \mathbf{K}^{mc} + \frac{1}{T} \mathbf{f}_{L_{\text{middle}}+1} \mathbf{K}^l[0, \mathcal{D}^{kc}], & \mathbf{V}^{mc} &\leftarrow \mathbf{V}^{mc} + \frac{1}{T} \mathbf{f}_{L_{\text{middle}}+1} \mathbf{V}^l[0, \mathcal{D}^{vc}] \\ \mathbf{f}_{t+1} &= \begin{bmatrix} 0 & 1 & \cdots & \cos \frac{2\pi(N-1)L_{\text{middle}}}{T} & \sin \frac{2\pi(N-1)L_{\text{middle}}}{T} \end{bmatrix}^\top \end{aligned} \quad (6)$$

To eliminate intermediate read-write cost in decompression, we try to implement a custom kernel, **FlashFourierAttention**, using Triton (Tillet et al., 2019), integrating the decompression into standard FlashAttention2 (Dao, 2024) and FlashDecoding (Dao et al., 2023). FlashFourierAttention loads compressed intermediate states once and decompresses at corresponding sequence positions during iterative KV cache loading, which is further detailed in Appendix D.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

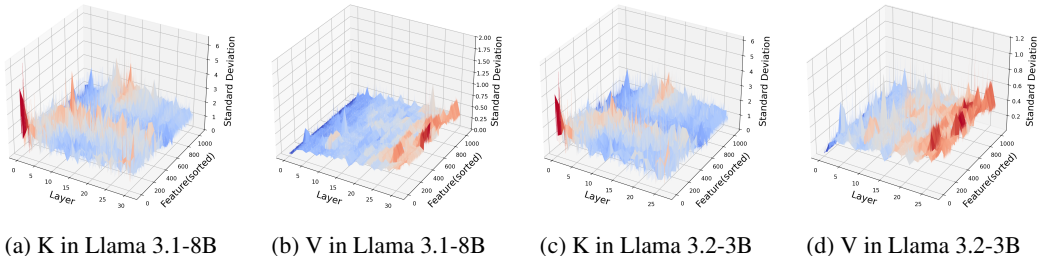


Figure 3: Visualization of standard deviation of KV cache in different layers in Llama 3.1-8B and Llama 3.2-3B. The feature dimensions are sorted based on the indices in each head.

3.4 FINE-GRAINED COMPRESSION SCHEMA

In FourierAttention, a crucial point lies in how to select the dimension to be compressed. To address this, we directly compress and decompress all KV caches, prioritizing dimensions with smaller mean-squared error in reconstruction to a fixed length. Based on further observations of the KV cache, we adopt a fine-grained compression schema, where more dimensions of the V cache and lower-layer caches are compressed to a fixed length. We analyze the standard deviation of KV cache dimensions along the temporal direction across different layers and find that for both Llama 3.1-8B and Llama 3.2-3B, as shown in Figure 7. The standard deviation of the K cache is consistently higher than that of the V cache, and the standard deviation in upper layers exceeds that of lower layers. Consequently, we compress more dimensions of the smoother V cache and lower-layer caches to a fixed length, while retaining more K cache and upper-layer caches to extend with sequence length. Thus, FourierAttention exhibits an asymmetric, inverted-pyramid compression pattern.

Interestingly, this differs from most KV cache compression approaches. Works like Cai et al. (2024) and Xing et al. (2024) suggest preserving more KV caches in lower layers, as attention becomes sparser in upper layers. However, in FourierAttention, the optimization criterion is whether the dimension can be well reconstructed. Since caches in upper layers exhibit more oscillatory features due to more deterministic predictions, we retain more dimensions to maintain output stability.

4 EXPERIMENT

4.1 SETUP

We conduct experiments on Llama 3.1-8B (Dubey et al., 2024) and Llama 3.2-3B (Meta, 2024a). For all models, we set the length of initial tokens L_{init} to 4, the length of local tokens L_{local} to 1024, and the number of states $N = 512$, namely $N' = 1024$. We evaluate the reconstruction loss using the prompt portion of the 32k Needle-In-A-Haystack benchmark in OpenCompass (Contributors, 2023). As mentioned earlier, we employ an asymmetric inverted pyramid compression strategy: for the first 4 layers, we compress 90% of K dimensions and 95% of V dimensions; for the last 8 layers, 50% of K and 70% of V; and for the remaining layers, 80% of both K and V. Overall, 76% KV caches are compressed to a fixed length. All experiments are performed on an NVIDIA H100 GPU with FP16 precision and accelerated with FlashAttention2 (Dao, 2024).

4.2 LONG-CONTEXT EVALUATION

We evaluate our method against other KV cache optimization approaches with two long-context benchmarks in OpenCompass (Contributors, 2023), LongBench (Bai et al., 2023), and Needle-In-A-Haystack (NIAH) (Kamradt, 2023; Li et al., 2024a), with a truncation context length of 32K. We compare with StreamingLLM (Xiao et al., 2024), SnapKV (Li et al., 2024b), PyramidKV (Cai et al., 2024), and Palu (Chang et al., 2024), covering both token eviction and feature compression. For fair comparison, we retain 4 initial tokens and 1024 local tokens in StreamingLLM, additionally keep 1024 recalled middle tokens, matching our compressed dimension count, in SnapKV and PyramidKV, and compress KV feature dimensions to 70% in Palu.

