

CAKE: CASCADING AND ADAPTIVE KV CACHE EVICTION WITH LAYER PREFERENCES

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs)’ proficiency in handling long sequences boosts **Key-value (KV)** caching demand. Recent efforts to evict KV cache have alleviated the inference burden, but they often fail to allocate resources rationally across layers with different attention patterns. In this paper, we introduce Cascading and Adaptive KV cache Eviction (**CAKE**), a method that significantly improves LLM inference efficiency by optimizing KV cache eviction through an adaptive cache allocation strategy implemented via a cascading cache management and an innovative eviction indicator. We approach KV cache eviction as a “cake-slicing problem,” assessing each layer’s KV cache needs by considering attention dynamics in both spatial and temporal dimensions. During prompt prefilling, CAKE allocates rational cache size for layers by analyzing layer-specific KV cache preferences and manages the memory budgets with the guidance of these preferences in a cascading manner. This approach allows for a global view of cache size allocation, distributing resources **adaptively** based on the diverse attention mechanisms across layers. Also, we’ve designed a new eviction indicator that considers the shifting importance of tokens over time, addressing a limitation in existing methods that often overlook temporal dynamics. Our comprehensive experiments on the LongBench and NeedleBench datasets show that CAKE is capable of preserving the performance of models when retaining only 3.2% KV cache and consistently outperforms current baselines across various models and memory constraints, especially in low-memory situations. Moreover, CAKE outperforms full cache with FlashAttention-2 implementation, achieving 10 \times faster decoding for 128K-token sequences and maintaining consistent decoding speed across sequence lengths.

1 INTRODUCTION

Large language models (LLMs) (Zhao et al., 2023; Achiam et al., 2023; Dubey et al., 2024; Anthropic, 2024; AI, 2024) have enhanced their long-text processing capabilities, improving performance in multi-turn dialogues (Chiang et al., 2023), document summarization (Zhang et al., 2024a), question answering (Kamaloo et al., 2023), and information retrieval (Liu et al., 2024c). New models such as GPT-4 (Achiam et al., 2023), Claude 3.5 (Anthropic, 2024), LLaMA 3.1 (Dubey et al., 2024) and Mistral Large 2 (AI, 2024) have extended token processing capacities beyond 128K. Expanded contexts necessitate a linear increase in key-value (KV) cache size, resulting in heavier inference-time memory burdens. Shazeer (2019); Ainslie et al. (2023) partially address this issue by merging key-value heads during the training phase. However, optimizing key-value cache without additional training is crucial for efficient inference of long contexts under memory constraints, particularly in typical deployment scenarios where the model structure is fixed.

One way to maintain a manageable KV cache size on the fly is to remove some KV pairs (Xiao et al., 2023; Zhang et al., 2024c; Li et al., 2024b). The idea is to eliminate less important KV pairs based on certain rules. Although recent methods have enhanced pair selection for removal, they typically assign uniform cache sizes across layers, disregarding layer-specific requirements. This approach can impair performance under memory constraints. Recent research is trying to improve this by looking at how attention works in different layers (Yang et al., 2024a; Cai et al., 2024). These methods, which depend on prior observations, might not work well for all input types and models. Another idea is to adjust the cache size based on the current attention pattern (Wan et al., 2024). This can help, but it may not always optimize overall performance.

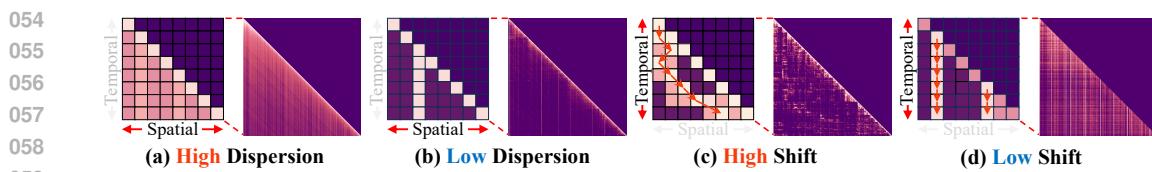


Figure 1: Variation in spatial (a, b) and temporal (c, d) characteristics of attention patterns. We provide toy examples (left) and real examples from Mistral’s different layers (right) for illustration. For more detailed analysis and visualization of attention dynamics, please refer to Appendix J.

Allocating optimal cache sizes for different layers with a fixed memory budget is akin to “cake-slicing.” The crux of this challenge is assessing each layer’s affinity for KV cache, which is mirrored in its attention mechanisms. Our research delves into these mechanisms, examining spatial and temporal aspects across layers. Spatially, attention distribution across tokens can differ markedly between layers. Some layers may spread attention broadly, while others concentrate on specific tokens, as depicted in Figure 1(a) and Figure 1(b). Temporally, attention hotspots in certain layers shift over time in Figure 1(c), whereas in others, they remain constant in Figure 1(d). The variability of attention pattern necessitates a dynamic approach to memory allocation that goes beyond uniform (Zhang et al., 2024c; Liu et al., 2024d; Ren & Zhu, 2024; Li et al., 2024b), depicted in Figure 2(a), or fixed-pattern methods (Yang et al., 2024a), drawn in Figure 2(b). They lack the flexibility to adapt to the nuanced attention dynamics. To enhance cache efficiency, a comprehensive strategy is required, factoring in each layer’s attention pattern and KV cache preference. Moreover, conventional methods that evict KV pairs relying solely on static spatial attention scores often overlook the temporal evolution of attention, missing critical insights into how token relevance changes in various contexts.

In this paper, we introduce Cascading and Adaptive KV cache Eviction (CAKE), enhancing KV cache eviction by leveraging proposed preference-prioritized adaptive cache allocation strategy coupled with an innovative token eviction strategy. It is progressive and dynamic, employing a novel metric that assesses each layer’s cache size preference, taking into account both the dispersion of spatial attention and the shifts in temporal attention, as outlined in Figure 2(c). To optimize cache allocation from a holistic viewpoint while keeping peak memory usage within desired limits, we have designed a cascading memory management system. This system dynamically adjusts the KV cache during the prompt prefilling phase, eliminating the need to store all KV cache for all layers. For token eviction, acknowledging the limitations of existing methods, we introduce an eviction indicator that considers sustained importance and attention variability, thereby minimizing the adverse effects of eviction on subsequent decoding steps. Extensive experiments have been conducted using various LLM architectures on the LongBench and NeedleBench benchmarks, encompassing nine major task categories. Results demonstrate CAKE’s superior performance across diverse memory scenarios within these benchmarks. The allocation strategy in CAKE, which can complement and enhance existing KV cache eviction methods, offers a versatile solution for improved cache management. Furthermore, compared to full cache FlashAttention-2 implementation, CAKE significantly reduces memory consumption while simultaneously enhancing LLM throughput.

Our contributions are: (1) An analysis of attention dynamics revealing spatial dispersion and temporal shifts, leading to a layer-specific metric for cache size requirements. (2) An adaptive cache allocation strategy that optimizes overall cache allocation based on layer preferences. (3) A cascading cache management method that dynamically adjusts the KV cache during the prefilling stage, achieving memory usage efficiency comparable to uniform strategies without sacrificing eviction performance. (4) A new eviction indicator that considers the sustained importance and variability of tokens, enhancing token eviction performance.

2 BACKGROUND AND RELATED WORKS

Basics of KV Cache Operations. We revisit the fundamentals of KV caching. We focus on a single attention head, characterized as weight matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D \times D}$, where D denotes the model’s hidden dimension. Considering a prompt embedding $\mathbf{X} \in \mathbb{R}^{S \times D}$, with S representing the sequence length, an attention module has two phases: prompt prefilling and token decoding.

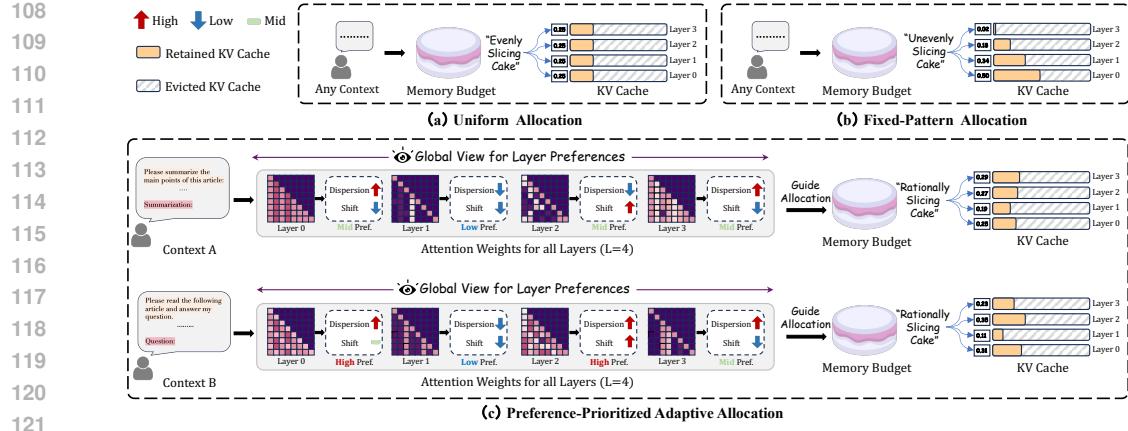


Figure 2: Illustration of CAKE compared with existing cache allocation strategies. (a) Uniform cache allocation (Xiao et al., 2023; Zhang et al., 2024c; Li et al., 2024b); (b) Fixed-shape cache allocation (Cai et al., 2024; Yang et al., 2024a); (c) Preference-prioritized adaptive cache allocation used in CAKE. Compared to (a) and (b), CAKE adjusts allocation ratios across different layers with layer preferences, adapting to various contexts and models, with given memory budgets.

Prompt Prefilling: The query, key, and value states are initially calculated as follows:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V. \quad (1)$$

Then, the output of the attention module is determined as:

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A}\mathbf{V}, \quad (2)$$

with attention weights $\mathbf{A} = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D}}) \in \mathbb{R}^{S \times S}$. The key-value states \mathbf{K} and \mathbf{V} are then stored in a cache, establishing the KV cache.

Token Decoding: The KV cache is used and updated to produce new tokens. For each decoding step i , let the new token embedding be $\mathbf{x}_i \in \mathbb{R}^{1 \times D}$. To circumvent repeated key-value projections, the KV cache is concatenated with the newly generated KV pair $\mathbf{x}_i\mathbf{W}_K, \mathbf{x}_i\mathbf{W}_V \in \mathbb{R}^{1 \times D}$ for the current attention computation and to refresh the KV cache:

$$\mathbf{K} = \text{Concat}(\mathbf{K}, \mathbf{x}_i\mathbf{W}_K), \quad \mathbf{V} = \text{Concat}(\mathbf{V}, \mathbf{x}_i\mathbf{W}_V). \quad (3)$$

Though the KV cache alleviates the considerable computational demands of attention mechanisms, its linear growth with sequence length poses a challenge for extremely long input or output sequences. Efficiently managing the KV cache within a finite cache budget remains a formidable issue, especially in contexts that demand longer context lengths.

KV Cache Eviction. KV cache eviction during model inference enhances computational efficiency without altering the attention mechanism. This optimization relies on strategic cache reduction techniques. Early approaches to KV cache eviction focus on specific parts of the input sequence. For instance, StreamingLLM (Xiao et al., 2023) and LM-Infinite (Han et al., 2023) prioritize the retention of the first and last tokens. However, this strategy risks ignoring potentially important tokens in the middle of the sequence. Recognizing this limitation, subsequent research introduces more sophisticated indicators to filter unnecessary KV cache entries, such as using cumulative attention scores (Zhang et al., 2024c; Liu et al., 2024d), last token attention scores (Oren et al., 2024), mean attention scores (Ren & Zhu, 2024), or clustering recent attention scores (Li et al., 2024b). While these methods offer more nuanced approaches to cache eviction, they often apply the uniform allocation strategy, which may not maximize cache utilization. Recent FastGen (Ge et al., 2023) chooses which tokens to keep based on special markers, but it doesn't limit the total cache size. PyramidInfer (Yang et al., 2024a) and PyramidKV (Cai et al., 2024) allocate the budget in a pyramid-shaped manner, while D2O (Wan et al., 2024) adjusts cache size based on the current layer's attention density. While these methods have shown promise, they often rely on predefined cache allocation strategies across all layers or are just based on the local layer attention, which does not fully capture the complex dynamics of attention mechanisms and lacks a global perspective on layer preferences for the cache. In contrast, our work addresses these limitations by considering layer-specific cache preferences with a comprehensive and global perspective.

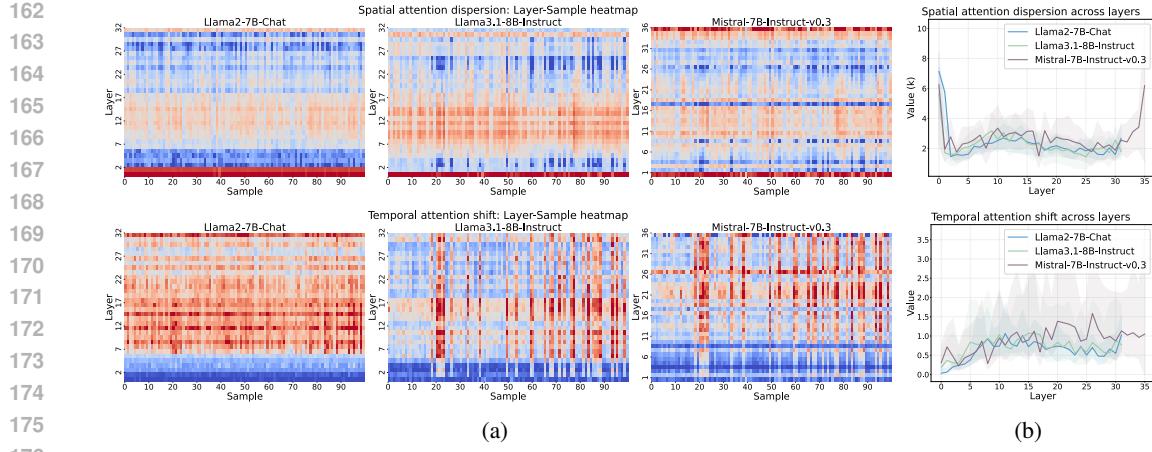


Figure 3: Analysis of attention dynamics. (a) Heatmap for spatial attention dispersion (upper) and temporal attention shift (lower), red color with high value, while blue for low value. The x-axis represents samples, and the y-axis represents layers. (b) Variation of spatial attention dispersion (upper) and temporal attention shift (lower) across layers and models. The experimental data is derived from the LongBench dataset (Bai et al., 2023).

3 INSIGHTS INTO ATTENTION DYNAMICS

Recent studies (Xiao et al., 2023; Zhang et al., 2024c) have shown that only a few KV cache elements are crucial during token decoding. Building on this, our analysis identifies two pivotal attention traits: *spatial attention dispersion*, showing how one token’s attention spreads across others, and *temporal attention shift*, tracking the evolution of highly-attended tokens over time.

Recall some basics about the attention weights \mathbf{A} . Each $\mathbf{A}[i, j]$ shows how much the i -th query token attends to the j -th key token. We analyze the matrix spatially and temporally to explore attention dynamics in decoder-only models. (1) **Spatial Analysis:** Row i , $\mathbf{A}[i, :]$, shows how the i -th token’s attention is distributed across tokens in one step. Examining $\mathbf{A}[i, :]$ reveals the attention landscape at that step. This helps reveal how the model prioritizes information at each generation step. (2) **Temporal Analysis:** Column j , $\mathbf{A}[:, j]$, shows how attention to the j -th token changes over steps. Each row in $\mathbf{A}[:, j]$ represents a time step, showing the j -th token’s evolving attention. This dimension is key to tracking the model’s shifting focus during sequence generation.

In this paper, we quantify the spatial attention dispersion by the entropy of $\mathbf{A}[i, :]$, defined as:

$$H(\mathbf{A}) = - \sum_{i=0}^{S-1} \mathbf{A}[i, :] \log(\mathbf{A}[i, :])^T, \quad (4)$$

where the inner product is conducted between each row of \mathbf{A} and its element-wise logarithm. Rows show higher values with even attention distribution and lower values with focused attention. Aggregating these values measures the evenness of attention distribution per token.

Accordingly, the temporal attention shift is characterized by the variance of $\mathbf{A}[:, j]$, expressed as:

$$V(\mathbf{A}) = \sum_{j=0}^{S-1} \text{Var}(\mathbf{A}[:, j]), \quad (5)$$

which computes the variance for each column of \mathbf{A} . Higher values indicate significant attention shifts across positions, while lower values suggest stable attention. Summing the variances column-wise captures the attention mechanism’s temporal dynamics.

We analyze attention dynamics across multiple LLM layers using the LongBench dataset’s long-context samples (Bai et al., 2023). Figure 3 (upper) shows attention dispersion changes across layers, and Figure 3 (lower) shows attention focus shifts. The results show significant variations in attention dispersion and shift across layers, models, and contexts. Despite recent KV cache advances (Xiao et al., 2023; Zhang et al., 2024c), our findings highlight LLM attention complexity, indicating a need for tailored KV cache management strategies to effectively handle the dynamic nature of attention.

216 **4 METHODOLOGY**

217 **4.1 PREFERENCE-PRIORITIZED ADAPTIVE ALLOCATION**

220 Given the variability in attention mechanisms across layers, models, and contexts, uniform or fixed-
 221 pattern cache allocation strategies prove inefficient, especially under limited memory budgets. To
 222 address these challenges, we introduce the preference-prioritized adaptive allocation strategy. It
 223 considers each layer’s unique characteristics and adaptively allocates cache sizes from a global view
 224 of attention patterns. We define a preference metric for each layer’s KV cache requirements, con-
 225 sidering both the spatial dispersion and temporal shift of attention:

226
$$\mathcal{P} = \mathcal{H}^{\frac{1}{\tau_1}} \cdot \mathcal{V}^{\frac{1}{\tau_2}}, \quad \mathcal{H} = H(\mathbf{A}[-S_w :, : -S_w]), \quad \mathcal{V} = V(\mathbf{A}[-S_w :, : -S_w]), \quad (6)$$

228 where \mathcal{P} is the preference score for an attention layer’s cache size. As dispersed attention (high \mathcal{H})
 229 requires retaining broader contexts, while frequent shifts (high \mathcal{V}) demand support for complex tem-
 230 poral patterns, layers with high preference score \mathcal{P} benefit more from larger KV cache to maintain
 231 performance. Accounting for the varying importance of attention dispersion and shift across models
 232 and memory constraints, we introduce temperature parameters τ_1 and τ_2 to adjust their influence
 233 on \mathcal{P} . We focus on the submatrix $\mathbf{A}[-S_w :, : -S_w]$ of \mathbf{A} , representing a recent window of size
 234 S_w . This focus on the recent window chunk of prefilling attention weights is inspired by recent re-
 235 search (Li et al., 2024b; Yang et al., 2024a). These studies have demonstrated that decoding patterns
 236 can be effectively captured by analyzing the attention patterns of recent queries.

237 We implement the preference-prioritized adaptive allocation strategy using layer-specific preference
 238 scores. These scores adjust dynamically to varying models and contexts, ensuring adaptability. Each
 239 layer’s cache size is set by normalizing the preference scores and allocating the total memory budget
 240 accordingly. This results in an adaptive cache size set $\mathbf{B} = \{B_l\}_{l=0}^{L-1}$, calculated as follows:

241
$$B_l = \frac{\mathcal{P}_l}{\sum_{k=0}^{L-1} \mathcal{P}_k} \cdot B_{\text{total}}, \quad (7)$$

244 where \mathcal{P}_l is the preference score for layer l , and B_{total} is the total cache budget. This method
 245 optimizes cache size allocation for each layer’s unique characteristics and the input context.

246 **4.2 PREFERENCE-GUIDED CASCADING CACHE MANAGEMENT**

249 Our preference-prioritized adaptive allocation strategy, though optimal, initially requires a com-
 250 prehensive view of prefilling attention weights across all layers. This leads to the peak memory usage
 251 of $\mathcal{O}(S \cdot L)$, which may exceed the cache budget B_{total} .

252 To overcome this, we further develop the preference-guided cascading cache management, as illus-
 253 trated in Figure 4. It dynamically maintains the cache budget during prefilling, effectively reducing
 254 peak memory usage to the target. Algorithm 1 outlines our cache management procedure, which par-
 255 titions the prefilling process into L stages, one for each model layer, and cascadingly manages cache
 256 memory guided by preference scores. At each stage, as a new layer completes its prefilling compu-
 257 tation, the budget is redistributed based on the current obtained preference scores (Algorithm 1, lines
 258 3–7), and KV caches are updated among all processed layers according to the redistributed budget
 259 (Algorithm 1, lines 9–14, which, while presented sequentially, can be parallelized as KV caches are
 260 pre-computed, reducing the time complexity to that of a single iteration). **This approach ensures**
 261 **constant memory usage during the prefilling stage while maintaining equivalent cache distribution**
 262 **and eviction results as the standard strategy. Our algorithm is supported by the following theoretical**
 263 **results (Proposition 1 and Theorem 1, proofs provided in Appendix D).**

264 **Proposition 1.** *For any layer $l \in [L]$, the allocated budget size decreases monotonically from stage
 265 l to $L - 1$:*

266
$$B_l^{(m+1)} < B_l^{(m)}, m \in [l, L - 1], \quad (8)$$

267 where $B_l^{(m)}$ is the redistributed cache budget for layer l at stage m , calculated based on the current
 268 obtained preference score \mathcal{P} :

269
$$B_l^{(m)} = \frac{\mathcal{P}_l}{\sum_{k=0}^m \mathcal{P}_k} \cdot B_{\text{total}}. \quad (9)$$

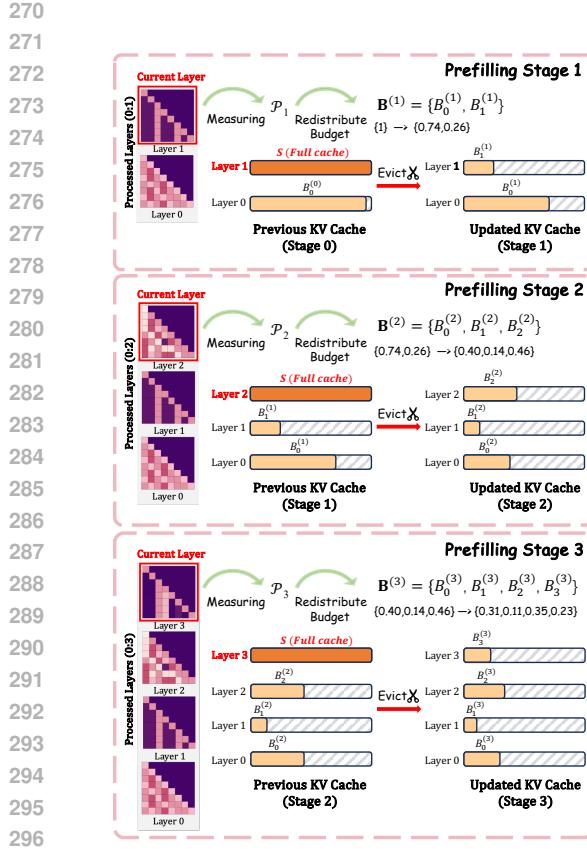


Figure 4: Illustration of preference-guided cascading cache management algorithm during three stages.

Algorithm 1 Preference-guided Cascading Cache Management

Input: Total cache budget B_{total} , Attention Weights $\mathbf{A}_m \in \mathbb{R}^{S \times S}$, Matrices $\mathbf{K}_m, \mathbf{V}_m \in \mathbb{R}^{S \times D}$, Number of layers L

Output: Retained cache set $\mathbf{C} = \{\{\hat{\mathbf{K}}_m^{(L-1)}, \hat{\mathbf{V}}_m^{(L-1)}\} : m \in [L]\}$

- 1: Initialize $\mathbf{C}^{(0)}, \mathbf{S}^{(0)}, \mathbf{B}^{(0)}, \mathbf{P}^{(0)} \leftarrow \emptyset$
- 2: **for** stage $m = 0$ to $L - 1$ **do**
- 3: $\mathbf{C}^{(m)} \leftarrow \mathbf{C}^{(m-1)} \cup (\{\mathbf{K}_m, \mathbf{V}_m\})$
- 4: Compute preference score \mathcal{P}_m for layer m using Eq.(6)
- 5: Compute indicator \mathbf{I}_m ($\mathbf{I}_m^{(m)}$) using \mathbf{A}_m from a chosen eviction method
- 6: $\mathbf{P}^{(m)} \leftarrow \mathbf{P}^{(m-1)} \cup (\mathcal{P}_m)$,
- 7: $\mathbf{S}^{(m)} \leftarrow \mathbf{S}^{(m-1)} \cup (\mathbf{I}_m^{(m)})$
- 8: $\mathbf{B}^{(m)} \leftarrow \{B_i^{(m)} : i \in [m]\}$, where $B_i^{(m)}$ is computed using Eq.(9)
- 9: // Manage current KV Caches
- 10: **for** layer $l = 0$ to m **do**
- 11: $\{\mathbf{K}_l, \mathbf{V}_l\}^{(m)} \leftarrow \mathbf{C}^{(m)}[l]$,
 $B_l^{(m)} \leftarrow \mathbf{B}^{(m)}[l]$, $\mathbf{I}_l^{(m)} \leftarrow \mathbf{S}^{(m)}[l]$
// Update cache for layer l
- 12: $\{\hat{\mathbf{K}}_l, \hat{\mathbf{V}}_l\}^{(m)}, \hat{\mathbf{I}}_l^{(m)} \leftarrow \text{EVICT}(\{\mathbf{K}_l, \mathbf{V}_l\}^{(m)}, B_l^{(m)}, \mathbf{I}_l^{(m)})$
- 13: $\mathbf{C}^{(m)}[l] \leftarrow \{\mathbf{K}_l, \mathbf{V}_l\}^{(m+1)} = \{\hat{\mathbf{K}}_l, \hat{\mathbf{V}}_l\}^{(m)}$,
 $\mathbf{S}^{(m)}[l] \leftarrow \mathbf{I}_l^{(m+1)} = \hat{\mathbf{I}}_l^{(m)}$
- 14: **end for**
- 15: **end for**
- 16: **return** Retained cache set $\mathbf{C} = \{\{\hat{\mathbf{K}}_m^{(L-1)}, \hat{\mathbf{V}}_m^{(L-1)}\} : m \in [L]\}$

Proposition 1 guarantees that the cache budget allocated to each layer decreases monotonically as stages process, ensuring that the KV cache $\{\mathbf{K}_l, \mathbf{V}_l\}^{(m+1)}$ is always a proper subset of the previous stage's cache $\{\mathbf{K}_l, \mathbf{V}_l\}^{(m)}$, regardless of the eviction method used.

Before presenting Theorem 1, let's define some operational details related to our algorithm. Given indicator vector $\mathbf{I}_l \in \mathbb{R}^S$, it measures token importance for layer l . The specific calculation of our eviction indicator will be detailed in the subsequent section. For each stage $m \in [L]$, the eviction operation, denoted as $\text{EVICT}(\cdot)$, retains KV pairs $\hat{\mathbf{K}}_l^{(m)}, \hat{\mathbf{V}}_l^{(m)} \in \mathbb{R}^{B_l^{(m)} \times D}$ from the current cache $\mathbf{K}_l^{(m)}, \mathbf{V}_l^{(m)} \in \mathbb{R}^{B_l^{(m-1)} \times D}$ corresponding to the top- $B_l^{(m)}$ positions with the highest scores in $\mathbf{I}_l^{(m)} \in \mathbb{R}^{B_l^{(m-1)}}$. The indicator vector $\mathbf{I}_l^{(m)}$ is also updated to $\hat{\mathbf{I}}_l^{(m)} \in \mathbb{R}^{B_l^{(m)}}$, retaining only the elements corresponding to the preserved positions. The updated cache and indicator $\{\hat{\mathbf{K}}_l, \hat{\mathbf{V}}_l\}^{(m)}$, $\hat{\mathbf{I}}_l^{(m)}$ will serve as the basis $\{\mathbf{K}_l, \mathbf{V}_l\}^{(m+1)}$ and $\mathbf{I}_l^{(m+1)}$ for further updates in the next stage $m + 1$.

Theorem 1. For any layer $l \in [L]$, the KV cache $\{\mathbf{K}_l, \mathbf{V}_l\}^{(L-1)}$ obtained through cascading eviction from stage l to stage $L - 1$ is equivalent to the result of applying the eviction operation once on the full KV cache $\{\mathbf{K}_l, \mathbf{V}_l\} \in \mathbb{R}^{S \times D}$ with the cache budget B_l computed in Eq.(7):

$$\{\mathbf{K}_l, \mathbf{V}_l\}^{(L-1)} = \text{EVICT}(\{\mathbf{K}_l, \mathbf{V}_l\}^{(m)}, B_l^{(m)}, \mathbf{I}_l^{(m)})_{m=l}^{L-1} = \text{EVICT}(\{\mathbf{K}_l, \mathbf{V}_l\}, B_l, \mathbf{I}_l). \quad (10)$$

Theorem 1 demonstrates that our preference-guided cascading cache management achieves equivalent KV cache eviction results to the vanilla preference-prioritized adaptive allocation strategy, despite its cascading operation. This method is theoretically compatible with existing KV eviction techniques, a claim that will be empirically validated in Section 5.5.

324 4.3 ATTENTION-SHIFT TOLERANT EVICTION INDICATOR
325

326 Current top eviction indicators, such as accumulated or mean attention scores (Zhang et al., 2024c;
327 Liu et al., 2024d; Ren & Zhu, 2024), simplify tokens’ attention profiles into single values, effec-
328 tively identifying important tokens. However, this approach may overlook tokens whose importance
329 fluctuates, as it fails to capture the dynamics of shifting attention. This could result in prematurely
330 evicting tokens from the cache, affecting the model’s future ability to retrieve relevant information.

331 We propose a robust eviction strategy tolerant of attention shifts. We use a multi-faceted indicator
332 considering sustained importance and attention variability. For layer l , the eviction indicator $\mathbf{I}_l \in \mathbb{R}^S$
333 is computed, where each element $\mathbf{I}_l[n]$ signifies the importance of token n , computed as:

$$334 \quad \mathbf{I}_l[n] = \begin{cases} \text{Mean}(\mathbf{A}_l[-S_w : , n]) + \gamma \cdot \text{Var}(\mathbf{A}_l[-S_w : , n]), & \text{if } n < S - S_w, \\ \Omega, & \text{otherwise,} \end{cases} \quad (11)$$

335 where $\text{Mean}(\cdot)$ and $\text{Var}(\cdot)$ measure sustained importance and attention variability, respectively, with
336 γ adjusting their influence. Inspired by recent research (Li et al., 2024b; Yang et al., 2024a), we
337 assign an arbitrarily large value Ω to ensure the preservation of the most recent W tokens. Also, a
338 pooling layer clusters $\mathbf{I}_l[: S - S_w]$ to maintain context and prevent information fragmentation.

339 The eviction operation $\text{EVICT}(\{\mathbf{K}_l, \mathbf{V}_l\}, B_l, \mathbf{I}_l)$ is then executed, retaining the KV pairs $\{\hat{\mathbf{K}}_l, \hat{\mathbf{V}}_l\}$
340 that correspond to the top- B_l scores in \mathbf{I}_l :

$$341 \quad \hat{\mathbf{K}}_l = \mathbf{K}_l[\mathbf{D}_l, :], \quad \hat{\mathbf{V}}_l = \mathbf{V}_l[\mathbf{D}_l, :], \quad \mathbf{D}_l = \text{TopK}(\mathbf{I}_l, B_l), \quad (12)$$

342 where TopK selects the indices $\mathbf{D}_l \in \mathbb{R}^{B_l}$ of the top B_l scores in \mathbf{I}_l .
343

344 5 EXPERIMENTATION

345 5.1 EXPERIMENTAL SETTINGS

346 **Backbone LLMs.** Our experiments involve five major open-source LLMs (7B-70B parameters)
347 capable of handling 4k-128k token contexts. These models feature two representative attention
348 structures: (1) *Multi-head attention*: including Llama2-Chat (Touvron et al., 2023) and Gemma-
349 Instruct (Team et al., 2024). (2) *Grouped-query attention*: including Llama3-Instruct (Dubey et al.,
350 2024), Mistral-v0.3 (Jiang et al., 2023), and Qwen2.5-Instruct (Team, 2024).

351 **Baseline Methods.** We assess CAKE’s performance against five baselines: (1) *Uniform Allocation*:
352 StreamingLLM (Xiao et al., 2023) balances initial and recent tokens, H2O (Zhang et al., 2024c)
353 uses cumulative attention, TOVA (Oren et al., 2024) adopts last-token attention, and SnapKV (Li
354 et al., 2024b) clusters recent attention. (2) *Non-Uniform Allocation*: PyramidKV (Cai et al., 2024)
355 employs a fixed pyramid-shaped allocation strategy with SnapKV’s eviction indicator.
356

357 **Evaluating Tasks.** To evaluate CAKE’s performance across various memory budgets, we use two
358 carefully designed benchmarks: (1) *LongBench* (Bai et al., 2023): Focuses on long-context under-
359 standing, encompassing 16 datasets in six categories: Single/Multi-Document QA, Summarization,
360 Few-shot Learning, Synthetic Tasks, and Code Completion. (2) *NeedleBench* (Li et al., 2024a):
361 Tests retrieval and reasoning in complex contexts through three subtasks: Single-Needle Retrieval,
362 Multi-Needle Retrieval, and Multi-Needle Reasoning.
363

364 **Implementation Details.** CAKE is compared to baseline methods across different total memory
365 budgets ranging from $64L$ to $2048L$. For uniform allocation, each layer has a fixed KV cache
366 size. Non-uniform methods like PyramidKV and CAKE vary cache sizes per layer while keeping
367 total memory constant. Tokens are evicted during both prefilling and decoding to maintain memory
368 usage. Experiments are run on NVIDIA A100 80GB GPUs. For specifics, see the Appendix C.
369

370 5.2 EVALUATIONS ON LONGBENCH DATASET

371 We benchmark CAKE against other methods on 16 datasets. Figure 5 displays the performance
372 of various methods under different memory constraints across three models. SnapKV surpasses
373 legacy methods like StreamingLLM and H2O by using an observation window for memory reten-
374 tion. TOVA, relying solely on the last token, falls short by ignoring the broader context. Despite
375

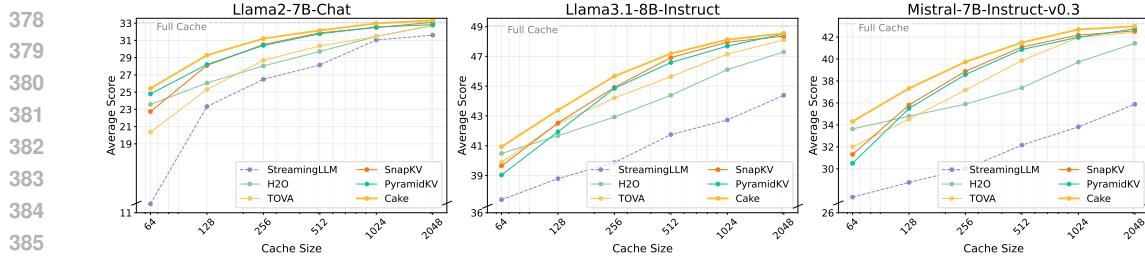


Figure 5: Average score among 16 datasets of LongBench under different cache budgets.