	Single-Doc			Multi-Doc			Summary			Few-shot			Synthetic			Code		Avg.
	NQ	Qsp	MF	HQ	WQ	Msq	GR	QS	MN	TR	TQ	SS	PC	PR	LCC	Re-P		
Llama-3.1-8B	13.2	20.2	32.8	12.0	13.6	8.7	29.7	25.1	0.9	73.5	91.0	47.3	0.8	26.8	72.0	69.3	39.5	
+ SLM	7.9	13.9	15.6	7.8	10.1	4.5	19.9	21.5	9.9	61.5	84.7	43.5	1.3	6.2	58.4	56.4	31.5	
+ SnapKV	12.7	19.8	32.5	12.0	13.8	8.6	29.2	24.9	12.6	73.0	91.0	46.5	0.8	26.8	60.0	59.7	37.1	
+ PyramidKV	18.5	19.8	32.5	12.1	13.8	8.7	29.7	24.9	12.4	73.0	90.1	46.6	0.8	26.8	60.0	59.4	<u>37.3</u>	
+ Palu	4.5	18.0	21.6	9.5	11.3	5.2	17.3	6.9	9.0	68.5	83.4	32.3	0.6	14.6	56.7	54.6	30.7	
+ FA (ours)	15.6	19.8	32.7	12.0	13.5	7.9	24.1	24.0	0.7	73.0	91.2	46.0	1.1	26.9	71.4	66.3	38.6	
Llama-3.2-3B	10.3	21.7	35.5	9.6	12.8	6.8	30.2	23.8	28.2	70.0	87.2	38.2	0.0	7.0	70.0	66.4	38.0	
+ SLM	9.1	17.6	21.6	7.1	9.8	4.0	19.0	21.3	23.2	53.0	84.3	39.8	1.4	6.5	55.5	53.5	31.2	
+ SnapKV	9.4	21.0	35.0	9.5	12.8	6.6	29.5	23.4	27.8	69.5	86.4	38.3	0.0	6.8	58.2	56.4	34.9	
+ PyramidKV	8.9	21.4	35.7	9.5	12.8	6.8	30.4	23.6	28.2	69.5	86.9	38.7	0.0	6.8	58.7	57.3	<u>35.2</u>	
+ Palu	2.0	19.2	20.4	5.8	10.3	2.7	13.4	4.1	14.0	57.0	47.4	21.5	1.5	3.1	55.1	49.5	25.5	
+ FA (ours)	12.8	21.1	35.8	9.8	11.6	6.5	23.0	23.8	24.1	69.0	87.0	39.4	0.0	6.2	69.8	63.3	37.0	

Table 1: Results of LLaMA Series (Dubey et al., 2024; Meta, 2024b) on LongBench (Bai et al., 2023). Our FourierAttention (FA) achieves consistent superiority over StreamingLLM (SLM), SnapKV, PyramidKV, and Palu and shows the closest performance with LLMs with full attention.

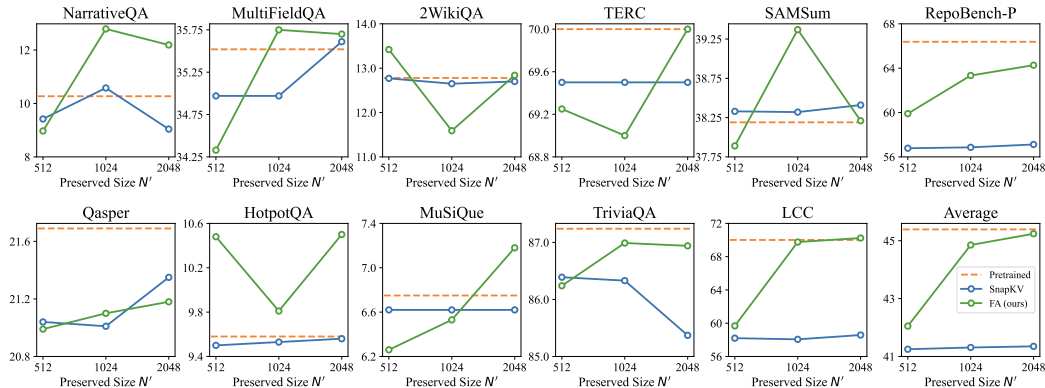


Figure 4: Comparison of different preserved sizes N' between FourierAttention and SnapKV in Llama 3.2-3B (Meta, 2024a). FourierAttention outperforms SnapKV on average across different N' .

For LongBench as shown in Tables 1, our FourierAttention achieves performance closest to the original model and shows consistent superiority over other cache optimization methods on both Llama 3.2-3B and Llama 3.1-8B. For the NIAH task as shown in Figure 5 and 6, we similarly achieve performance closest to the original pretrained models at 32k context length. While SnapKV and PyramidKV are theoretically suitable for retrieval tasks like NIAH, they still exhibit recall errors. Though Palu maintains stable attention approximation under moderate compression, 30-50%, they show significant performance degradation at 75% compression due to insufficient granular analysis of KV cache features. In contrast, our FourierAttention optimizes compression by identifying and preserving KV dimensions insensitive to compression, thereby maximally retaining the long-context capabilities and demonstrating superiority across both models and benchmarks.

On Llama 3.2-3B, we also verify the downstream task performance of FourierAttention with different values of N' and SnapKV with the corresponding intermediate recall sizes, as shown in Figure 4. We find that, regardless of the value of N' , FourierAttention exhibits sufficient robustness and outperforms SnapKV on average across different recall sizes. We chose $N' = 1024$ in our paper, which is a parameter that balances performance and speed.

4.3 EFFICIENCY VALIDATION

In addition to comparisons in downstream performance, we also conduct an efficiency comparison of FourierAttention with SnapKV and Palu, two representative cache optimization methods. First, in

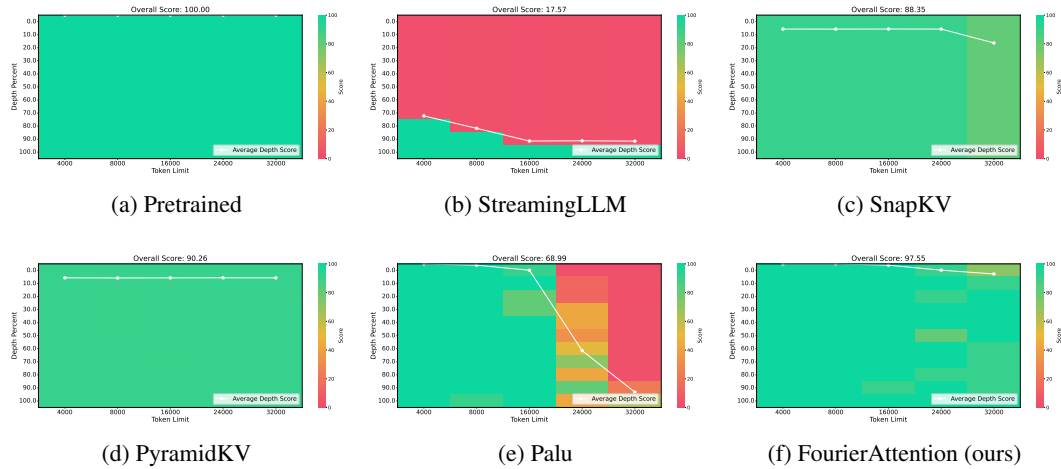


Figure 5: Results of Llama 3.1-8B (Dubey et al., 2024) on Needle-In-A-Haystack (Kamradt, 2023). FourierAttention achieves the highest average score over StreamingLLM, SnapKV, PyramidKV, and Palu, and shows the closest performance with LLMs with full attention.

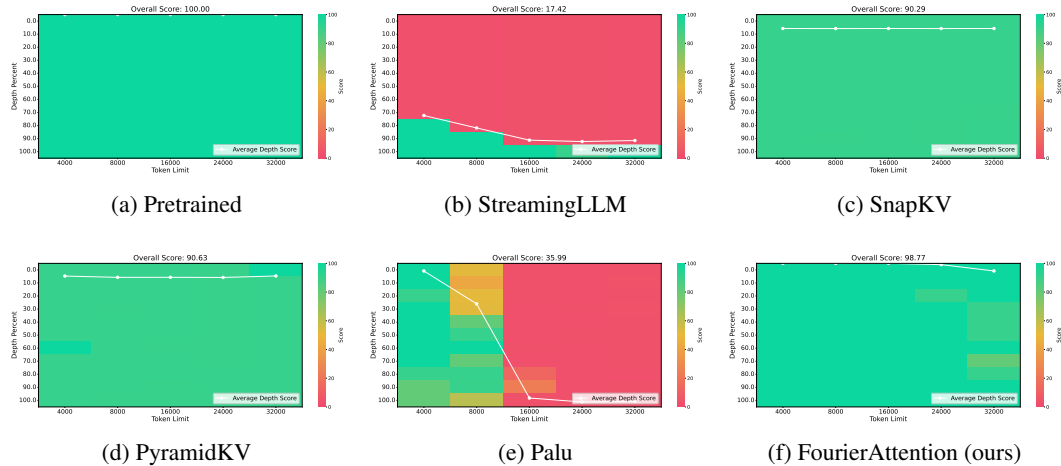


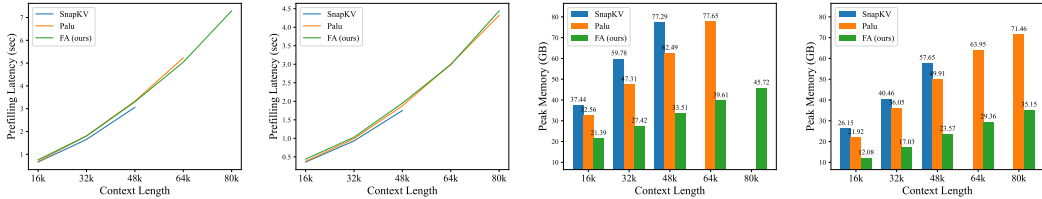
Figure 6: Results of Llama 3.2-3B (Meta, 2024a) on Needle-In-A-Haystack (Kamradt, 2023). FourierAttention achieves the highest average score over StreamingLLM, SnapKV, PyramidKV, and Palu, and shows the closest performance with LLMs with full attention.

terms of storage, thanks to the custom Triton kernel, FlashFourierAttention, detailed in Appendix D, compared with SnapKV and Palu, we achieve a clear advantage in memory efficiency and 80k context length inference of Llama 3.1-8B and Llama 3.2-3B on a single H100. In the prefilling phase, since we only add one Fourier transform, matrix multiplication, compared with the original attention and use the efficient FFT built-in operator, we have achieved a latency close to that of existing efficient approaches such as Palu and SnapKV in long contexts ranging from 16k to 80k. In addition, thanks to the accelerated operators during the decoding phase, we achieve throughput comparable to SnapKV.

5 DISCUSSION

5.1 CHOICE OF BASIS FUNCTIONS

Although FourierAttention employs HiPPO-FourierT for compression, Gu et al. (2020) proposes and claims polynomial basis functions like LegT with superior performance. While maintaining identical sliding window sizes, we compare LegT and FourierT in reconstructing KV caches from LLMs. As illustrated in Figure 8, we evaluate their reconstruction effects on 4 randomly selected KV cache



(a) Llama 3.1-8B Latency (b) Llama 3.2-3B Latency (c) Llama 3.1-8B Memory (d) Llama 3.2-3B Memory

Figure 7: Visualization of standard deviation of KV cache in different layers in Llama 3.1-8B and Llama 3.2-3B. The feature dimensions are sorted based on the indices in each head.

	NIAH				AGG				QA		Avg.			
	SK1	SK2	SK3	MK1	MK2	MK3	MV	MQ	CWE	FWE		VT	SQ	HP
Llama 3.2-3B	100.0	100.0	100.0	99.0	100.0	99.0	100.0	99.8	60.0	89.7	97.0	77.0	53.0	90.3
+ FourierT	100.0	100.0	98.0	99.0	99.0	100.0	96.0	97.8	57.9	80.7	76.6	81.0	51.0	87.5
+ LegT	55.0	82.0	42.0	89.0	93.0	50.0	93.8	84.0	13.5	58.3	76.0	36.0	21.0	61.0
+ uniform	100.0	100.0	99.0	99.0	99.0	98.0	93.0	98.8	59.7	77.3	78.2	80.0	50.0	87.1
+ KV inv.	100.0	100.0	100.0	99.0	99.0	100.0	93.5	98.8	54.3	79.7	72.6	80.0	53.0	86.9
+ layer inv.	100.0	98.0	88.0	98.0	97.0	94.0	80.0	96.8	41.3	69.3	61.4	80.0	50.0	81.1

Table 2: Validation of basis function and compression schema in Llama 3.2-3B based on RULER in 4k context length. Besides AGG tasks, FourierAttention achieve performance close to original model.

dimensions from layer 0 of Llama 3.2. Under equivalent state dimensions¹, FourierT consistently achieves lower reconstruction loss than LegT.