Table 1: Performance comparison over 16 datasets of LongBench. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code	
	NtvrQA	Qasper	MF-en	HopqaQA	2WikiQA	Musique	GovReport	QMSum	MultiNews	TREC	TriviaQA	SAMSUM	PCount	PR-en	Lcc	RB-p	Avg.			
Llama2-7B-Chat, $B_{\text{total}} = 128L$																				
StreamingLLM	13.57	13.33	22.39	20.16	17.12	6.38	16.72	19.58	15.37	31.50	74.84	33.00	1.67	4.00	43.21	40.45	23.33			
H2O	<u>14.56</u>	14.59	26.52	20.05	<u>28.73</u>	5.00	16.20	20.30	20.34	38.00	72.75	37.10	2.33	3.50	49.68	47.37	26.06			
TOVA	13.56	15.46	20.28	22.22	25.01	4.79	15.72	19.11	15.58	42.50	79.58	37.06	<u>3.45</u>	4.61	46.34	39.70	25.31			
SnapKV	13.70	16.27	<u>32.52</u>	22.76	28.59	<u>7.56</u>	18.87	19.96	20.05	43.50	79.90	36.74	2.79	7.00	51.86	47.41	28.09			
PyramidKV	13.45	16.61	31.21	<u>23.74</u>	28.43	6.89	<u>18.88</u>	<u>20.56</u>	19.81	44.50	80.32	37.82	2.29	<u>8.00</u>	51.19	47.56	28.20			
CAKE	15.15	16.38	<u>32.17</u>	<u>25.25</u>	31.09	7.00	19.47	21.07	20.36	48.50	80.66	<u>37.77</u>	4.42	9.00	<u>51.48</u>	48.90	29.29			
Llama2-7B-Chat, $B_{\text{total}} = 1024L$																				
StreamingLLM	13.73	16.75	26.28	25.97	25.67	5.87	22.21	19.27	24.07	61.00	80.88	40.90	2.28	1.00	56.71	48.79	29.46			
H2O	16.58	17.67	32.04	26.01	28.73	7.81	22.87	20.99	25.05	60.00	83.02	40.06	2.90	3.50	57.57	52.02	31.05			
TOVA	15.00	17.32	32.97	27.48	<u>31.71</u>	7.04	21.86	20.54	24.78	63.50	83.35	41.78	2.63	<u>8.00</u>	56.69	52.11	31.46			
SnapKV	<u>18.05</u>	19.73	38.07	<u>27.57</u>	30.86	<u>8.41</u>	<u>23.41</u>	20.79	25.22	63.50	82.13	40.98	<u>3.33</u>	7.50	58.20	<u>52.23</u>	32.50			
PyramidKV	17.94	<u>20.68</u>	38.85	27.25	29.99	7.91	23.17	21.55	25.45	63.50	82.97	40.79	<u>3.57</u>	7.50	57.72	51.93	<u>32.55</u>			
CAKE	18.35	20.93	<u>38.49</u>	<u>27.99</u>	<u>31.15</u>	8.59	24.29	<u>21.06</u>	<u>25.40</u>	<u>63.50</u>	<u>83.33</u>	<u>41.30</u>	3.11	9.50	<u>57.82</u>	52.93	32.98			
Llama3.1-8B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	23.01	22.50	31.00	46.79	40.52	24.76	18.38	20.89	16.96	40.50	86.00	38.55	7.00	99.50	57.06	47.16	38.79			
H2O	26.18	24.83	44.08	51.33	43.29	27.78	20.79	22.36	<u>20.69</u>	41.50	89.88	<u>40.94</u>	6.20	99.00	58.12	49.82	41.67			
TOVA	29.44	27.03	45.20	<u>54.15</u>	<u>44.87</u>	28.87	22.21	20.54	20.52	53.50	92.27	40.99	6.17	99.50	51.14	43.15	42.47			
SnapKV	26.40	27.67	<u>47.83</u>	<u>53.40</u>	44.20	<u>28.68</u>	20.56	<u>23.11</u>	20.13	45.50	89.36	39.98	6.33	99.00	57.69	50.69	42.53			
PyramidKV	25.57	27.40	46.36	53.14	44.51	27.16	20.74	22.26	19.90	45.50	87.06	40.13	<u>6.50</u>	99.50	56.98	48.26	41.94			
CAKE	27.41	32.01	49.58	52.99	45.76	28.32	<u>21.85</u>	23.15	21.56	47.50	89.61	40.90	6.33	99.50	<u>57.92</u>	49.83	43.39			
Llama3.1-8B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	26.64	30.77	35.59	47.31	42.03	24.17	25.81	21.31	25.66	63.50	88.84	42.76	6.50	88.00	61.36	53.47	42.73			
H2O	29.57	36.15	45.94	54.43	44.81	29.04	27.64	23.31	45.47	62.00	91.83	<u>43.14</u>	6.36	99.00	<u>62.74</u>	55.39	46.11			
TOVA	30.66	40.95	51.09	54.58	46.51	30.62	<u>28.12</u>	23.61	26.24	68.00	91.49	43.80	5.92	99.50	60.73	52.64	47.15			
SnapKV	30.95	44.74	<u>52.58</u>	55.09	46.83	30.37	27.87	<u>24.57</u>	25.99	68.00	92.03	42.60	6.50	99.50	63.00	<u>56.50</u>	47.95			
PyramidKV	30.54	43.64	<u>52.73</u>	55.29	46.29	31.28	27.53	24.50	26.00	68.00	<u>92.09</u>	41.75	6.05	99.50	62.35	55.44	47.69			
CAKE	<u>30.88</u>	44.95	52.38	<u>55.49</u>	46.99	30.82	28.68	24.91	<u>26.39</u>	69.00	91.94	42.60	6.00	99.50	62.65	56.89	48.13			
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 128L$																				
StreamingLLM	16.91	21.51	24.85	34.14	26.99	16.64	15.67	18.61	14.40	43.50	83.26	37.00	0.00	23.00	41.63	42.12	28.76			
H2O	21.25	26.66	35.13	38.82	<u>29.80</u>	18.88	<u>21.00</u>	19.50	<u>18.63</u>	41.00	87.64	<u>38.25</u>	1.50	67.00	44.35	47.29	34.79			
TOVA	22.47	24.26	37.22	42.26	28.85	19.97	19.40	18.70	17.86	63.00	88.98	37.71	0.50	56.50	37.42	37.22	34.52			
SnapKV	21.02	27.26	41.25	<u>45.15</u>	29.23	22.75	20.47	<u>20.17</u>	17.75	42.50	87.28	38.01	0.50	69.50	44.48	45.69	<u>35.81</u>			
PyramidKV	21.73	26.60	<u>41.46</u>	43.20	29.32	21.47	20.23	19.82	17.46	40.00	87.64	37.11	<u>1.05</u>	69.00	42.55	42.98	35.14			
CAKE	<u>22.31</u>	29.15	<u>43.51</u>	<u>44.51</u>	30.36	22.85	21.56	20.47	18.96	47.00	88.60	39.36	1.00	76.50	44.96	<u>46.19</u>	37.33			
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 1024L$																				
StreamingLLM	20.96	28.05	30.03	37.06	27.56	16.03	24.03	19.07	22.79	67.00	87.61	40.96	1.50	21.50	48.05	48.87	33.82			
H2O	23.78	31.63	41.31	43.24	31.07	20.43	26.74	20.41	<u>23.93</u>	67.50	88.84	42.62	1.50	72.00	50.60	49.87	39.72			
TOVA	26.97	34.51	45.58	44.32	<u>32.58</u>	22.83	<u>26.91</u>	20.75	23.49	75.00	88.66	43.17	1.00	<u>91.00</u>	47.51	47.06	41.96			
SnapKV	<u>26.63</u>	35.78	48.11	<u>45.75</u>	32.20	23.37	26.71	21.84	23.18	70.50	88.61	41.37	0.50	88.00	<u>50.60</u>	51.79	42.18			
PyramidKV	25.51	<u>36.02</u>	47.72	44.74	33.16	23.91	26.55	<u>21.83</u>	23.27	70.50	88.41	40.94	1.00	87.00	50.17	50.93	41.98			
CAKE	26.09	36.34	48.11	<u>45.97</u>	32.39	<u>23.49</u>	27.56	21.45	24.03	<u>72.50</u>	88.61	<u>42.71</u>	0.00	91.50	51.06	<u>51.25</u>	42.69			

similarities to SnapKV, PyramidKV underperforms in many scenarios, possibly due to its limited adaptability. CAKE, however, outperforms others by dynamically allocating memory to each layer based on its specific needs, especially under tight memory conditions. Notably, CAKE achieves performance comparable to full cache models while only utilizing a small cache budget ($B_{\text{total}} > 512L$).

For detailed analysis, see Table 1 for results across two memory scenarios: a low setting ($B_{\text{total}} = 128L$) and a high setting ($B_{\text{total}} = 1024L$). Full results under different constraints are detailed in Appendix F.1. CAKE consistently ranks among the top performers or closely matches them across various tasks. In the low memory scenario with the Llama2 model, CAKE leads or finishes second in all datasets. CAKE's main advantage is its scenario-tailored memory distribution strategy, ensuring balanced performance. We also provide additional evaluations on diverse architectures like Qwen2.5 and Gemma as well as larger models ranging from 13B to 70B parameters in Appendix F.2 and F.3.

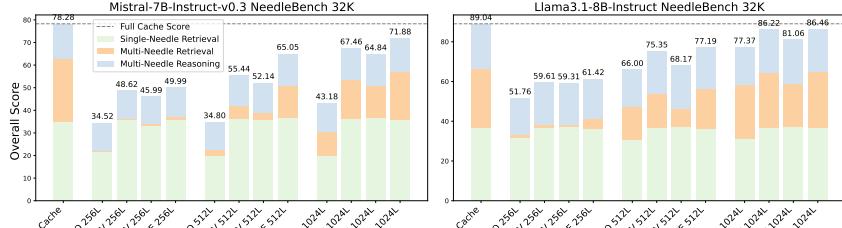


Figure 6: Performance comparison on NeedleBench 32K.

5.3 EVALUATIONS ON NEEDLEBENCH DATASET

We also test CAKE on retrieval and reasoning tasks using the Needlebench dataset, with the results depicted in Figure 6. Consistent with the LongBench outcomes, CAKE shows superior overall performance. Notably, CAKE excels in more complex tasks akin to real-world scenarios, such as Multi-Retrieval and Multi-Reasoning. CAKE efficiently utilizes only 3.2% of the cache size while preserving the model’s capabilities for processing 32K-token-long contexts. This performance indicates CAKE’s proficiency in managing and leveraging information across diverse contexts, despite their length or complexity. In the simpler Single-retrieval task, both CAKE and SnapKV ($B_{\text{total}} = 1024L$) surpass the performance achieved with the full cache. This suggests that pruning specific cache information may eliminate superfluous details, allowing generation to focus on the most salient aspects and enhance its overall efficacy. For a more in-depth examination of the experimental outcomes, including detailed statistics and additional graphs, please refer to Appendix G.

5.4 EVALUATION OF MEMORY AND THROUGHPUT

We assess the effectiveness of our method in reducing memory consumption and enhancing time efficiency during LLM inference by analyzing peak memory usage and decoding latency on Mistral-7B-Instruct-v0.3 [implemented with FlashAttention-2](#) (Dao, 2023). Our comparison includes full cache, and three KV cache eviction methods, all with $B_{\text{total}} = 1024L$: SnapKV (uniform allocation strategy), PyramidKV (fixed-pattern allocation strategy), and our proposed CAKE method.

Peak Memory Usage. As depicted in the left panel of Figure 7, CAKE exhibits substantial memory-saving capabilities, comparable to other KV cache eviction methods such as uniform allocation (SnapKV) and fixed-pattern non-uniform allocation (PyramidKV). All these approaches maintain a fixed-size KV cache. When compared to the full cache implementation, CAKE achieves an impressive reduction in peak memory usage of approximate 48.63% with a 128K context length.

Throughput Analysis. Regarding decoding latency, the right panel of Figure 7 shows that both uniform and non-uniform eviction methods exhibit comparable inference performance. As input length increases, decoding latency for the full cache method grows significantly due to escalating computational demands and I/O latency bottlenecks. In contrast, CAKE maintains a relatively stable decoding speed by preserving a fixed amount of KV cache, resulting in significantly lower latency compared to the full cache, particularly for longer sequences. It is noteworthy that CAKE demonstrates remarkable efficiency, achieving over 10× speedup in decoding latency compared to the full cache approach when processing sequences with 128K context length.

5.5 COMPATIBILITY WITH EXISTING KV EVICTION METHODS

Our preference-prioritized adaptive allocation strategy demonstrates strong compatibility with existing eviction indicators. We evaluate its performance using two representative methods: H2O, which serves as a classical indicator, and SnapKV, which represents a new advanced indicator. The experiments are conducted on LongBench datasets using Llama2-7B-Chat under B_{total} of 128L and 512L. We report the average performance across six tasks. In Figure 8, when compared with vanilla uniform cache allocation, methods equipped with our allocation strategy consistently improve performance across nearly all tasks. Importantly, they achieve significant overall performance gains. This comprehensive improvement across different eviction methods and tasks demonstrates the versatility and effectiveness of our allocation approach. It stresses the potential of our strategy as a generalizable framework for optimizing KV cache eviction, enhancing various existing methods.

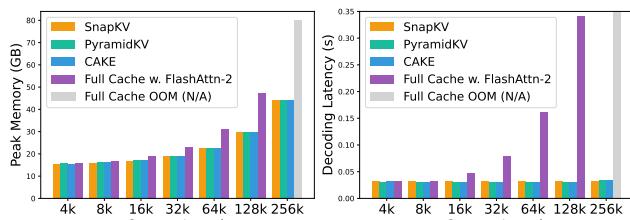


Figure 7: Peak memory usage and decoding latency on A100 80GB.



Figure 8: Performance on methods w/ or w/o preference-prioritized cache allocation strategy, abbreviated as “P2A.”

Table 2: Performance comparison with different allocation strategies. “P2A” denotes our preference-prioritized adaptive allocation.

Strategy	Avg.
Uniform	28.36
Pyramid	28.69
Random	28.3±0.2
P2A	29.29

Table 3: Performance comparison with different eviction indicators.

Indicator	Avg.
Mean	28.82
Var	28.68
Mean · Var	28.6
Mean + Var	29.29

5.6 ABLATION STUDIES

In this section, we present a series of ablation studies on LongBench to evaluate the effectiveness of our proposed allocation strategy and eviction indicator. We use Llama2-7B-Chat with cache budget $B_{\text{total}} = 128L$ as the default setting. Additional ablation studies are provided in Appendix I.

Effectiveness of Proposed Allocation Strategy. To assess our preference-prioritized adaptive allocation strategy, we compare it with various cache allocation strategies, including uniform, pyramid-shaped, and random allocation. As shown in Table 2, the pyramid-shape allocation strategy performs better than the uniform allocation strategy in this case. However, it may become ineffective in many scenarios, as we verified in the previous experimental section. Notably, the random method even outperforms the uniform strategy in some cases. This indicates that the conservative allocation strategy of uniform distribution fails to fully utilize the limited memory budget. Our allocation strategy consistently outperforms all baselines, as our method comprehensively measures the KV cache preference by capturing complex layer-specific attention patterns both spatially and temporally.

Effectiveness of Proposed Eviction Indicator. We further investigate effective eviction indicators to identify the most crucial tokens for preserving model performance. We compare different metrics for the eviction indicator: mean only, variance only, and combinations of mean and variance (multiplication and addition). The results in Table 3 suggest that, compared to variance only, tokens selected by mean attention are more conducive to maintaining performance. While the multiplicative combination of mean and variance slightly underperforms the individual metrics, the additive combination achieves the best performance overall. As we analyze, using mean attention effectively captures tokens with long-term importance, ensuring the retention of consistently relevant information. On the other hand, incorporating variance helps identify positions with the most significant changes, thus aiding in maintaining the attention distribution. The additive combination optimally balances these two aspects, leading to more informed eviction decisions.

6 CONCLUSION

In this paper, we propose Cascading and Adaptive KV cache Eviction (CAKE), a novel approach for optimizing KV cache evicting in LLMs. CAKE dynamically allocates cache sizes by leveraging a global view of layer-specific attention patterns, using cascading cache management guided by layer preferences. Additionally, CAKE introduces a new eviction indicator that accounts for both the long-term influence and temporal variability of token importance, allowing for more informed token selection. Experiments on the LongBench and NeedleBench benchmarks highlight CAKE’s superior performance across different models and memory constraints, especially in low-memory scenarios. Our method enhances both LLM performance on long-context tasks and inference efficiency.

540 REFERENCES
541

- 542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Mistral AI. Large enough, 2024. URL [https://mistral.ai/news/
546 mistral-large-2407/](https://mistral.ai/news/mistral-large-2407/).
- 547 Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit
548 Sanghi. Gqa: Training generalized multi-query transformer models from multi-head check-
549 points. *arXiv preprint arXiv:2305.13245*, 2023.
- 550 Silas Alberti, Niclas Dern, Laura Thesing, and Gitta Kutyniok. Sumformer: Universal approxi-
551 mation for efficient transformers. In *Topological, Algebraic and Geometric Learning Workshops
552 2023*, pp. 72–86. PMLR, 2023.
- 553 Anthropic. The claude 3 model family: Opus, sonnet, haiku, March 2024. URL [https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/
Model_Card_Claude_3.pdf](https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/
554 Model_Card_Claude_3.pdf).
- 555 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,
556 Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long
557 context understanding. *arXiv preprint arXiv:2308.14508*, 2023.
- 558 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
559 *arXiv preprint arXiv:2004.05150*, 2020.
- 560 Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong,
561 Baobao Chang, Junjie Hu, et al. Pyramidkv: Dynamic kv cache compression based on pyramidal
562 information funneling. *arXiv preprint arXiv:2406.02069*, 2024.
- 563 Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, and Christopher Ré. Scatterbrain: Uni-
564 fying sparse and low-rank attention. *Advances in Neural Information Processing Systems*, 34:
565 17413–17426, 2021.
- 566 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
567 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, et al. Vicuna: An open-source chatbot
568 impressing GPT-4 with 90%* ChatGPT quality, March 2023. URL [https://lmsys.org/
blog/2023-03-30-vicuna/](https://lmsys.org/
569 blog/2023-03-30-vicuna/).
- 570 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse
571 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 572 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas
573 Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention
574 with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- 575 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning, 2023. URL
576 <https://arxiv.org/abs/2307.08691>.
- 577 Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. Get more
578 with less: Synthesizing recurrence with kv cache compression for efficient llm inference. *arXiv
579 preprint arXiv:2402.09398*, 2024.
- 580 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
581 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
582 *arXiv preprint arXiv:2407.21783*, 2024.
- 583 Qichen Fu, Minsik Cho, Thomas Merth, Sachin Mehta, Mohammad Rastegari, and Mahyar Na-
584 jibi. Lazyllm: Dynamic token pruning for efficient long context llm inference. *arXiv preprint
585 arXiv:2407.14057*, 2024.