We further evaluate FourierT and LegT compression on Llama 3.2-3B using more discriminative RULER benchmark (Hsieh et al., 2024) in 4k context length. For fair comparison, we employ the same method to identify dimensions suitable for LegT compression and apply an identical compression schema. Results in Table 2 show FourierT still performs better, demonstrating that FourierT offers better parallelizability for compression efficiency and performance in downstream evaluation. Nevertheless, we must acknowledge that FourierAttention performs slightly behind the pre-trained model on aggregation (AGG) tasks, such as VT. We attribute this to the fact that aggregation tasks rely on statistical information and demand finer input details. By discarding high frequencies during compression, FourierAttention loses these details and thus degrades on such tasks.

5.2 ABLATION ON COMPRESSION SCHEMA

As mentioned in Section 3.4, we propose a more fine-grained compression scheme based on additional observations of the KV cache. As shown in the Table 2, we compare three approaches: uniform compression across all layers and between KV (uniform), inverted KV compression schema by K-priority over V (KV inv.), and inverted layer-wise compression schema by upper-layer priority over lower-layer (layer inv.). Results demonstrate that our original V-priority and lower-layer-priority compression schema achieves superior performance on RULER in a 4k context length. This further illustrates that frequency-based sequence-wise KV cache compression exhibits different optimization characteristics compared to conventional KV token eviction (Cai et al., 2024; Xing et al., 2024).

5.3 COMPRESSED DIMENSION DISTRIBUTION

Finally, we analyze the compressed dimensions selected by our FourierAttention. We count the number of each dimension selected for compression, averaged across attention heads in different layers, grouped every 16 dimensions. Results in Figure 9 show that in both Llama 3.1-8B and Llama 3.2-3B, starting from layer 2, lower dimensions are more frequently compressed while the upper dimensions tend to preserve complete temporal information in our FourierAttention. This

¹FourierT uses lower-order basis functions since FourierAttention’s state size is twice the number of states

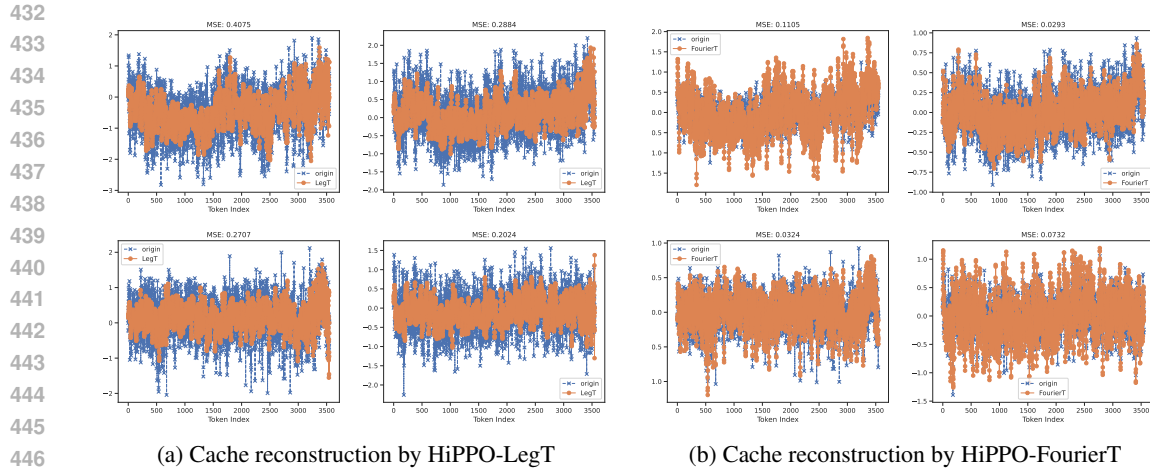


Figure 8: Visualization of KV cache reconstruction in Llama 3.2-3B for different basis functions, LegT and FourierT under HiPPO framework (Gu et al., 2020). FourierT outperforms LegT in cache reconstruction with lower reconstruction loss.

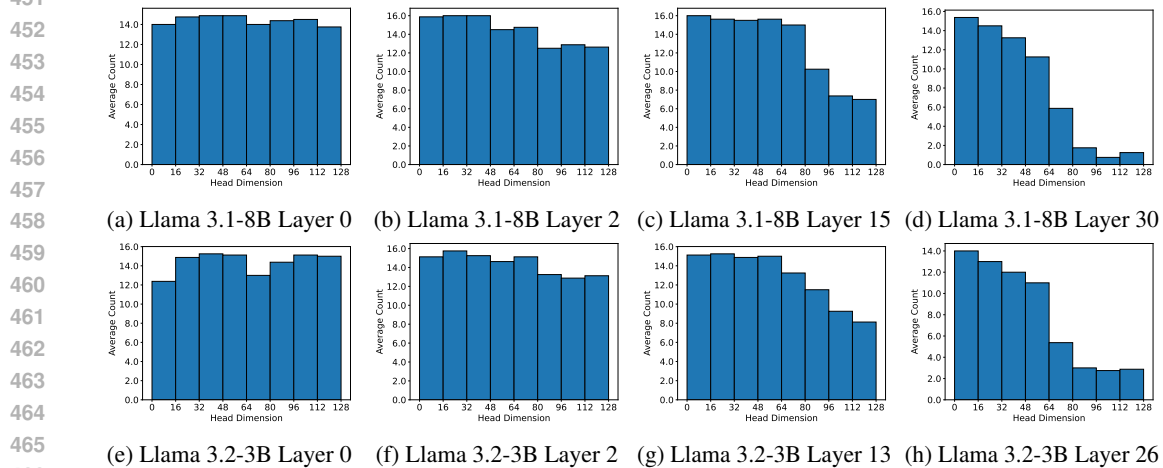


Figure 9: The statistics of each dimension selected for compression, averaged across attention heads in different layers, grouped every 16 dimensions, in Llama 3.1-8B and Llama 3.2-3B.

phenomenon is more evident in upper layers, where fewer dimensions are chosen to be compressed. As illustrated in Figure 2, these uncompressed upper dimensions primarily contribute to forming attention sinks and capturing long-context semantic relationships, thus requiring complete retention, whereas other dimensions can be stored with limited length.

6 CONCLUSION

We propose FourierAttention, a novel KV cache optimization approach that compresses long-context-insensitive dimensions without sacrificing contextual awareness based on an interesting phenomenon in transformer head dimensions, that lower dimensions capture local features, while upper ones capture long-context dependencies. Inspired by HiPPO, we optimize the long-context-insensitive KV cache through a translated Fourier transform into fixed-length states in the prefilling phase and reconstruct the KV cache in the decoding phase. FourierAttention shows the best performance on the LLaMA Series in LongBench and NIAH on average. We are trying to improve the efficiency of FourierAttention through a customized Triton-based kernel, FlashFourierAttention, eliminating intermediate read-write operations and effectively reducing memory overhead.

486 ETHICAL STATEMENT
487

488 This research adheres to established ethical standards. To the best of our knowledge, our study does
489 not process any sensitive personal data, does not involve any human subjects, and does not target
490 any ethically risky applications. All experiments and analyses are conducted in line with recognized
491 guidelines, ensuring integrity, transparency, and reliability.

492
493 REPRODUCIBILITY STATEMENT
494

495 To ensure the reproducibility of and to support the open-source community, we will publicly release
496 FourierAttention, its trained checkpoints, and the complete training and evaluation code, especially
497 the custom Triton kernel. We expect these as a reference for future work on long-context LLMs,
498 facilitating innovation and advancing progress in this field.

500 REFERENCES
501

502 Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit
503 Shanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints.
504 *arXiv preprint arXiv:2305.13245*, 2023.

505
506 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao
507 Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long context
508 understanding. *arXiv preprint arXiv:2308.14508*, 2023.

509 Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong,
510 Baobao Chang, Junjie Hu, et al. Pyramidkv: Dynamic kv cache compression based on pyramidal
511 information funneling. *arXiv preprint arXiv:2406.02069*, 2024.

512
513 Chi-Chih Chang, Wei-Cheng Lin, Chien-Yu Lin, Chong-Yan Chen, Yu-Fang Hu, Pei-Shuo Wang,
514 Ning-Chi Huang, Luis Ceze, Mohamed S Abdelfattah, and Kai-Chiang Wu. Palu: Compressing
515 kv-cache with low-rank projection. *arXiv preprint arXiv:2407.21118*, 2024.

516 OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models.
517 <https://github.com/open-compass/opencompass>, 2023.

518
519 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The*
520 *Twelfth International Conference on Learning Representations*, 2024.

521 Tri Dao, Daniel Haziza, Francisco Massa, and Grigory Sizov. Flash-decoding for long-context in-
522 ference., 2023. URL [https://crfm.stanford.edu/2023/10/12/flashdecoding.](https://crfm.stanford.edu/2023/10/12/flashdecoding.html)
523 [html](https://crfm.stanford.edu/2023/10/12/flashdecoding.html).

524
525 Haojie Duanmu, Zhihang Yuan, Xiuhong Li, Jiangfei Duan, Xingcheng Zhang, and Dahua Lin.
526 Skvq: Sliding-window key and value cache quantization for large language models. *arXiv preprint*
527 *arXiv:2405.06219*, 2024.

528 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
529 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
530 *arXiv preprint arXiv:2407.21783*, 2024.

531
532 Yao Fu. Challenges in deploying long-context transformers: A theoretical peak performance analysis.
533 *arXiv preprint arXiv:2405.08944*, 2024.

534 Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory
535 with optimal polynomial projections. *Advances in Neural Information Processing Systems*, 33:
536 1474–1487, 2020.

537
538 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
539 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- 540 Ziwei He, Meng Yang, Minwei Feng, Jingcheng Yin, Xinbing Wang, Jingwen Leng, and Zhouhan
541 Lin. Fourier transformer: Fast long range modeling by removing sequence redundancy with fft
542 operator. *arXiv preprint arXiv:2305.15099*, 2023.
- 543 Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao,
544 Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with
545 kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024.
- 547 Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang
548 Zhang, and Boris Ginsburg. Ruler: What’s the real context size of your long-context language
549 models? *arXiv preprint arXiv:2404.06654*, 2024.
- 550 Greg Kamradt. Needle in a haystack - pressure testing llms. [https://github.com/
551 gkamradt/LLMTest_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack), 2023.
- 553 Mo Li, Songyang Zhang, Yunxin Liu, and Kai Chen. Needlebench: Can llms do retrieval and
554 reasoning in 1 million context window? *arXiv preprint arXiv:2407.11963*, 2024a.
- 555 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai,
556 Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before generation.
557 *arXiv preprint arXiv:2404.14469*, 2024b.
- 559 Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong
560 Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-
561 experts language model. *arXiv preprint arXiv:2405.04434*, 2024a.
- 562 Xiaoran Liu, Hang Yan, Chenxin An, Xipeng Qiu, and Dahua Lin. Scaling laws of rope-based
563 extrapolation. In *The Twelfth International Conference on Learning Representations*, 2024b.
- 565 Xiaoran Liu, Ruixiao Li, Mianqiu Huang, Zhigeng Liu, Yuerong Song, Qipeng Guo, Siyang He, Qiqi
566 Wang, Linlin Li, Qun Liu, et al. Thus spake long-context large language model. *arXiv preprint
567 arXiv:2502.17129*, 2025.
- 568 Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi
569 Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint
570 arXiv:2402.02750*, 2024c.
- 572 Xin Men, Mingyu Xu, Bingning Wang, Qingyu Zhang, Hongyu Lin, Xianpei Han, and Weipeng
573 Chen. Base of rope bounds context length. *arXiv preprint arXiv:2405.14591*, 2024.
- 574 AI Meta. Introducing meta llama 3: The most capable openly available llm to date. *Meta AI.*, 2024a.
- 576 AI Meta. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI.*,
577 2024b.
- 578 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 580 OpenAI. O1: Openai’s first model, 2024. URL <https://openai.com/o1/>. Accessed: 2024-
581 12-25.
- 582 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context win-
583 dow extension of large language models. In *The Twelfth International Conference on Learning
584 Representations*, 2024.
- 586 Utkarsh Saxena, Gobinda Saha, Sakshi Choudhary, and Kaushik Roy. Eigen attention: Attention in
587 low-rank space for kv cache compression. *arXiv preprint arXiv:2408.05646*, 2024.
- 588 Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li, Qinyuan Cheng, Xiangyang Liu, Hang Yan,
589 Yunfan Shao, Qiong Tang, Shiduo Zhang, Xingjian Zhao, Ke Chen, Yining Zheng, Zhejian Zhou,
590 Ruixiao Li, Jun Zhan, Yunhua Zhou, Linyang Li, Xiaogui Yang, Lingling Wu, Zhangyue Yin,
591 Xuanjing Huang, Yu-Gang Jiang, and Xipeng Qiu. Moss: An open conversational large language
592 model. *Machine Intelligence Research*, 2024. ISSN 2731-5398. doi: 10.1007/s11633-024-1502-8.
593 URL <https://github.com/OpenMOSS/MOSS>.