- 594 Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. Model tells
 595 you what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*,
 596 2023.
- 597 Chi Han, Qifan Wang, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. Lm-infinite: Simple
 598 on-the-fly length generalization for large language models. *arXiv preprint arXiv:2308.16137*,
 599 2023.
- 600 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
 601 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
 602 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 603 Ehsan Kamalloo, Nouha Dziri, Charles LA Clarke, and Davood Rafiei. Evaluating open-domain
 604 question answering in the era of large language models. *arXiv preprint arXiv:2305.06984*, 2023.
- 605 Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo
 606 Zhao. Gear: An efficient kv cache compression recipefor near-lossless generative inference of
 607 llm. *arXiv preprint arXiv:2403.05527*, 2024.
- 608 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are
 609 rnns: Fast autoregressive transformers with linear attention. In *International conference on ma-*
610 chine learning, pp. 5156–5165. PMLR, 2020.
- 611 Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv*
612 preprint arXiv:2001.04451, 2020.
- 613 Mo Li, Songyang Zhang, Yunxin Liu, and Kai Chen. Needlebench: Can llms do retrieval and
 614 reasoning in 1 million context window? *arXiv preprint arXiv:2407.11963*, 2024a.
- 615 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle
 616 Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before
 617 generation. *arXiv preprint arXiv:2404.14469*, 2024b.
- 618 Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong
 619 Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-
 620 of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024a.
- 621 Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. Mini-
 622 cache: Kv cache compression in depth dimension for large language models. *arXiv preprint*
623 arXiv:2405.14366, 2024b.
- 624 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
 625 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*
626 Association for Computational Linguistics, 12:157–173, 2024c.
- 627 Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios
 628 Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance
 629 hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing*
630 Systems, 36, 2024d.
- 631 Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi
 632 Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint*
633 arXiv:2402.02750, 2024e.
- 634 Piotr Nawrot, Adrian Łaćucki, Marcin Chochowski, David Tarjan, and Edoardo M Ponti. Dy-
 635 namic memory compression: Retrofitting llms for accelerated inference. *arXiv preprint*
636 arXiv:2403.09636, 2024.
- 637 Matanel Oren, Michael Hassid, Yossi Adi, and Roy Schwartz. Transformers are multi-state rnns.
 638 *arXiv preprint arXiv:2401.06104*, 2024.
- 639 Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and
 640 Bo Dai. Combiner: Full attention transformer with sparse computation cost. *Advances in Neural*
641 Information Processing Systems, 34:22470–22482, 2021.

- 648 Siyu Ren and Kenny Q Zhu. On the efficacy of eviction policy for key-value constrained generative
 649 language model inference. *arXiv preprint arXiv:2402.06262*, 2024.
 650
- 651 Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse
 652 attention with routing transformers. *Transactions of the Association for Computational Linguistics*,
 653 9:53–68, 2021.
- 654 Noam Shazeer. Fast transformer decoding: One write-head is all you need. *arXiv preprint
 655 arXiv:1911.02150*, 2019.
 656
- 657 Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wenhui Wang, Shuming Ma, Quanlu Zhang, Jianyong
 658 Wang, and Furu Wei. You only cache once: Decoder-decoder architectures for language models.
 659 *arXiv preprint arXiv:2405.05254*, 2024.
- 660 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
 661 Pathak, Laurent Sifre, Morgane Rivière, Minir Sanjay Kale, Juliette Love, et al. Gemma: Open
 662 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
 663
- 664 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL <https://qwenlm.github.io/blog/qwen2.5/>.
 665
- 666 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
 667 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
 668 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
 669
- 670 Zhongwei Wan, Xinjian Wu, Yu Zhang, Yi Xin, Chaofan Tao, Zhihong Zhu, Xin Wang, Siqi Luo,
 671 Jing Xiong, and Mi Zhang. D2o: Dynamic discriminative operations for efficient generative
 672 inference of large language models. *arXiv preprint arXiv:2406.13035*, 2024.
- 673 Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention
 674 with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
 675
- 676 Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan
 677 Liu, Song Han, and Maosong Sun. Inflm: Unveiling the intrinsic capacity of llms for under-
 678 standing extremely long sequences with training-free memory. *arXiv preprint arXiv:2402.04617*,
 679 2024.
- 680 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
 681 language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
 682
- 683 Yuhui Xu, Zhanming Jie, Hanze Dong, Lei Wang, Xudong Lu, Aojun Zhou, Amrita Saha, Caiming
 684 Xiong, and Doyen Sahoo. Think: Thinner key cache by query-driven pruning. *arXiv preprint
 685 arXiv:2407.21018*, 2024.
- 686 Dongjie Yang, XiaoDong Han, Yan Gao, Yao Hu, Shilin Zhang, and Hai Zhao. Pyramidinfer: Pyra-
 687 mid kv cache compression for high-throughput llm inference. *arXiv preprint arXiv:2405.12532*,
 688 2024a.
- 689 June Yong Yang, Byeongwook Kim, Jeongin Bae, Beomseok Kwon, Gunho Park, Eunho Yang,
 690 Se Jung Kwon, and Dongsoo Lee. No token left behind: Reliable kv cache compression via
 691 importance-aware mixed precision quantization. *arXiv preprint arXiv:2402.18096*, 2024b.
 692
- 693 Yuxuan Yue, Zhihang Yuan, Haojie Duanmu, Sifan Zhou, Jianlong Wu, and Liqiang Nie. Wkvquant:
 694 Quantizing weight and key/value cache for large language models gains more. *arXiv preprint
 695 arXiv:2402.12065*, 2024.
- 696 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago
 697 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for
 698 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
 699
- 700 Tianyi Zhang, Faisal Ladha, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B
 701 Hashimoto. Benchmarking large language models for news summarization. *Transactions of the
 702 Association for Computational Linguistics*, 12:39–57, 2024a.

- 702 Yuxin Zhang, Yuxuan Du, Gen Luo, Yunshan Zhong, Zhenyu Zhang, Shiwei Liu, and Rongrong Ji.
703 Cam: Cache merging for memory-efficient llms inference. In *Forty-first International Conference*
704 *on Machine Learning*, 2024b.
- 705 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,
706 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient gen-
707 erative inference of large language models. *Advances in Neural Information Processing Systems*,
708 36, 2024c.
- 709 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
710 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*
711 *preprint arXiv:2303.18223*, 2023.
- 712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 **A ADDITIONAL RELATED WORKS**

757
Efficient Attention. Numerous strategies have been devised to tackle the high time and space demands of the attention layer. These strategies include low-rank techniques to transform the original 758 attention mechanism into one with linear or near-linear complexity (Choromanski et al., 2020; Wang 759 et al., 2020; Katharopoulos et al., 2020; Alberti et al., 2023), sparse attention to confine the attention 760 to specific, predetermined or adaptively learned patterns (Child et al., 2019; Kitaev et al., 2020; Roy 761 et al., 2021), and their integration for a more comprehensive solution (Beltagy et al., 2020; Zaheer 762 et al., 2020; Chen et al., 2021; Ren et al., 2021). Another group minimizes KV cache size. This 763 is achieved by either re-engineering the architecture of multi-head attention or by innovating the 764 update process of the KV cache. Multi-Query Attention (MQA) (Shazeer, 2019) and Group-Query 765 Attention (GQA) (Ainslie et al., 2023) reduce KV heads by sharing KV pairs. Multi-head Latent 766 Attention (MLA) (Liu et al., 2024a) furthers efficiency by compressing keys and values into latent 767 vectors, while YOCO (Sun et al., 2024) employs cross-attention to recycle the shared KV cache. Dy- 768 namic Memory Compression (DMC) (Nawrot et al., 2024) dynamically consolidates the KV cache 769 during the training process, thereby optimizing resource utilization.
770

771
Orthogonal KV Cache Compression. Beyond the eviction of the KV cache, numerous strategies 772 have been investigated for its compression. Specifically, quantization methods (Liu et al., 2024e; Yue 773 et al., 2024; Kang et al., 2024) have been applied to condense the KV cache into a lower-bitwidth 774 format. Moreover, a fusion of KV cache eviction with quantization (Yang et al., 2024b), pruning Xu 775 et al. (2024), low-rank factorization (Dong et al., 2024), or cache merging (Zhang et al., 2024b; Liu 776 et al., 2024b; Wan et al., 2024) has been explored to achieve a more compact form of the cache, 777 preserving some of the discarded data. Additionally, there are initiatives to transfer the evicted KV 778 cache to auxiliary memory resources for potential redeployment in computation (Xiao et al., 2024; 779 Fu et al., 2024). It is crucial to acknowledge that these compression techniques are compatible 780 with KV cache eviction and, when integrated, can amplify their individual benefits, leading to more 781 robust solutions for easing the KV cache load.
782

783 **B LIMITATION AND FUTURE WORK**
784

785 CAKE introduces a novel perspective on KV cache management by considering the model’s vary- 786 ing attention patterns across layers and comprehensively measuring layer-specific preferences for 787 KV cache allocation. This approach enables effective adaptation to different models and input con- 788 texts, optimizing the use of limited memory resources. Despite these strengths, several limitations 789 and areas for future work remain. While CAKE improves cache management at the layer level, it 790 does not address finer-grained dynamics within layers. Future work could explore integrating head- 791 level attention patterns with layer-specific variations, potentially offering more nuanced control over 792 cache allocation. Additionally, although we demonstrate CAKE’s compatibility with existing evic- 793 tion methods and quantization methods, there is potential to further enhance memory efficiency by 794 combining it with other KV cache optimization techniques such as cache merging (Liu et al., 2024b) 795 and pruning (Xu et al., 2024). This integration could significantly expand CAKE’s applicability and 796 resource-saving capabilities. In the future, we will continue to explore strategies for efficient LLMs, 797 focusing on both memory efficiency and computational optimization to further push the boundaries 798 of resource-constrained inference.
799

800 **C MORE IMPLEMENTATION DETAILS**
801

802 In this section, we provide more detailed implementation information for CAKE. Our algorithm is 803 divided into two phases: prompt prefilling and token decoding. During prompt prefilling, we utilize 804 preference-guided cascading cache management (Section 4.2) based on the preference-prioritized 805 adaptive allocation strategy (Section 4.1) to allocate appropriate cache sizes for different layers. 806 During preference-guided cascading cache management, we manage the KV cache by updating it 807 using eviction operation based on the attention-shift tolerant eviction indicator (Section 4.3). For the 808 preference-prioritized adaptive allocation strategy, the temperature parameters τ_1 and τ_2 are set to 809 $0.2 \sim 2$ and $0.4 \sim 3$ respectively, optimized through grid search. For our experimental models, we 810 aggregate each quantized value to represent its layer characteristics. For the attention-shift tolerant 811 eviction indicator, we fix $\gamma = 200$. Following SnapKV (Li et al., 2024b), we use an observation
812

810 window of size $S_w = 32$ and utilize a pooling layer clustering information. In the decoding stage,
 811 once the attention mechanisms for layers have been well established, we fix the budget size for each
 812 layer according to the allocation results obtained from the preference-prioritized adaptive allocation.
 813 **Listing 1 provides PyTorch-style pseudo-code for CAKE’s implementation with FlashAttention.**
 814

815 Listing 1: Implementation of CAKE in pseudo PyTorch style.

```

1  class CakeCache(Cache):
2      def __init__(self):
3          self.pref_scores = []
4          self.evict_indicators = []
5          self.layers_budget = []
6          ...
7          # Layer 0 to L-1
8      def attention_forward(self, hidden_states, ..., past_key_value: Optional[CakeCache] = None):
9          ...
10         # Compute query_states, key_states, value_states
11         bsz, q_len, _ = hidden_states.size()
12         ...
13         # Compute local attention to the current layer
14         local_attn = compute_local_attention(query_states[:, :-self.window_size:, :], key_states,
15                                             value_states)
16         observed_attn = local_attn[:, :, :-self.window_size:, :-self.window_size]
17         if prefill:
18             # Calculate preference score
19             dispersion = calculate_attn_dispersion(observed_attn)
20             shift = calculate_attn_shift(observed_attn)
21             past_key_value.pref_scores[layer] = (dispersion **(1/tau1) * shift **(1/tau2))
22             # Calculate eviction indicator
23             attn_mean = mean(observed_attn)
24             attn_var = var(observed_attn)
25             indicator = attn_mean + self.gamma * attn_var
26             indicator = poolId(indicator)
27             # Update
28             past_key_value.pref_scores.append(pref_score)
29             past_key_value.evict_indicators(pref_score, indicator)
30             past_key_value = cake(past_key_value, q_len, self.layer_id)
31         else:
32             indicator = poolId(observed_attn)
33             past_key_value.update_score(indicator)
34             ...
35             # Compute attention output using flash_attention
36             attn_output = flash_attention_forward(query_states, key_states, value_states, ...)
37         return attn_output, past_key_value
38
39     def cake(self, past_key_value, seq_len, current_layer_idx):
40         if seq_len <= self.cache_size - self.window_size:
41             return past_key_values
42             # Cache budget management
43             if prefill:
44                 pref_scores = past_key_values.pref_scores
45                 layer_budgets = [pref_score / sum(pref_scores) * self.total_budget for pref_score in
46 pref_scores]
47                 # The following loop can be executed in parallel for each layer
48                 for layer_idx, budget_size in enumerate(layer_budgets):
49                     budget_size = min(budget_size, seq_len - self.window_size)
50                     past_key_value = evict_layer_kvcache(past_key_value, layer_idx, budget_size)
51                     past_key_value.layers_budget[layer_idx] = budget_size
52             else:
53                 budget_size = past_key_values.layer_budget[current_layer_idx]
54                 past_key_values = self.evict_layer_kvcache(past_key_values, current_layer_idx,
55                 budget_size)
56             return past_key_values
57
58     def evict_layer_kvcache(self, past_key_values, layer_idx, budget_size):
59         bsz, num_key_value_heads, seq_len, head_dim = past_key_values.key_cache[layer_idx].shape
60         num_key_value_groups = self.num_heads // num_key_value_heads
61         indicator = past_key_values.evict_indicators[layer_idx]
62         indices = indicator.topk(budget, dim=-1).indices
63         # Update the eviction indicators
64         past_key_values.evict_indicators[layer_idx] = compress_kv(indicator, indices, self.
65         window_size)
66         indices = indices.unsqueeze(-1).expand(-1, -1, -1, head_dim)
67         # Update the KV cache
68         past_key_values.key_cache[layer_idx] = compress_kv(key_states, indices, self.window_size)
69         past_key_values.value_cache[layer_idx] = compress_kv(value_states, indices, self.
70         window_size)
71
72         return past_key_values

```

864 **D PROOFS**
 865

866 *Proof of Proposition 1.* Given the budget calculation formula in Eq.(9), we have:
 867

$$868 \quad B_l^{(m)} = \frac{\mathcal{P}_l}{\sum_{k=0}^m \mathcal{P}_k} \cdot B_{\text{total}}, \quad (13)$$

$$870 \quad B_l^{(m+1)} = \frac{\mathcal{P}_l}{\sum_{k=0}^{m+1} \mathcal{P}_k} \cdot B_{\text{total}}. \quad (14)$$

872 And the inequality Eq.(8) is equivalent to:
 873

$$874 \quad \frac{\mathcal{P}_l}{\sum_{k=0}^{m+1} \mathcal{P}_k} < \frac{\mathcal{P}_l}{\sum_{k=0}^m \mathcal{P}_k}. \quad (15)$$

877 Since preference scores are always positive, this inequality holds. Therefore, we complete the proof
 878 that the allocated budget size monotonically decreases. \square
 879

880 *Proof of Theorem 1.* Recall that the eviction operation on layer l essentially retains KV pairs based
 881 on top- B_l selection indices \mathbf{D}_l on the values of \mathbf{I}_l , where B_l is the target budget size. Therefore,
 882 proving Eq.(10) is equivalent to demonstrating:

$$883 \quad \text{TopK}(\mathbf{I}_l^{(m)}, B_l^{(m)})_{m=l}^{L-1} = \text{TopK}(\mathbf{I}_l, B_l). \quad (16)$$

885 The left-hand side of the above equation can be expanded as a cascading process:
 886

$$\begin{aligned} 887 \quad \mathbf{D}_l^{(l)} &= \text{TopK}(\mathbf{I}_l^{(l)}, B_l^{(l)}), \quad \mathbf{D}_l^{(l)} \in \mathbb{R}^{B_l^{(l)}}, \\ 888 \quad \mathbf{D}_l^{(l+1)} &= \text{TopK}(\mathbf{I}_l^{(l+1)}, B_l^{(l+1)}), \quad \mathbf{D}_l^{(l+1)} \in \mathbb{R}^{B_l^{(l+1)}}, \\ 889 \quad &\vdots \\ 891 \quad \mathbf{D}_l^{(L-1)} &= \text{TopK}(\mathbf{I}_l^{(L-1)}, B_l^{(L-1)}), \quad \mathbf{D}_l^{(L-1)} \in \mathbb{R}^{B_l^{(L-1)}}. \end{aligned} \quad (17)$$