594 Jamba Team, Barak Lenz, Alan Arazi, Amir Bergman, Avshalom Manevich, Barak Peleg, Ben
595 Aviram, Chen Almagor, Clara Fridman, Dan Padnos, et al. Jamba-1.5: Hybrid transformer-mamba
596 models at scale. *arXiv preprint arXiv:2408.12570*, 2024.

597
598 Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler
599 for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International*
600 *Workshop on Machine Learning and Programming Languages*, pp. 10–19, 2019.

601 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
602 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Ad-*
603 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
604 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
605 [file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).

606
607 Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Jian Tong, Haodong
608 Duan, Qipeng Guo, Jiaqi Wang, et al. Videorope: What makes for good video rotary position
609 embedding? *arXiv preprint arXiv:2502.05173*, 2025.

610 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
611 language models with attention sinks. In *The Twelfth International Conference on Learning*
612 *Representations*, 2024.

613
614 Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui
615 He, Jiaqi Wang, Feng Wu, et al. Pyramiddrop: Accelerating your large vision-language models via
616 pyramid visual redundancy reduction. *arXiv preprint arXiv:2410.17247*, 2024.

617 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
618 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*,
619 2025.

620
621 Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie,
622 Y. X. Wei, Lean Wang, Zhiping Xiao, Yuqing Wang, Chong Ruan, Ming Zhang, Wenfeng Liang,
623 and Wangding Zeng. Native sparse attention: Hardware-aligned and natively trainable sparse
624 attention. *arXiv preprint arXiv:2502.11089*, 2025.

625
626 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,
627 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient
628 generative inference of large language models. *Advances in Neural Information Processing*
629 *Systems*, 36:34661–34710, 2023.

630 A USE OF LARGE LANGUAGE MODELS

631
632 We only use Large Language Models for language-centric assistance, namely, for grammar, style,
633 and clarity, ensuring that no component of research ideation, experimental design, or scientific
634 contribution is either influenced or generated by LLM outputs.

635 B MORE OBSERVATION ON HETEROGENEITY

636
637 Regarding the observation on K cache, there exists a relation between this dimension bifurcation with
638 RoPE. This is a very correct insight. In Figure 2, we compare the attention score components of the
639 first 70 dimensions and the last 58 dimensions of all QK states in Llama 3.1-8B and Llama 3.2-3B,
640 as well as the effects on the downstream NIAH task after adding noise. The choice of the first 70
641 dimensions refers to the concept of critical dimension in Liu et al. (2024b), which corresponds to the
642 number of dimensions with complete position information observed during the pre-training phase in
643 LLMs with RoPE. This dimension d_{extra} can be calculated based on the head dimension d , the initial
644 pre-training length T_{train} , and the initial rotary base β .

$$645 \quad d_{\text{extra}} = 2 \left\lceil \frac{d}{2} \log_{\beta} \frac{T_{\text{train}}}{2\pi} \right\rceil$$

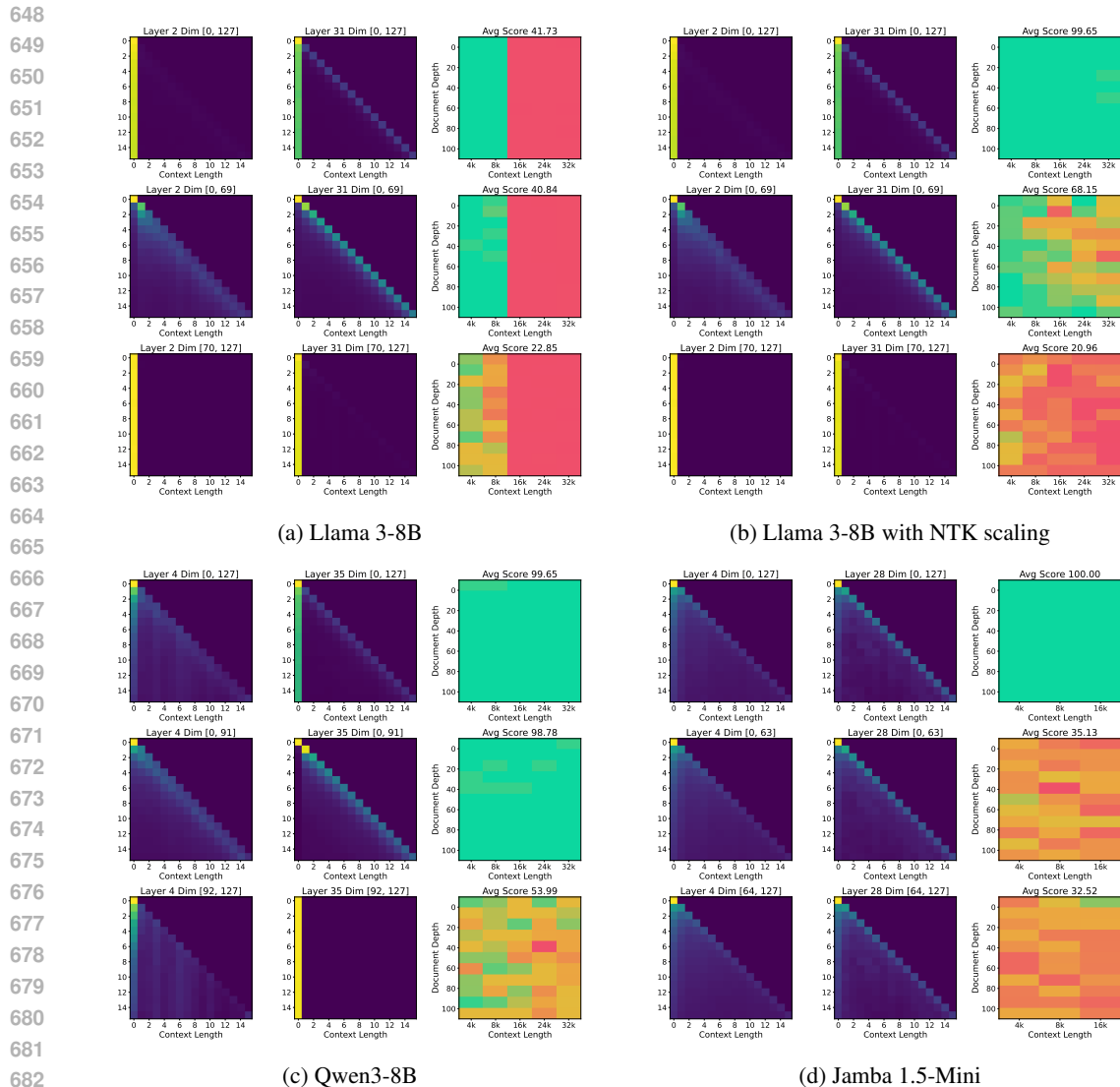


Figure 10: Visualization of the average attention score and its components in other RoPE-based LLMs, including Llama 3-8B (Meta, 2024b), as well as its NTK-scaled model, and Qwen3-8B, (Meta, 2024a), and NoPE-based LLMs, such as Jamba 1.5-Mini (Team et al., 2024). In RoPE-based LLMs, the component of the lower dimensions corresponds to the local branch in StreamingLLM (Xiao et al., 2024), while that of the upper dimensions corresponds to the global branch. This reveals the different functions of different dimensions in the attention mechanism. It can be further validated that adding Gaussian noise to the lower dimensions has little effect on NIAH performance, but adding noise to the upper dimensions will harm the performance remarkably. In NoPE-based LLMs, the components of the lower and upper dimensions do not have heterogeneous features in RoPE-based LLMs. Besides, adding Gaussian noise to the lower or upper dimensions will harm the NIAH performance equally.

Some studies have found that the self-attention components before and after the critical dimension have different characteristics in long contexts (Liu et al., 2024b; Wei et al., 2025), and when the rotation angle is large, LLMs are not good at characterizing long-context features (Men et al., 2024). Therefore, we choose the critical dimension as the split and find that different dimensions before and after have different impacts on long-context tasks. For LLaMA3 as well as Llama 3.1 and Llama 3.2, $d = 128$, $T_{\text{train}} = 8192$, $\beta = 500000$, and the calculated critical dimension size is 70.

We have also observed similar conclusions on Llama 3-8B with short contexts as shown in Figure 10a. For the first 70 dimensions and the last 58 dimensions of Llama 3-8B, adding noise to the first 70

702 dimensions has almost no impact, while adding noise to the last 58 dimensions leads to a significant
 703 drop in performance within the pre-training length (8k). As for the case of a larger rotation angle
 704 base, we scale up the rotary base of Llama 3-8B by $13\times$ to support a 32k context. Similarly, as shown
 705 in Figure 10b, for the NTK-extrapolated Llama 3-8B, adding noise to the first 70 dimensions and the
 706 last 58 dimensions, respectively, we find that noise in the first 70 dimensions has a weaker impact,
 707 while noise in the last 58 dimensions has a greater impact.

708 We have also verified our observations on the Qwen3-3B (Yang et al., 2025), as shown in Figure 10c.
 709 Since the initial pre-training length and rotation angle base of Qwen3 are 4096 and 10000, respectively,
 710 the critical dimension $d_{\text{extra}} = 92$. Then we add noise to the first 92 dimensions and the last 36
 711 dimensions, respectively, and find that noise in the upper dimensions has a greater impact on NIAH.

712 To prove that RoPE indeed causes the above phenomenon, we also conduct experiments on the hybrid
 713 model Jamba (Team et al., 2024) without position encoding, as shown in Figure 10d. Since only 1/8
 714 of the layers in Jamba use self-attention, we increase the noise level to $\mathcal{N}(0, 16)$. We find that for the
 715 first 64 dimensions and the last 64 dimensions, adding noise respectively produces effects that are
 716 quite close, and not as significant as in LLMs with RoPE. At the same time, we have also observed
 717 that in the above RoPE-based LLM, the upper dimensions contribute to attention sink, while the
 718 lower dimensions contribute to local attention. These phenomena have not been presented on Jamba.

720 C MATHEMATICAL DERIVATION

721
 722 We will present the derivation process of HiPPO-FourierT in this section. Firstly, according to the
 723 HiPPO framework (Gu et al., 2020), we give the definitions of the basis functions and the measure
 724 functions. We use the Fourier bases as the basis functions $g_n(x, t)$, and adopt the translated average
 725 measure $\mu(x, t)$ within a fixed window. The formulas are as follows, where the length of the fixed
 726 window is equal to the maximum context length T supported by the model.

$$727 \quad g_n(x, t) = e^{i\frac{2\pi nx}{T}}, \quad \mu(x, t) = \frac{1}{T} \mathbb{I}_{[t-T, t]} \quad (7)$$

729 Based on this, the input signal $f(t)$, which is the KV cache to be compressed, can be obtained as the
 730 projection $c_n(t)$ of the basis functions $g_n(x, t)$ under the measure $\mu(x, t)$.

$$731 \quad c_n(t) = \langle f_{\leq t}(x), g_n(x, t) \rangle_{\mu(x, t)} \\
 732 \quad = \int_0^t f(x) \cdot g_n(x, t) \cdot \mu(x, t) dx \quad (8)$$