893 For the vanilla method, \mathbf{D}_l is obtained by a single TopK operation:
 894

$$895 \quad \mathbf{D}_l = \text{TopK}(\mathbf{I}_l, B_l). \quad (18)$$

896 Given these definitions, our goal simplifies to proving that:
 897

$$898 \quad \mathbf{D}_l^{(L-1)} = \mathbf{D}_l. \quad (19)$$

900 To prove this equality, we first show that these two sets have the same number of elements. Note that
 901 $|\mathbf{D}_l^{(L-1)}| = B_l^{(L-1)}$ and $|\mathbf{D}_l| = B_l$. At the final step L , both methods have access to the preference
 902 scores of all layers, and the budget calculation for layer l is identical in both cases:

$$903 \quad B_l^{(L-1)} = B_l = \frac{\mathcal{P}_l}{\sum_{k=0}^{L-1} \mathcal{P}_k} \cdot B_{\text{total}}, \quad (20)$$

906 where \mathcal{P}_l is the preference score for layer l , and B_{total} is the total cache budget. Thus, $|\mathbf{D}_l^{(L-1)}| =$
 907 $|\mathbf{D}_l|$. Next, we further prove that elements in $\mathbf{D}_l^{(L-1)}$ and \mathbf{D}_l are identical. Given that $B_l^{(m)}$ is
 908 decreasing as m increases, for any stage $m \in [l, L-1]$, we have:

$$909 \quad \text{TopK}(\mathbf{I}_l^{(m)}, B_l^{(m)}) \supset \text{TopK}(\mathbf{I}_l, B_l^{(m+1)}) = \text{TopK}(\mathbf{I}_l^{(m)}, B_l^{(m+1)}). \quad (21)$$

911 Therefore, we have:
 912

$$913 \quad \mathbf{D}_l^{(l)} \supset \mathbf{D}_l^{(l+1)} \supset \dots \supset \mathbf{D}_l^{(L-1)} = \text{TopK}(\mathbf{I}_l^{(L-1)}, B_l^{(L-1)}) = \text{TopK}(\mathbf{I}_l, B_l). \quad (22)$$

914 In the final stage, set $\mathbf{D}_l^{(L-1)}$ contains exactly the $B_l^{(L-1)} = B_l$ highest-scoring elements from \mathbf{I}_l ,
 915 which is identical to \mathbf{D}_l . Thus, we have proven that $\mathbf{D}_l^{(L-1)} = \mathbf{D}_l$, which completes the proof of
 916 Theorem 1. \square
 917

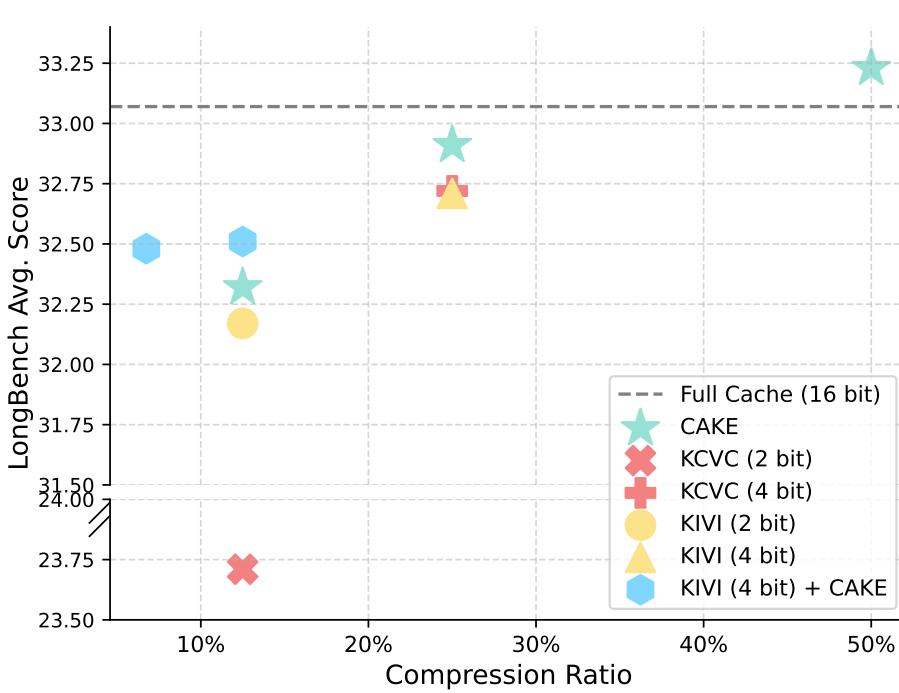


Figure 9: Comparison of CAKE, KV cache quantization methods (KIVI and KCVC), and their combination (KIVI with CAKE) on Llama2-7B-Chat using LongBench dataset.

E COMPARISON AND COMPATIBILITY WITH QUANTIZATION METHODS

While CAKE focuses on selectively dropping less important KV pairs to reduce memory footprint, quantization approaches aim to compress the cache by reducing the bit precision of stored values. These two strategies are orthogonal and potentially complementary in nature.

In this section, we evaluate CAKE against two quantization approaches: KIVI (Liu et al., 2024e), a state-of-the-art method using asymmetric KCVT quantization (quantizes Key cache per-channel and Value cache per-token), and a computationally efficient channel-wise quantization variant for both key and value caches (denoted as KCVC). As is shown in Figure 9, under matched compression ratios, CAKE outperforms full cache (33.23 vs. 33.07) and both quantization methods at 25% and 12.5% compression ratios. Notably, when combining KIVI-INT4 with CAKE, it demonstrates significant advantages in maintaining performance at high compression ratios. For instance, KIVI-INT4 with CAKE achieves 32.51 at 12.5% compression ratio, outperforming KIVI-INT2 (32.17) while using the same storage overhead. Even at a very aggressive 6.25% compression ratio, the combined approach still maintains strong performance (32.48), showing the effectiveness of integrating eviction and quantization strategies.

F EXTENDED EXPERIMENTAL RESULTS ON LONGBENCH

In this section, we provide comprehensive experimental results on LongBench (Bai et al., 2023), a benchmark focused on long-context understanding with input lengths ranging from 1,235 to 18,409 tokens. We present detailed performance evaluations across cache budgets from 64L to 2048L for three base models: Llama2-7B-Chat (Touvron et al., 2023), Llama3.1-8B-Instruct (Dubey et al., 2024), and Mistral-7B-Instruct-v0.3 (Jiang et al., 2023) (Appendix F.1). To demonstrate CAKE’s generalizability, we conduct additional experiments on two more LLM architectures: Qwen2.5-7B-Instruct (Team, 2024) and Gemma-7B-Instruct (Team et al., 2024) (Appendix F.2). Furthermore, we evaluate CAKE’s scalability on larger models ranging from 13B to 70B parameters, including Llama2-13B-Chat, Qwen2.5-32B-Instruct, and Llama3-70B-Instruct (Appendix F.3).

972
973
974 Table 4: LongBench Task Type.
975
976
977
978
979
980
981
982

Task Type	#English Task	#Code Task
Multi-document QA	3	-
Single-document QA	3	-
Summarization	3	-
Few-shot learning	3	-
Synthetic Tasks	2	-
Code Completion	-	2

983 LongBench is a long context dataset comprised of six task categories covering key long-text applica-
 984 tion scenarios: single-document QA, multi-document QA, summarization, few-shot learning, syn-
 985 thetic tasks and code completion. Our experiments are conducted on all of its 14 English subtasks
 986 and 2 code subtasks, spanning throughout all six categories (Table 4), with input lengths ranging
 987 from 1,235 to 18,409 tokens.

988 F.1 DETAILED PERFORMANCE ANALYSIS ACROSS CACHE BUDGETS

989 Tables 5, 6 and 7 present the detailed LongBench results of CAKE and comparative methods applied
 990 to Llama2-7B-Chat, Llama3.1-8B-Instruct, and Mistral-7B-Instruct-v0.3, respectively.

991 **Results on Llama2-7B-Chat.** As shown in Table 5, CAKE consistently outperforms other KV
 992 eviction methods in terms of average score. Additionally, CAKE achieves top or second across each
 993 individual subtask and at all cache sizes. Notice that when using with $512L$ or more total budget
 994 size, CAKE achieves evaluation scores that are comparable with full cache across all subtasks.

995 **Results on Llama3.1-8B-Instruct.** The evaluation results are detailed in Table 6. The Llama3.1
 996 series models have all incorporated the Grouped Query Attention (GQA) mechanism, and are able
 997 to support longer context length. The results show that CAKE continues to perform outstandingly
 998 across all subtasks in models based on the GQA structure.

999 **Results on Mistral-7B-Instruct-v0.3.** Table 7 contains the results of Mistral-7B-Instruct-v0.3.
 1000 Around all 16 subtasks, CAKE shows consistent advantage throughout all the cache size budgets
 1001 we have tested. Especially when the budget is limited, say $64L$, CAKE, equipped with the robust
 1002 indicator and strategy, expresses relatively strong resistance to score decline.

1003 The results as a whole demonstrate that CAKE, compared with other methods, consistently outper-
 1004 forms across all task types in LongBench when applied to the tested models and under various cache
 1005 size budgets ranging from $64L$ to $2048L$. This demonstrates CAKE’s effectiveness and broad ap-
 1006 plicability in KV cache-efficient long-context processing across various domains and across widely
 1007 used open-source LLMs.

1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

Table 5: Performance comparison over 16 datasets of LongBench on Llama2-7B-Chat for cache budgets from $64L$ to $2048L$. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code	
	NrrQA	Qasper	MF-en	HolpiQA	2WikiMQA	Masique	GovReport	QMSum	MultiNews	TREC	TriviaQA	SAMSum	PCount	PR-en	Lcc	RB-P	Avg.			
	Llama2-7B-Chat, $B_{\text{total}} = \text{Full}$																			
Full	20.68	20.19	36.69	27.83	31.45	8.21	27.09	20.71	26.24	64.0	83.09	41.41	2.94	7.75	58.51	52.29	33.07			
Llama2-7B-Chat, $B_{\text{total}} = 64L$																				
StreamingLLM	7.35	<u>16.62</u>	6.84	10.73	14.68	2.84	13.34	18.42	12.76	17.50	19.38	11.44	0.25	4.50	20.87	14.91	12.03			
H2O	12.77	13.76	21.00	19.16	26.53	4.47	14.09	<u>20.01</u>	18.49	32.00	70.92	34.79	2.75	3.00	42.09	<u>41.36</u>	23.57			
TOVA	8.64	14.06	12.61	13.71	13.56	3.26	11.92	16.59	13.73	36.00	<u>74.88</u>	31.88	2.15	4.00	36.72	31.94	20.35			
SnapKV	<u>12.63</u>	14.95	19.93	15.48	22.94	4.78	<u>15.54</u>	19.50	15.10	34.50	72.02	34.16	2.00	4.12	40.03	36.29	22.75			
PyramidKV	10.94	15.53	28.43	23.53	28.82	<u>5.10</u>	15.33	19.79	15.29	34.50	76.16	<u>34.84</u>	<u>2.88</u>	2.56	<u>44.82</u>	38.16	24.79			
CAKE	12.62	16.64	<u>28.04</u>	21.46	<u>28.00</u>	5.19	17.20	20.16	<u>17.23</u>	<u>35.00</u>	74.30	35.02	<u>3.33</u>	4.50	45.86	42.33	25.43			
Llama2-7B-Chat, $B_{\text{total}} = 128L$																				
StreamingLLM	13.57	13.33	22.39	20.16	17.12	6.38	16.72	19.58	15.37	31.50	74.84	33.00	1.67	4.00	43.21	40.45	23.33			
H2O	<u>14.56</u>	14.59	26.52	20.05	<u>28.73</u>	5.00	16.20	20.30	<u>20.34</u>	38.00	72.75	37.10	2.33	3.50	49.68	47.37	26.06			
TOVA	13.56	15.46	20.28	22.22	25.01	4.79	15.72	19.11	15.58	42.50	79.58	37.06	<u>3.45</u>	4.61	46.34	39.70	25.31			
SnapKV	13.70	16.27	32.52	22.76	28.59	7.56	18.87	19.96	20.05	43.50	79.90	36.74	2.79	7.00	51.86	47.41	28.09			
PyramidKV	13.45	16.61	31.21	<u>23.74</u>	28.43	6.89	<u>18.88</u>	<u>20.56</u>	19.81	<u>44.50</u>	<u>80.32</u>	37.82	2.29	<u>8.00</u>	51.19	<u>47.56</u>	<u>28.20</u>			
CAKE	15.15	<u>16.38</u>	<u>32.17</u>	25.25	31.09	<u>7.00</u>	19.47	21.07	20.36	48.50	80.66	<u>37.77</u>	4.42	9.00	<u>51.48</u>	48.90	29.29			
Llama2-7B-Chat, $B_{\text{total}} = 256L$																				
StreamingLLM	13.20	13.45	24.36	21.50	24.08	5.74	17.28	18.97	18.12	44.00	78.43	37.55	1.58	4.50	53.85	47.03	26.48			
H2O	<u>15.82</u>	14.72	27.26	21.07	27.00	6.32	19.75	20.29	<u>22.23</u>	45.50	79.64	37.93	<u>2.93</u>	4.00	53.91	50.33	28.04			
TOVA	13.78	14.69	23.81	<u>25.92</u>	28.93	7.24	18.00	19.48	19.04	<u>58.00</u>	<u>82.33</u>	39.33	2.69	7.00	51.42	47.16	28.68			
SnapKV	15.47	<u>15.94</u>	<u>34.51</u>	25.17	<u>30.10</u>	9.30	<u>20.44</u>	20.89	21.78	57.50	81.71	<u>39.22</u>	2.80	<u>7.50</u>	<u>55.93</u>	<u>50.10</u>	<u>30.52</u>			
PyramidKV	15.13	15.86	33.93	25.58	29.51	<u>9.23</u>	20.35	21.22	22.01	<u>58.00</u>	82.39	38.47	2.15	<u>7.50</u>	55.81	49.53	30.42			
CAKE	15.93	16.80	36.25	26.26	30.54	9.01	20.79	<u>20.98</u>	22.56	61.00	82.00	39.01	3.24	8.00	56.97	50.04	31.21			
Llama2-7B-Chat, $B_{\text{total}} = 512L$																				
StreamingLLM	12.61	14.49	25.87	24.21	23.86	5.40	19.51	19.19	21.43	56.00	80.56	38.92	1.70	2.50	55.57	48.56	28.15			
H2O	15.86	15.89	30.14	22.14	28.89	7.01	21.40	20.89	23.54	52.50	81.81	<u>40.55</u>	3.19	5.00	55.68	51.07	29.72			
TOVA	13.63	15.47	26.26	<u>26.55</u>	<u>31.25</u>	7.63	19.81	19.88	22.19	62.50	83.10	40.63	2.72	<u>8.00</u>	55.83	50.28	30.36			
SnapKV	<u>16.42</u>	17.19	<u>36.52</u>	26.32	31.15	<u>9.09</u>	21.33	21.11	23.58	64.00	81.82	39.29	3.33	8.50	<u>58.00</u>	52.46	<u>31.88</u>			
PyramidKV	16.91	<u>18.30</u>	36.02	26.74	29.53	9.13	<u>21.53</u>	<u>21.22</u>	<u>23.72</u>	63.50	81.88	39.76	3.00	<u>8.00</u>	58.02	51.14	31.77			
CAKE	16.36	19.31	37.82	26.51	31.34	8.10	21.77	21.42	24.36	64.00	<u>82.63</u>	39.62	3.33	<u>8.00</u>	57.73	<u>52.33</u>	32.16			
Llama2-7B-Chat, $B_{\text{total}} = 1024L$																				
StreamingLLM	13.73	16.75	26.28	25.97	25.67	5.87	22.21	19.27	24.07	61.00	80.88	40.90	2.28	1.00	56.71	48.79	29.46			
H2O	16.58	17.67	32.04	26.01	28.73	7.81	22.87	20.99	25.05	60.00	83.02	40.06	2.90	3.50	57.57	52.02	31.05			
TOVA	15.00	17.32	29.57	27.48	31.71	7.04	21.86	20.54	24.78	63.50	83.35	41.78	2.63	<u>8.00</u>	56.69	52.11	31.46			
SnapKV	<u>18.05</u>	19.73	38.07	<u>27.57</u>	30.86	<u>8.41</u>	<u>23.41</u>	20.79	25.22	63.50	82.13	40.98	<u>3.33</u>	7.50	58.20	<u>52.23</u>	32.50			
PyramidKV	17.94	<u>20.68</u>	38.85	27.25	29.99	7.91	23.17	21.55	25.45	63.50	82.97	40.79	3.57	7.50	57.72	51.93	<u>32.55</u>			
CAKE	18.35	20.93	<u>38.49</u>	27.99	<u>31.15</u>	8.59	24.29	<u>21.06</u>	<u>25.40</u>	63.50	<u>83.33</u>	<u>41.30</u>	3.11	9.50	<u>57.82</u>	52.93	32.98			
Llama2-7B-Chat, $B_{\text{total}} = 2048L$																				
StreamingLLM	16.95	17.98	32.25	<u>28.15</u>	30.31	7.66	23.79	19.97	25.75	63.50	83.17	40.36	2.05	5.17	57.36	51.30	31.61			
H2O	18.35	20.26	36.71	28.65	30.85	8.33	25.05	<u>21.20</u>	26.31	63.00	83.32	40.78	2.60	7.00	58.14	52.92	32.72			
TOVA	18.73	20.52	35.38	28.10	30.91	<u>8.87</u>	24.86	20.87	<u>26.22</u>	64.00	83.13	41.99	2.75	7.00	58.07	<u>52.77</u>	32.76			
SnapKV	<u>18.80</u>	22.18	<u>38.18</u>	27.72	<u>30.92</u>	8.37	25.87	21.35	25.92	64.00	83.13	40.95	<u>3.32</u>	<u>8.50</u>	<u>58.18</u>	<u>52.23</u>	33.10			
PyramidKV	18.51	22.68	37.04	27.76	27.84	7.95	26.42	20.95	26.08	64.00	<u>83.26</u>	40.72	3.05	7.75	58.80	52.55	32.83			
CAKE	19.25	22.16	38.70	27.72	31.58	8.91	<u>26.19</u>	20.70	26.19	64.00	83.25	<u>41.02</u>	3.86	9.00	57.92	52.76	33.33			