736 Then, we differentiate the state $c_n(t)$.

$$737 \quad \frac{dc_n(t)}{dt} = \int_0^t f(x) \cdot \frac{\partial g_n(x, t)}{\partial t} \cdot \mu(x, t) dx + \\
 738 \quad \int_0^t f(x) \cdot g_n(x, t) \cdot \frac{\partial \mu(x, t)}{\partial t} dx \quad (9)$$

744 Since $\frac{\partial g_n(x, t)}{\partial t} = 0$ and $\frac{\partial \mu(x, t)}{\partial t} = \frac{\delta_t}{T} - \frac{\delta_{t-T}}{T}$, the differentiation of $c_n(t)$ is simplified as follows.

$$745 \quad \frac{dc_n(t)}{dt} = 0 + \int_0^t f(x) \cdot g_n(x, t) \cdot \frac{\delta_t}{T} dx - 0 \\
 746 \quad = \frac{f(t)}{T} \cdot e^{i\frac{2\pi nt}{T}} \quad (10)$$

750 Considering that the actual storage of LLM is real numbers, we calculate the derivatives of the real
 751 and imaginary parts, respectively.

$$752 \quad \text{Re} \left[\frac{dc_n(t)}{dt} \right] = \frac{f(t)}{T} \cos \frac{2\pi nt}{T}, \\
 753 \quad \text{Im} \left[\frac{dc_n(t)}{dt} \right] = \frac{f(t)}{T} \sin \frac{2\pi nt}{T} \quad (11)$$

After discretizing them respectively, we obtain the final state update equations. Regarding the choice of discretization strategy, since the \mathbf{A} matrix in HiPPO-FourierT is an identity matrix, the results obtained by different discretization methods in HiPPO-FourierT are only different in step size. We choose the simplest forward Euler discretization as follows.

$$c_{t+1}^{(n)} = \begin{cases} c_t^{(n)} + \frac{f_t}{T} \cos \frac{2\pi nt}{T} & n = 2m \\ c_t^{(n)} + \frac{f_t}{T} \sin \frac{2\pi nt}{T} & n = 2m + 1 \end{cases} \quad (12)$$

$$m = 0, 1, \dots, N - 1$$

Then we can derive the compression, update, and decompression functions as shown in Equation 4 in our paper,

$$c_{t+1} = \frac{1}{T} \mathcal{F}_{N \times t} f_t$$

$$\mathcal{F}_{k \times t} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 0 \\ 1 & \cos \frac{2\pi}{T} & \dots & \cos \frac{2\pi t}{T} \\ 0 & \sin \frac{2\pi}{T} & \dots & \sin \frac{2\pi t}{T} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \cos \frac{2\pi(N-1)}{T} & \dots & \cos \frac{2\pi(N-1)t}{T} \\ 0 & \sin \frac{2\pi(N-1)}{T} & \dots & \sin \frac{2\pi(N-1)t}{T} \end{bmatrix}, \quad (13)$$

where

$$c_{t+1} = [c_{t+1}^{(0)}, c_{t+1}^{(1)}, \dots, c_{t+1}^{(2k-2)}, c_{t+1}^{(2k-1)}]^\top$$

$$f_t = [f_0, f_1, \dots, f_t]^\top$$

Regarding discretization, since the \mathbf{A} matrix in HiPPO-FourierT is a zero matrix, the results obtained by different discretization methods in HiPPO-FourierT are only different in step size. We finally chose the simplest forward Euler discretization. In the comparative experiments, we have compared HiPPO-LegT with the bilinear discretization method, which is the best-performing discretization method reported in Gu et al. (2020).

D TRITON IMPLEMENTAION DETAILS

Regarding Triton implementation details, since PyTorch itself already provides efficient support for Fourier transforms, we only need to design and integrate a custom kernel for the decoding stage, focusing on solving the decompression calculations involved in the decoding stage, and wrapping the decompression within the attention calculations to avoid the read-write of full-size KV cache. In this process, we applied three techniques: FlashDecoding, dimension reordering, and compression delay.

FlashDecoding is a fundamental method for decoding acceleration (Dao et al., 2023). Building on the iterative calculation of self-attention in FlashAttention, it targets the parallel computation of attention in the decoding phase. During decoding, it first segments and computes in parallel, and within each segment, it calculates the attention numerators, denominators, and max attention scores in a block-cycling manner. Finally, it merges the results from different segments. Since the inverse Fourier transform is essentially a matrix multiplication that can be split and computed in parallel along the sequential dimension, the decompression logic of FourierAttention can naturally be adapted to FlashDecoding. Each segment loads the complete compressed state, first decompresses the K cache within the corresponding block, concatenates it with the retained part, multiplies QK, and then gets the attention distribution. After that, it decompresses the corresponding V cache, concatenates it with the retained part, and outputs the attention result within the block.

Secondly, the dimension indices of the KV cache that need to be compressed and those that need to be fully retained, denoted as \mathcal{D}^{ku} , \mathcal{D}^{kc} , \mathcal{D}^{vu} , \mathcal{D}^{vc} , are theoretically discrete. However, to improve the loading speed, following the approach in Duanmu et al. (2024), we reorder the dimensions of the QK and VO matrices according to these indices. Specifically, the dimensions in \mathcal{D}^{ku} , \mathcal{D}^{vu} that do not need to be compressed are placed at the lower part of the QK features and VO features of each head, respectively. Conversely, the dimensions in \mathcal{D}^{kc} , \mathcal{D}^{vc} that need to be compressed are placed at the

810 upper part of the QK features and VO features of each head. In addition, we also record the average
811 value and standard deviation of each dimension in the original and reconstructed cache, and use these
812 to standardize, thus maintaining the numeric stability.

813 Finally, during the decoding process, theoretically, we need to continuously compress the KV cache
814 that exceeds the local length into an intermediate state. However, if the compression process is called
815 every time during decoding, it would lead to a significant waste of computation. Therefore, we set a
816 delayed compression decoding step, and only compress the exceeded local part of the KV cache once
817 after decoding for several steps, reducing the impact of compression on latency.

818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863