Table 6: Performance comparison over 16 datasets of LongBench on Llama3.1-8B-Instruct for cache budgets from $64L$ to $2048L$. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code	
	NnQA	Qasper	MF-en	HopQA	2WikiQA	Musique	GovReport	QMSum	MultNews	TREC	TriviaQA	SAMSum	pCount	PR-en	Lcc	RB-P	Avg.			
Llama3.1-8B-Instruct, $B_{\text{total}} = \text{Full}$																				
Full	31.06	45.43	53.78	55.04	47.14	31.29	34.87	25.33	27.49	72.5	91.65	43.81	6.0	99.5	63.36	56.65	49.06			
Llama3.1-8B-Instruct, $B_{\text{total}} = 64L$																				
StreamingLLM	22.25	22.52	30.62	46.45	40.02	23.96	16.92	20.95	15.18	38.50	82.71	34.97	6.50	99.50	52.01	45.07	37.38			
H2O	27.35	<u>24.95</u>	40.70	51.75	43.57	26.72	<u>20.23</u>	22.06	19.35	39.00	88.30	38.69	6.10	98.00	54.48	46.41	<u>40.48</u>			
TOVA	<u>27.14</u>	23.83	<u>41.82</u>	53.54	43.07	<u>27.42</u>	20.06	16.67	18.11	47.50	89.82	<u>39.12</u>	6.50	98.00	47.41	38.46	39.90			
SnapKV	26.09	24.06	41.22	51.32	45.11	25.71	17.38	<u>21.47</u>	15.81	38.50	84.98	37.36	6.50	99.50	53.16	46.09	39.64			
PyramidKV	24.71	24.08	39.66	50.56	43.50	25.58	17.14	20.50	15.96	39.00	83.54	37.27	6.50	98.50	52.38	45.54	39.03			
CAKE	27.03	26.56	42.96	<u>52.42</u>	<u>45.00</u>	27.73	21.29	18.57	<u>18.26</u>	<u>40.50</u>	<u>88.42</u>	39.31	6.07	99.00	<u>54.38</u>	47.35	40.93			
Llama3.1-8B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	23.01	22.50	31.00	46.79	40.52	24.76	18.38	20.89	16.96	40.50	86.00	38.55	7.00	99.50	57.06	47.16	38.79			
H2O	26.18	24.83	44.08	51.33	43.29	27.78	20.79	22.36	<u>20.69</u>	41.50	<u>89.88</u>	<u>40.94</u>	6.20	99.00	58.12	49.82	41.67			
TOVA	29.44	27.03	45.20	54.15	<u>44.87</u>	28.87	22.21	20.54	20.52	53.50	92.27	40.99	6.17	99.50	51.14	43.15	42.47			
SnapKV	26.40	<u>27.67</u>	<u>47.83</u>	<u>53.40</u>	44.20	<u>28.68</u>	20.56	<u>23.11</u>	20.13	45.50	89.36	39.98	6.33	99.00	57.69	50.69	<u>42.53</u>			
PyramidKV	25.57	27.40	46.36	53.14	44.51	27.16	20.74	22.26	19.90	45.50	87.06	40.13	<u>6.50</u>	99.50	56.98	48.26	41.94			
CAKE	<u>27.41</u>	32.01	49.58	52.99	45.76	28.32	<u>21.85</u>	23.15	21.56	<u>47.50</u>	89.61	40.90	6.33	99.50	<u>57.92</u>	<u>49.83</u>	43.39			
Llama3.1-8B-Instruct, $B_{\text{total}} = 256L$																				
StreamingLLM	23.57	24.73	30.31	46.42	39.19	24.23	20.61	21.10	19.46	45.50	87.50	41.01	6.50	99.50	59.55	48.70	39.87			
H2O	25.93	27.16	43.86	52.23	43.20	28.28	22.10	22.37	22.75	47.50	91.23	<u>41.99</u>	<u>6.28</u>	99.00	61.02	<u>52.02</u>	42.93			
TOVA	<u>29.70</u>	30.85	48.80	54.36	46.71	30.00	<u>23.21</u>	22.18	<u>22.88</u>	57.50	91.83	42.77	6.05	99.50	54.83	46.15	44.21			
SnapKV	27.22	35.25	50.96	54.50	<u>45.73</u>	28.56	22.75	<u>23.80</u>	22.83	55.50	<u>91.28</u>	41.42	6.03	99.50	60.64	52.80	<u>44.92</u>			
PyramidKV	28.73	<u>36.14</u>	<u>51.01</u>	54.57	44.55	<u>29.58</u>	22.68	23.33	22.34	56.00	90.37	41.09	6.25	99.50	60.03	51.45	44.85			
CAKE	29.75	38.74	52.27	<u>54.52</u>	45.39	29.57	24.15	24.24	<u>23.77</u>	<u>57.00</u>	90.97	41.82	6.10	99.50	<u>61.01</u>	51.90	45.67			
Llama3.1-8B-Instruct, $B_{\text{total}} = 512L$																				
StreamingLLM	25.64	27.48	33.30	47.36	40.06	24.80	23.16	20.80	22.85	57.50	87.69	42.08	6.50	97.00	60.51	51.28	41.75			
H2O	27.76	29.01	44.75	52.78	44.31	29.22	24.71	23.11	24.56	54.50	91.38	42.10	<u>6.36</u>	99.00	62.30	<u>54.33</u>	44.39			
TOVA	30.11	35.77	49.88	<u>54.58</u>	45.86	29.73	<u>25.45</u>	22.73	<u>25.00</u>	<u>63.50</u>	<u>91.90</u>	43.33	6.00	99.50	58.38	48.58	45.64			
SnapKV	<u>30.76</u>	42.03	<u>52.13</u>	54.15	46.14	30.51	24.98	24.24	24.65	64.00	92.05	42.04	6.08	99.50	62.62	54.90	<u>46.92</u>			
PyramidKV	30.47	<u>42.15</u>	52.17	54.67	45.25	<u>30.60</u>	25.00	<u>24.33</u>	24.51	62.50	91.24	41.67	5.95	99.50	61.58	53.89	46.59			
CAKE	31.82	42.99	51.65	54.37	46.89	30.73	26.36	24.94	25.27	<u>63.50</u>	91.54	<u>42.52</u>	6.33	99.50	<u>62.31</u>	54.30	47.19			
Llama3.1-8B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	26.64	30.77	35.59	47.31	42.03	24.17	25.81	21.31	25.66	63.50	88.84	42.76	6.50	88.00	61.36	53.47	42.73			
H2O	29.57	36.15	45.94	54.43	44.81	29.04	27.64	23.31	26.47	62.00	91.83	<u>43.14</u>	6.36	99.00	<u>62.74</u>	55.39	46.11			
TOVA	30.66	40.95	51.09	54.58	46.51	30.62	<u>28.12</u>	23.61	26.24	<u>68.00</u>	91.49	43.80	5.92	99.50	60.73	52.64	47.15			
SnapKV	30.95	<u>44.74</u>	<u>52.58</u>	55.09	<u>46.83</u>	30.37	27.87	<u>24.57</u>	25.99	<u>68.00</u>	<u>92.03</u>	42.60	6.50	99.50	63.00	<u>56.50</u>	<u>47.95</u>			
PyramidKV	30.54	43.64	52.73	<u>55.29</u>	46.29	31.28	27.53	24.50	26.00	<u>68.00</u>	92.09	41.75	6.05	99.50	62.35	55.44	47.69			
CAKE	<u>30.88</u>	44.95	52.38	55.49	46.99	<u>30.82</u>	28.68	24.91	<u>26.39</u>	69.00	91.94	42.60	6.00	99.50	62.65	56.89	48.13			
Llama3.1-8B-Instruct, $B_{\text{total}} = 2048L$																				
StreamingLLM	27.40	36.91	37.85	49.23	44.66	24.31	28.57	21.67	27.12	67.50	90.98	42.49	6.12	87.00	63.06	55.39	44.39			
H2O	29.65	39.53	48.64	54.23	46.50	29.28	29.97	23.68	27.21	68.00	<u>91.48</u>	43.06	<u>6.11</u>	99.50	63.06	56.91	47.30			
TOVA	30.30	43.40	52.78	55.12	46.38	30.62	<u>30.80</u>	24.10	27.09	<u>70.50</u>	<u>91.48</u>	43.63	6.00	99.50	62.76	55.32	48.11			
SnapKV	<u>30.99</u>	<u>45.06</u>	53.15	55.25	46.56	<u>30.78</u>	30.24	24.63	<u>27.32</u>	<u>70.50</u>	<u>91.48</u>	42.37	6.00	99.50	63.23	56.66	48.36			
PyramidKV	31.13	<u>45.06</u>	53.80	55.78	<u>46.59</u>	30.89	30.25	24.82	27.35	71.00	91.65	42.62	6.00	99.50	<u>63.27</u>	56.44	<u>48.51</u>			
CAKE	30.79	45.83	<u>53.57</u>	<u>55.56</u>	46.60	30.47	31.12	<u>24.67</u>	27.16	<u>70.50</u>	<u>91.48</u>	<u>43.48</u>	6.00	99.50	63.28	<u>56.84</u>	48.55			

1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Table 7: Performance comparison over 16 datasets of LongBench on Mistral-7B-Instruct for cache budgets from $64L$ to $2048L$. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code		
	NrrQA	Qasper	MF-en	HopQA	2WikiQA	Musique	GovReport	QMSum	MultNews	TREC	TriviaQA	SAMSum	pCount	PR-en	Lcc	RB-P	Avg.				
	Mistral-7B-Instruct-v0.3, $B_{\text{total}} = \text{Full}$																				
Full	28.26	39.33	49.47	45.29	33.21	24.57	32.25	22.16	24.58	75.5	88.36	43.78	0.5	83.0	49.79	51.3	43.21				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 64L$																					
StreamingLLM	17.59	21.81	23.39	34.93	26.73	15.14	13.11	18.29	12.59	37.50	82.27	33.97	0.00	28.00	36.66	36.79	27.42				
H2O	21.04	<u>26.57</u>	34.31	39.34	28.84	17.95	<u>20.11</u>	19.25	<u>16.64</u>	36.50	85.91	36.79	2.00	<u>70.00</u>	41.53	<u>41.28</u>	<u>33.63</u>				
TOVA	21.87	23.48	34.36	41.60	28.35	<u>19.46</u>	17.79	18.41	15.24	50.50	89.16	37.18	0.00	48.50	32.78	33.28	32.00				
SnapKV	20.63	24.03	<u>34.69</u>	<u>41.22</u>	29.94	18.04	15.42	18.82	12.62	36.50	86.44	36.70	<u>1.00</u>	46.50	37.77	40.65	31.31				
PyramidKV	20.33	21.96	29.01	40.13	28.59	18.17	15.34	18.46	12.55	36.50	86.39	36.65	0.00	46.00	38.64	39.30	30.50				
CAKE	<u>21.83</u>	26.85	35.19	39.70	<u>28.89</u>	19.70	20.12	<u>19.19</u>	17.21	<u>38.00</u>	<u>86.49</u>	<u>37.02</u>	0.50	75.00	<u>41.17</u>	42.18	34.31				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 128L$																					
StreamingLLM	16.91	21.51	24.85	34.14	26.99	16.64	15.67	18.61	14.40	43.50	83.26	37.00	0.00	23.00	41.63	42.12	28.76				
H2O	21.25	26.66	35.13	38.82	<u>29.80</u>	18.88	<u>21.00</u>	19.50	<u>18.63</u>	41.00	87.64	<u>38.25</u>	1.50	67.00	44.35	47.29	34.79				
TOVA	22.47	24.26	37.22	42.26	28.85	19.97	19.40	18.70	17.86	63.00	88.98	37.71	0.50	56.50	37.42	37.22	34.52				
SnapKV	21.02	<u>27.26</u>	41.25	45.15	29.23	<u>22.75</u>	20.47	<u>20.17</u>	17.75	42.50	87.28	38.01	0.50	<u>69.50</u>	<u>44.48</u>	45.69	<u>35.81</u>				
PyramidKV	21.73	26.60	<u>41.46</u>	43.20	29.32	21.47	20.23	19.82	17.46	40.00	87.64	37.11	<u>1.05</u>	69.00	42.55	42.98	35.14				
CAKE	<u>22.31</u>	29.15	43.51	<u>44.51</u>	30.36	22.85	21.56	20.47	18.96	<u>47.00</u>	<u>88.60</u>	39.36	1.00	76.50	44.96	<u>46.19</u>	37.33				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 256L$																					
StreamingLLM	18.10	23.41	25.99	34.55	26.78	16.25	17.97	18.82	16.56	52.50	84.26	39.42	1.00	16.50	44.76	44.99	30.12				
H2O	22.36	27.14	37.79	39.21	30.25	19.79	22.35	19.40	20.06	50.50	87.34	39.49	1.50	62.00	47.01	48.11	35.89				
TOVA	22.52	27.29	41.32	44.49	<u>31.08</u>	22.36	21.99	18.96	20.00	68.00	88.58	<u>39.76</u>	0.50	66.00	40.71	41.35	37.18				
SnapKV	<u>23.06</u>	<u>29.43</u>	44.61	<u>45.05</u>	30.39	<u>22.89</u>	22.32	20.77	<u>20.19</u>	56.00	88.21	39.34	0.00	83.00	48.09	48.72	<u>38.88</u>				
PyramidKV	22.66	28.85	<u>45.84</u>	44.30	30.43	21.78	<u>22.48</u>	20.97	19.94	53.00	<u>88.41</u>	38.88	<u>1.05</u>	<u>85.00</u>	46.11	47.78	38.59				
CAKE	24.01	30.78	45.88	45.06	31.45	23.43	22.81	<u>20.87</u>	20.92	<u>58.50</u>	88.29	40.16	0.50	87.00	<u>47.76</u>	<u>48.43</u>	39.74				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 512L$																					
StreamingLLM	19.93	25.83	28.25	36.21	27.60	15.41	21.23	18.67	19.50	64.00	86.45	39.75	1.00	16.50	46.70	47.57	32.16				
H2O	22.34	28.01	38.32	40.12	30.25	18.38	24.37	20.38	<u>21.80</u>	59.50	87.33	<u>41.30</u>	2.00	65.50	48.68	<u>49.76</u>	37.38				
TOVA	25.06	31.69	44.22	44.90	32.48	23.44	24.33	19.81	21.77	71.50	88.46	42.39	0.00	77.50	44.54	45.66	39.86				
SnapKV	24.88	32.92	47.31	<u>45.23</u>	31.78	23.08	<u>24.60</u>	21.80	21.52	65.50	88.61	40.43	0.50	<u>88.50</u>	<u>50.59</u>	50.22	<u>41.09</u>				
PyramidKV	24.15	<u>33.10</u>	46.62	44.16	<u>32.14</u>	23.37	24.22	21.15	21.29	<u>68.00</u>	<u>88.49</u>	40.27	<u>1.55</u>	87.50	48.74	49.17	40.87				
CAKE	25.82	34.24	<u>46.93</u>	45.34	31.71	<u>23.38</u>	25.81	<u>21.47</u>	22.41	68.00	88.46	41.27	0.50	90.00	<u>49.04</u>	49.67	41.50				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 1024L$																					
StreamingLLM	20.96	28.05	30.03	37.06	27.56	16.03	24.03	19.07	22.79	67.00	87.61	40.96	1.50	21.50	48.05	48.87	33.82				
H2O	23.78	31.63	41.31	43.24	31.07	20.43	26.74	20.41	<u>23.93</u>	67.50	88.84	42.62	1.50	72.00	<u>50.60</u>	49.87	39.72				
TOVA	26.97	34.51	45.58	44.32	<u>32.58</u>	22.83	<u>26.91</u>	20.75	23.49	75.00	<u>88.66</u>	43.17	1.00	<u>91.00</u>	47.51	47.06	41.96				
SnapKV	<u>26.63</u>	35.78	48.11	<u>45.75</u>	32.20	23.37	26.71	21.84	23.18	70.50	88.61	41.37	0.50	88.00	<u>50.60</u>	51.79	<u>42.18</u>				
PyramidKV	25.51	<u>36.02</u>	47.72	44.74	33.16	23.91	26.55	<u>21.83</u>	23.27	70.50	88.41	40.94	1.00	87.00	50.17	50.93	41.98				
CAKE	26.09	36.34	48.11	45.97	32.39	<u>23.49</u>	27.56	21.45	24.03	<u>72.50</u>	88.61	<u>42.71</u>	0.00	91.50	51.06	<u>51.25</u>	42.69				
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 2048L$																					
StreamingLLM	22.60	32.48	33.26	40.53	29.57	15.69	26.64	19.70	23.64	69.00	<u>88.41</u>	42.44	2.00	28.50	48.69	50.85	35.88				
H2O	25.76	34.28	44.42	44.41	32.17	20.14	28.79	21.23	24.38	71.00	<u>88.41</u>	43.49	<u>1.50</u>	82.00	<u>50.51</u>	50.39	41.43				
TOVA	27.48	37.20	47.47	45.25	<u>33.14</u>	23.73	<u>29.37</u>	21.73	24.52	75.00	88.66	43.33	1.00	82.00	49.08	50.18	42.45				
SnapKV	<u>26.43</u>	37.60	<u>48.27</u>	45.37	33.08	<u>24.02</u>	28.81	21.43	<u>24.40</u>	73.50	88.36	42.34	1.00	<u>85.50</u>	49.62	<u>51.28</u>	42.56				
PyramidKV	26.00	<u>37.63</u>	48.25	45.82	37.13	24.00	28.90	22.13	24.35	72.50	88.36	42.28	1.00	84.50	50.06	51.07	<u>42.75</u>				
CAKE	25.43	37.84	48.30	<u>45.60</u>	32.21	24.14	30.00	<u>21.83</u>	24.35	<u>74.50</u>	88.36	<u>43.38</u>	0.50	89.00	50.68	51.48	42.98				

1188
 1189 Table 8: Performance comparison over 16 datasets of LongBench on Qwen2.5-7B-Instruct and
 1190 Gemma-7B-Instruct. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code	
	NirvQA	Qasper	MF-en	HopqaQA	2WikiMQA	Musique	GovReport	QMSum	MulitNews	TREC	TriviaQA	SAMSum	PCount	PR-en	Lcc	RB-p	Avg.			
Qwen2.5-7B-Instruct, $B_{\text{total}} = \text{Full}$																				
Full	29.05	43.34	52.52	57.59	47.05	30.24	31.78	23.64	23.96	72.5	89.47	45.61	8.5	100.0	59.61	67.12	48.87			
Qwen2.5-7B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	18.80	23.57	26.71	40.69	37.20	16.24	15.70	18.01	12.61	42.50	81.22	40.10	8.50	25.00	46.77	47.47	31.32			
H2O	22.73	24.52	31.04	43.89	40.19	18.92	18.28	<u>19.82</u>	16.03	42.00	83.11	<u>41.88</u>	8.50	96.50	<u>50.39</u>	<u>55.05</u>	38.30			
TOVA	21.01	23.98	33.65	47.65	38.97	19.42	19.84	19.41	15.54	51.50	85.43	42.49	8.50	93.50	42.60	44.48	38.00			
SnapKV	<u>23.85</u>	<u>27.19</u>	<u>42.43</u>	<u>51.31</u>	<u>42.59</u>	<u>24.38</u>	18.19	19.76	14.54	42.50	84.24	41.10	9.00	97.00	49.45	53.95	40.09			
PyramidKV	19.56	23.39	37.72	48.33	38.33	19.53	15.10	19.32	12.05	42.00	83.29	38.84	9.00	84.50	47.87	50.39	36.83			
CAKE	25.99	33.00	46.22	54.57	45.55	25.59	<u>19.51</u>	20.95	<u>15.82</u>	43.00	84.88	40.99	8.50	96.50	50.62	55.27	41.68			
Qwen2.5-7B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	22.72	29.42	31.47	43.57	38.18	17.99	24.33	19.47	22.46	61.00	87.53	43.79	8.50	34.00	55.17	58.43	37.38			
H2O	26.45	34.94	40.49	48.63	42.02	22.27	25.67	20.90	22.41	59.00	87.83	<u>45.07</u>	8.50	98.48	59.77	63.88	44.14			
TOVA	26.36	37.97	47.88	55.32	<u>45.72</u>	28.38	<u>26.15</u>	21.64	<u>22.70</u>	69.50	88.64	44.23	8.50	100.00	56.66	60.07	46.23			
SnapKV	29.24	<u>41.61</u>	<u>50.93</u>	57.60	45.50	29.39	25.63	<u>23.06</u>	22.26	65.50	<u>88.92</u>	44.65	8.50	100.00	58.16	65.30	47.27			
PyramidKV	<u>29.34</u>	38.60	50.17	55.67	45.12	27.82	23.26	22.16	20.55	62.50	86.85	43.26	8.50	100.00	57.76	61.99	45.85			
CAKE	29.47	42.71	52.12	<u>56.11</u>	46.41	<u>29.13</u>	26.86	23.12	<u>22.72</u>	<u>67.50</u>	<u>89.23</u>	45.46	8.50	100.00	<u>59.11</u>	<u>64.79</u>	47.70			
Gemma-7B-Instruct, $B_{\text{total}} = \text{Full}$																				
Full	15.33	34.27	47.43	28.2	22.69	7.48	26.64	19.67	23.84	69.0	79.41	32.59	1.0	44.5	47.39	45.96	34.09			
Gemma-7B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	10.65	20.76	27.41	24.74	20.06	5.91	13.53	16.98	14.30	47.00	74.51	29.16	4.00	9.50	46.68	45.70	25.68			
H2O	12.98	27.27	35.97	27.70	23.58	<u>7.21</u>	<u>17.70</u>	18.17	<u>18.23</u>	44.00	79.12	<u>31.84</u>	1.63	46.50	47.28	48.62	30.49			
TOVA	<u>13.73</u>	27.38	42.38	<u>28.73</u>	23.34	7.80	17.06	17.76	17.61	57.50	79.65	31.35	1.50	49.00	41.85	43.22	31.24			
SnapKV	12.98	27.53	<u>42.90</u>	29.69	25.45	7.01	17.14	<u>18.29</u>	17.65	50.00	79.07	31.54	0.57	<u>48.50</u>	<u>47.64</u>	<u>49.20</u>	31.57			
PyramidKV	12.51	<u>28.02</u>	42.24	28.18	23.33	6.66	16.82	17.94	17.11	46.50	77.94	30.94	<u>2.00</u>	46.00	45.21	48.35	30.61			
CAKE	14.09	30.33	44.33	28.63	<u>24.47</u>	6.70	18.20	18.83	18.98	<u>53.00</u>	<u>79.54</u>	32.14	1.50	48.00	48.71	50.59	32.38			
Gemma-7B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	13.23	27.04	29.08	25.39	22.15	5.67	20.94	17.74	21.97	67.50	79.28	<u>33.64</u>	<u>1.50</u>	11.50	49.80	49.25	29.73			
H2O	13.39	30.25	41.23	28.75	<u>24.41</u>	7.31	<u>22.85</u>	18.68	22.73	66.50	79.78	33.12	1.57	44.50	<u>49.49</u>	48.55	33.32			
TOVA	13.58	32.92	45.60	<u>28.97</u>	24.50	7.96	22.19	18.95	<u>22.86</u>	73.00	79.94	32.61	<u>1.50</u>	44.50	48.74	47.29	<u>34.07</u>			
SnapKV	13.72	<u>33.76</u>	46.00	29.19	23.85	<u>7.49</u>	22.44	19.35	22.73	<u>71.50</u>	79.96	33.35	<u>1.50</u>	43.50	48.64	47.86	34.05			
PyramidKV	<u>14.32</u>	33.24	<u>46.35</u>	28.75	23.51	6.91	21.66	<u>19.37</u>	22.47	66.00	78.95	33.52	1.00	43.50	47.43	48.07	33.44			
CAKE	15.16	35.12	47.67	28.24	23.12	6.59	23.28	19.72	23.40	68.50	79.47	34.07	<u>1.50</u>	43.00	48.84	<u>49.20</u>	34.18			

F.2 EXTENDED EVALUATION ON ADDITIONAL MODEL ARCHITECTURES

To demonstrate the generalizability of CAKE across different model architectures, we conduct additional experiments on Qwen2.5-7B-Instruct and Gemma-7B-Instruct. The experiments are conducted on two memory scenarios: a low setting ($B_{\text{total}} = 128L$) and a high setting ($B_{\text{total}} = 1024L$). As shown in Table 8, similar to results on Llama and Mistral, CAKE consistently outperforms baseline methods across both low-memory and high-memory scenarios on Qwen and Gemma architectures. CAKE demonstrates significant advantages in low-memory settings by rationally allocating cache sizes and evicting KV pairs with more robust indicators. Under the high-budget setting, CAKE achieves superior performance, even exceeding the full-cache baseline on Gemma (34.18 vs. 34.09). These results further validate CAKE’s adaptability across different model architectures, maintaining its performance advantages regardless of the underlying model design.

1242 Table 9: Performance comparison over 16 datasets of LongBench on different models from 13B to
 1243 70B. The best result is highlighted in **bold**, the second best in underline.

Method	Single-Document QA				Multi-Document QA				Summarization				Few-shot Learning				Synthetic		Code	
	NtvyQA	Qasper	MF-en	HoppotQA	2WikiQA	Musique	GovReport	QMSum	MultiNews	TREC	TriviaQA	SAMSum	PCount	PR-en	Lcc	RB-P	Avg.			
Llama2-13B-Chat, $B_{\text{total}} = \text{Full}$																				
Full	14.39	17.07	27.52	12.66	13.21	5.01	27.52	20.89	26.61	68.5	87.75	42.44	2.27	15.25	48.27	49.87	29.95			
Llama2-13B-Chat, $B_{\text{total}} = 128L$																				
StreamingLLM	10.48	13.44	19.16	12.86	13.81	4.90	15.63	18.91	16.26	40.00	81.66	34.00	2.29	5.00	40.38	35.16	22.75			
H2O	12.74	17.63	<u>25.71</u>	<u>14.12</u>	13.29	3.84	19.35	19.71	21.51	40.00	80.62	<u>39.37</u>	2.03	10.50	40.75	41.31	25.15			
TOVA	11.18	13.55	19.65	13.61	13.20	3.95	17.46	18.17	16.51	58.00	88.77	38.91	2.36	7.50	40.10	37.21	25.01			
SnapKV	13.33	14.73	24.34	13.47	<u>15.10</u>	<u>4.68</u>	20.37	<u>20.01</u>	20.44	42.00	86.71	38.35	3.21	11.50	43.07	<u>42.28</u>	25.85			
PyramidKV	<u>13.04</u>	13.91	25.26	14.37	15.24	<u>4.68</u>	20.01	19.74	20.52	<u>43.50</u>	85.37	38.83	2.80	13.00	<u>43.57</u>	41.31	25.95			
CAKE	12.19	<u>17.17</u>	27.28	13.12	14.73	4.33	<u>20.04</u>	20.14	<u>20.74</u>	42.00	<u>87.75</u>	40.60	<u>2.96</u>	<u>12.50</u>	45.52	43.95	26.56			
Llama2-13B-Chat, $B_{\text{total}} = 1024L$																				
StreamingLLM	12.47	15.03	19.85	11.39	13.68	3.90	23.48	19.33	24.88	63.00	86.37	40.12	4.40	6.50	47.24	45.56	27.32			
H2O	<u>14.76</u>	19.89	29.10	14.84	14.05	4.41	23.38	20.18	25.03	62.50	82.14	41.85	1.13	12.50	47.34	47.07	28.76			
TOVA	13.69	16.10	24.41	12.79	13.43	4.96	23.81	20.22	24.99	69.00	<u>88.15</u>	42.18	<u>3.19</u>	13.25	46.71	<u>48.77</u>	29.10			
SnapKV	14.39	17.87	27.99	12.94	13.51	4.95	<u>24.65</u>	<u>20.74</u>	25.08	68.50	86.47	42.26	2.62	14.00	<u>48.27</u>	49.42	29.60			
PyramidKV	14.15	<u>18.30</u>	28.35	12.98	14.42	<u>5.10</u>	25.05	20.84	25.62	68.50	88.21	<u>42.52</u>	2.28	<u>15.00</u>	47.99	48.12	<u>29.84</u>			
CAKE	15.33	17.84	<u>29.05</u>	13.99	<u>14.30</u>	<u>5.46</u>	24.43	20.60	<u>25.29</u>	68.00	87.42	42.58	3.14	15.50	48.40	48.41	29.98			
Qwen2.5-32B-Instruct, $B_{\text{total}} = \text{Full}$																				
Full	29.25	45.24	51.9	63.31	61.5	38.15	30.33	23.38	23.01	73.5	87.86	45.36	12.0	100.0	51.3	38.17	48.39			
Qwen2.5-32B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	16.26	22.62	25.95	44.14	47.57	23.37	15.66	17.50	13.22	41.50	81.47	38.03	12.00	58.04	43.51	32.46	33.33			
H2O	<u>24.03</u>	24.48	31.10	52.53	52.38	30.84	17.00	<u>19.67</u>	16.27	43.00	83.64	41.41	12.22	94.58	<u>45.85</u>	34.35	38.96			
TOVA	25.53	<u>27.73</u>	38.07	<u>57.08</u>	54.92	<u>32.92</u>	19.28	17.47	15.08	60.50	77.55	<u>41.02</u>	11.81	97.58	38.05	26.75	40.08			
SnapKV	21.28	26.68	<u>40.43</u>	56.74	<u>55.61</u>	33.75	17.15	19.65	14.55	48.50	86.46	40.86	12.50	95.58	45.75	33.55	<u>40.56</u>			
PyramidKV	18.31	24.62	36.85	55.30	55.40	32.38	15.33	18.02	13.15	46.50	<u>84.50</u>	39.62	11.50	88.83	44.05	31.76	38.51			
CAKE	<u>24.03</u>	29.35	41.64	57.42	57.37	31.85	<u>18.54</u>	20.42	<u>16.26</u>	<u>50.00</u>	82.90	40.63	12.50	<u>97.42</u>	46.59	<u>33.86</u>	41.30			
Qwen2.5-32B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	21.42	29.11	31.07	46.44	47.31	26.17	23.76	19.07	21.68	63.50	86.12	42.61	12.50	63.08	49.84	36.13	38.74			
H2O	25.96	32.84	40.35	56.46	55.57	32.91	<u>25.06</u>	21.02	21.94	60.50	86.63	<u>44.85</u>	12.50	97.58	50.96	37.15	43.89			
TOVA	29.38	40.61	48.18	59.65	60.91	34.77	24.86	20.69	<u>22.03</u>	73.00	87.34	45.09	12.00	100.00	49.11	35.90	46.47			
SnapKV	29.93	<u>42.29</u>	<u>50.00</u>	<u>61.97</u>	<u>61.37</u>	37.22	25.01	<u>22.47</u>	21.94	69.50	87.31	44.42	12.50	100.00	51.54	37.32	<u>47.17</u>			
PyramidKV	29.39	41.56	49.31	63.12	61.52	38.25	22.47	21.87	18.36	70.50	87.78	44.01	11.83	100.00	50.27	<u>37.38</u>	46.73			
CAKE	30.77	44.49	50.73	61.59	60.60	<u>37.39</u>	25.86	<u>22.77</u>	22.18	<u>72.50</u>	87.78	43.84	12.12	100.00	<u>51.13</u>	37.71	47.59			
Llama3-70B-Instruct, $B_{\text{total}} = \text{Full}$																				
Full	26.08	47.26	49.73	49.61	54.84	27.51	31.01	22.48	27.99	73.5	92.88	45.21	11.0	68.5	37.72	67.28	45.79			
Llama3-70B-Instruct, $B_{\text{total}} = 128L$																				
StreamingLLM	22.33	28.13	28.07	42.85	50.73	22.18	17.75	20.02	18.23	44.00	82.95	40.05	11.50	67.50	44.44	59.76	37.53			
H2O	24.92	29.15	36.50	45.98	50.40	24.52	<u>19.38</u>	20.61	<u>21.40</u>	43.50	91.78	<u>42.90</u>	<u>10.55</u>	70.15	46.92	63.45	40.13			
TOVA	20.61	25.79	30.59	43.67	47.05	24.33	6.18	17.45	18.22	58.50	91.52	37.04	7.46	34.62	41.71	51.35	34.76			
SnapKV	<u>25.29</u>	33.37	<u>44.25</u>	48.01	<u>52.93</u>	<u>27.01</u>	18.90	<u>21.53</u>	20.80	48.50	<u>91.63</u>	41.53	7.50	68.00	45.87	64.16	<u>41.20</u>			
PyramidKV	25.87	<u>34.16</u>	41.26	48.55	52.24	24.19	19.17	21.20	21.02	46.50	91.43	41.90	9.00	<u>68.00</u>	<u>46.52</u>	63.02	40.88			
CAKE	25.21	<u>39.26</u>	44.39	47.47	54.48	27.22	21.26	<u>22.37</u>	22.25	<u>57.00</u>	91.28	43.69	10.00	67.50	44.41	64.16	42.62			
Llama3-70B-Instruct, $B_{\text{total}} = 1024L$																				
StreamingLLM	25.18	35.93	32.89	45.57	49.49	21.63	24.42	20.28	26.03	66.00	91.32	43.53	11.00	68.00	<u>42.61</u>	65.87	41.86			
H2O	26.66	41.54	44.56	47.47	53.84	26.26	26.05	21.52	26.75	67.00	92.88	44.90	11.00	68.00	40.43	69.09	44.25			
TOVA	27.03	45.73	47.61	49.77	<u>55.54</u>	27.02	26.04	21.59	26.58	<u>73.00</u>	92.95	46.24	11.00	68.00	39.76	62.59	45.03			
SnapKV	26.63	<u>45.96</u>	47.66	48.57	54.92	27.60	<u>26.31</u>	22.34	<u>26.75</u>	72.00	<u>92.88</u>	45.10	10.00	68.50	40.72	67.22	<u>45.20</u>			
PyramidKV	<u>26.99</u>	45.44	<u>48.01</u>	49.04	54.17	27.37	26.08	22.45	26.22	<u>73.00</u>	91.95	45.10	10.50	68.50	38.44	67.95	45.08			
CAKE	26.95	46.94	48.51	47.91	55.98	<u>27.46</u>	26.39	<u>22.44</u>	27.44	73.50	92.44	<u>45.60</u>	11.00	68.00	44.23	<u>68.47</u>	45.83			

F.3 EXTENDED EVALUATION ON LARGER-SCALE MODELS

We further extend our experiments to models ranging from 13B to 70B parameters, including Llama2-13B-Chat, Qwen2.5-32B-Instruct, and Llama3-70B-Instruct, under two settings: $B_{\text{total}} = 128L$ and $B_{\text{total}} = 1024L$. In Table 9, CAKE outperforms baseline methods across all model sizes. Notably, under the high memory setting ($B_{\text{total}} = 1024L$), CAKE achieves even better performance than full-cache settings for both Llama2-13B (29.98 vs. 29.95) and Llama3-70B (45.83 vs. 45.79), showing that our method scales well to larger models while maintaining its efficiency advantages.

1296 **G EXTENDED EXPERIMENTAL RESULTS ON NEEDLEBENCH**
 1297
 1298

1299 In this section, we provide more detailed experimental results on NeedleBench (Li et al., 2024a),
 1300 using Mistral-7B-Instruct-v0.3 (Jiang et al., 2023) and Llama3.1-8B-Instruct (Dubey et al., 2024)
 1301 as backbone LLMs. We examine three distinct tasks on NeedleBench: (1) *Single-Needle Retrieval*
 1302 (*S-RT*): Assessing precision in locating individual details within extensive texts. (2) *Multi-Needle*
 1303 *Retrieval (M-RT)*: Evaluating the retrieval of multiple related information pieces dispersed through-
 1304 out lengthy texts. (3) *Multi-Needle Reasoning (M-RS)*: Measuring complex reasoning abilities by
 1305 extracting multiple information elements from extended contexts to answer questions.

1306 All experiments are conducted rigorously following the original NeedleBench protocol: Levenshtein
 1307 distance is used to measure the similarity between predictions and references. Each case is repeated
 1308 ten times to ensure stable scores, and the results are weighted-averaged to obtain an overall score,
 1309 providing a balanced representation of each task.

1310 **Results on Mistral-7B-Instruct-v0.3.** The results presented in Table 10 and Table 11 demonstrate
 1311 that CAKE consistently outperforms other KV cache eviction methods across all evaluated tasks.
 1312 CAKE exhibits the smallest decrease in performance scores when compared with the full cache
 1313 baseline. The robustness of CAKE is evident across all tested cache sizes, with particularly no-
 1314 table improvements in Multi-Needle Retrieval tasks (Figure 10 and Figure 11). This superior per-
 1315 formance indicates CAKE’s enhanced ability to tolerate attention-shifts in long contexts, which can
 1316 be attributed to its eviction strategy that incorporates both sustained importance and attention vari-
 1317 ability. The proposed indicator and strategy used in CAKE proves to be effective in maintaining
 1318 performance even under constrained cache conditions.

1319 **Results on Llama3.1-8B-Instruct.** We further evaluate CAKE’s performance on Llama3.1-8B-
 1320 Instruct using NeedleBench 32K. The results presented in Table 12 illustrate CAKE’s comprehensive
 1321 advantages across all three tasks, and Figure 12 showcases CAKE’s outstanding performance on
 1322 multibench tasks, confirming its effectiveness on models with GQA architecture as well.

1323 On the whole, CAKE achieves better overall performance across the two models and various cache
 1324 sizes, especially for Multi-Needle Retrieval tasks. Yet, given extremely small cache size budgets,
 1325 Multi-Needle Retrieval tasks performance experiences a sharp decline in performance. This may be
 1326 attributed to the stringent cache limitations are inadequate to support the information demanded by
 1327 multiple needles. Nevertheless, compared to alternative approaches, CAKE demonstrates superior
 1328 capability in mitigating this issue. Future research should focus on exploring more sophisticated
 1329 methods for retaining KV pairs, as this remains a promising avenue for further improvement.

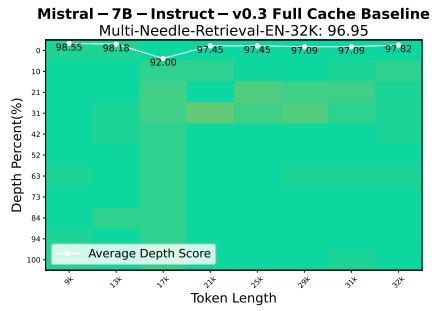
1330
 1331
 1332 Table 10: Performance comparison over three subtasks of NeedleBench 32K benchmark on Mistral-
 1333 7B-Instruct-v0.3. The best is highlighted in **bold**, the second best is in underline.

Method	Single-Retrieval			Multi-Retrieval			Multi-Reasoning			Overall
	ZH	EN	Overall	ZH	EN	Overall	ZH	EN	Overall	
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = \text{Full}$										
Full	91.72	83.01	87.37	90.36	96.95	93.66	55.64	45.89	50.77	78.28
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 1024L$										
H2O	59.20	39.75	49.48	<u>39.50</u>	32.91	36.20	36.95	46.58	41.76	43.18
SnapKV	<u>91.73</u>	<u>90.25</u>	<u>90.99</u>	<u>32.27</u>	<u>81.59</u>	<u>56.93</u>	<u>48.87</u>	44.34	<u>46.61</u>	<u>67.46</u>
PyramidKV	91.94	90.69	91.32	25.23	70.55	47.89	47.68	45.32	46.50	64.84
CAKE	89.43	89.89	89.66	51.68	90.32	71.00	52.50	45.57	49.04	71.87
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 512L$										
H2O	60.87	38.44	49.66	4.50	12.68	8.59	36.21	<u>46.15</u>	41.18	34.79
SnapKV	91.21	<u>90.49</u>	<u>90.85</u>	<u>4.95</u>	<u>33.32</u>	<u>19.14</u>	<u>43.88</u>	45.20	<u>44.54</u>	<u>55.44</u>
PyramidKV	90.06	88.42	89.24	3.50	17.91	10.70	42.43	45.77	44.10	52.14
CAKE	89.97	92.06	91.01	26.18	70.91	48.55	47.60	46.28	46.94	65.05
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 256L$										
H2O	66.54	40.69	53.61	0.23	<u>4.64</u>	<u>2.43</u>	34.61	47.74	41.17	34.53
SnapKV	<u>89.27</u>	<u>89.10</u>	<u>89.18</u>	<u>0.59</u>	2.95	1.77	<u>35.67</u>	47.10	<u>41.39</u>	<u>48.62</u>
PyramidKV	81.51	84.87	83.19	0.18	4.27	2.23	33.41	46.90	40.16	45.99
CAKE	90.33	89.83	90.08	<u>1.27</u>	<u>7.50</u>	<u>4.39</u>	<u>39.03</u>	<u>47.21</u>	<u>43.12</u>	<u>50.29</u>

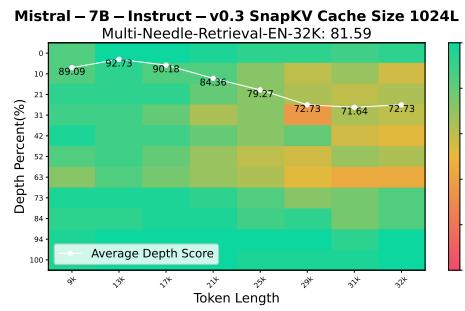
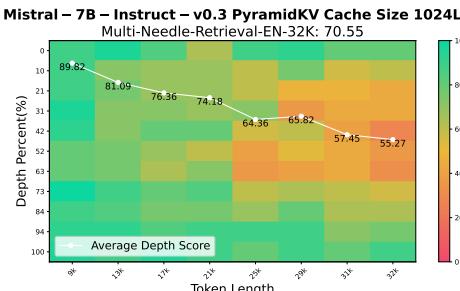
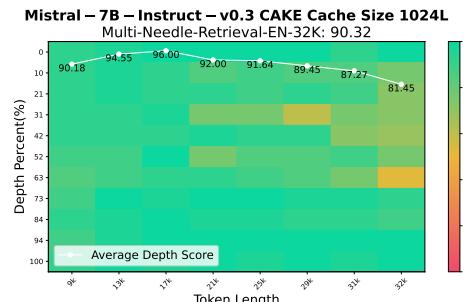
1350 Table 11: Performance comparison over three subtasks of NeedleBench 8K benchmark on Mistral-
 1351 7B-Instruct-v0.3. The best is highlighted in **bold**, the second best is in underlined.

1352

Method	Single-Retrieval			Multi-Retrieval			Multi-Reasoning			Overall
	ZH	EN	Overall	ZH	EN	Overall	ZH	EN	Overall	
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = \text{Full}$										
Full	98.33	84.96	91.65	98.25	96.85	97.55	63.59	57.14	60.37	84.03
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 1024L$										
H2O	77.93	54.86	66.39	5.15	42.20	23.67	43.55	57.39	50.47	48.80
SnapKV	97.66	<u>89.92</u>	<u>93.79</u>	<u>12.25</u>	50.65	31.45	<u>57.38</u>	<u>57.52</u>	<u>57.45</u>	64.19
PyramidKV	98.68	87.55	93.11	9.95	37.65	23.80	55.59	56.77	56.18	61.24
CAKE	98.20	90.27	94.24	53.80	87.25	70.53	57.75	57.57	57.66	76.15
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 512L$										
H2O	72.73	50.71	61.72	0.85	2.75	1.80	39.79	57.86	48.82	39.87
SnapKV	98.11	90.58	<u>94.35</u>	<u>2.25</u>	<u>3.85</u>	<u>3.05</u>	<u>52.71</u>	57.78	<u>55.25</u>	<u>55.23</u>
PyramidKV	97.45	87.87	92.66	1.25	3.00	2.12	49.23	58.52	53.87	53.86
CAKE	98.08	92.33	95.20	16.60	45.25	30.93	52.78	<u>58.48</u>	55.63	64.05
Mistral-7B-Instruct-v0.3, $B_{\text{total}} = 256L$										
H2O	74.35	48.40	61.37	0.20	0.20	0.20	37.38	57.87	47.62	38.90
SnapKV	97.47	91.89	94.68	<u>0.60</u>	0.50	<u>0.55</u>	<u>42.28</u>	<u>57.93</u>	<u>50.10</u>	<u>53.07</u>
PyramidKV	92.02	88.24	90.13	0.30	<u>0.55</u>	0.43	39.88	57.02	48.45	50.71
CAKE	97.22	91.93	<u>94.57</u>	1.00	1.45	1.22	44.21	58.73	51.47	53.64



(a) Mistral-7B-Instruct-v0.3 full cache baseline

(b) SnapKV: cache budget $B_{\text{total}} = 1024L$ (c) PyramidKV: cache budget $B_{\text{total}} = 1024L$ (d) CAKE: cache budget $B_{\text{total}} = 1024L$

1399 Figure 10: Performance comparison of Mistral-7B-Instruct-v0.3 on NeedleBench 32K Multi-Needle
 1400 Retrieval Task (EN) using different KV eviction methods under $B_{\text{total}} = 1024L$.
 1401

1402

1403

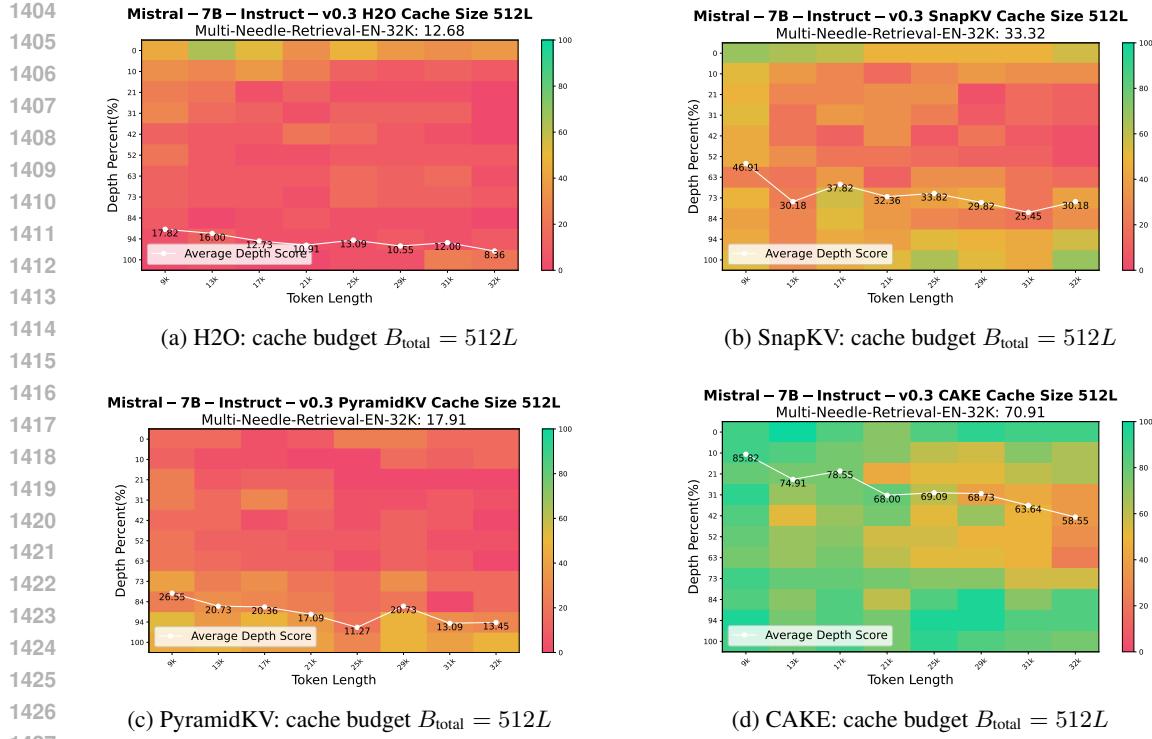
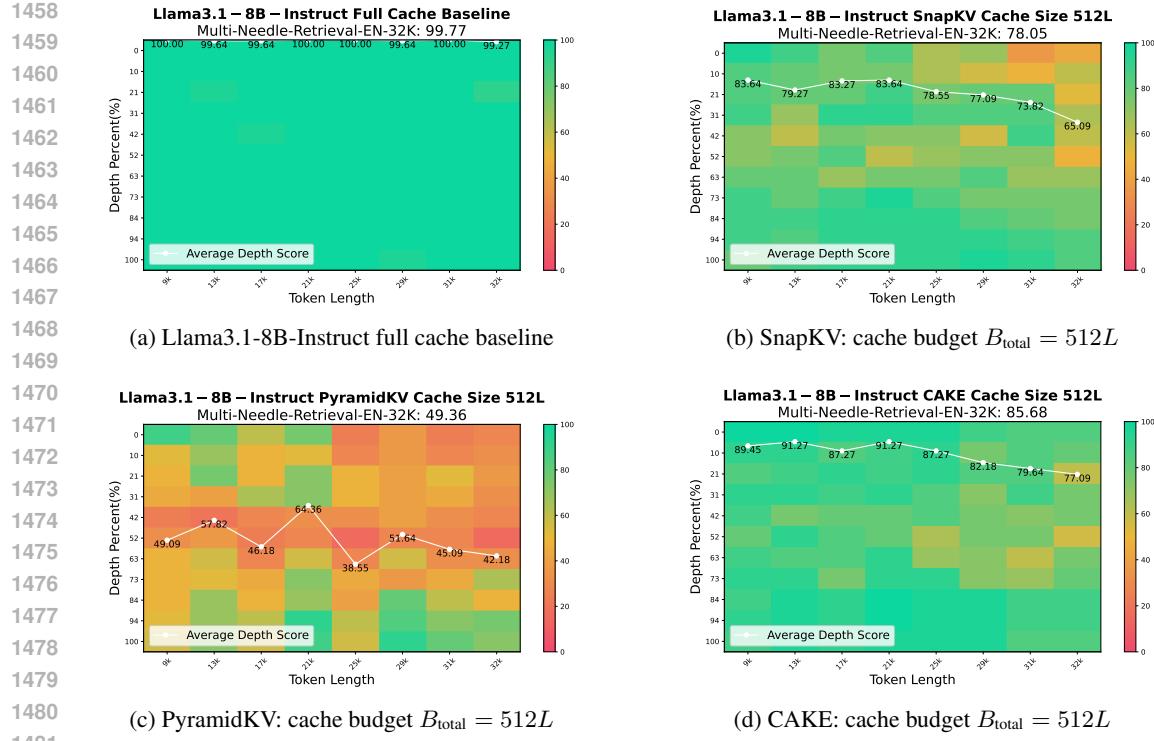


Figure 11: Performance comparison of Mistral-7B-Instruct-v0.3 on NeedleBench 32K Multi-Needle Retrieval Task (EN) using different KV eviction methods under $B_{\text{total}} = 512L$.

Table 12: Performance comparison over three subtasks of NeedleBench 32K benchmark on Llama3.1-8B Instruct. The best is highlighted in **bold**, the second best is in underline.

Method	Single-Retrieval			Multi-Retrieval			Multi-Reasoning			Overall
	ZH	EN	Overall	ZH	EN	Overall	ZH	EN	Overall	
Llama3.1-8B Instruct, $B_{\text{total}} = \text{Full}$										
Full	100.00	82.99	91.49	99.64	99.77	99.70	72.50	77.70	75.10	89.04
Llama3.1-8B Instruct, $B_{\text{total}} = 1024L$										
H2O	89.51	66.59	78.05	89.45	91.77	90.61	54.52	71.91	63.21	77.37
SnapKV	99.90	<u>84.21</u>	<u>92.06</u>	86.68	98.23	<u>92.45</u>	<u>71.11</u>	73.29	<u>72.20</u>	<u>86.22</u>
PyramidKV	100.00	85.53	92.76	56.77	88.55	72.66	71.59	76.16	73.87	81.07
CAKE	100.00	84.02	92.01	<u>89.41</u>	98.91	94.16	69.29	<u>73.42</u>	71.35	86.46
Llama3.1-8B Instruct, $B_{\text{total}} = 512L$										
H2O	88.14	66.08	77.11	30.05	<u>80.14</u>	55.09	50.97	73.20	62.08	66.00
SnapKV	99.90	<u>84.30</u>	<u>92.10</u>	<u>36.09</u>	78.05	<u>57.07</u>	69.20	73.41	<u>71.31</u>	<u>75.35</u>
PyramidKV	100.00	86.17	93.09	11.00	49.36	30.18	<u>68.72</u>	77.13	72.93	68.17
CAKE	100.00	82.07	91.04	47.36	85.68	66.52	67.67	71.15	69.41	77.19
Llama3.1-8B Instruct, $B_{\text{total}} = 256L$										
H2O	88.54	69.61	79.08	<u>4.95</u>	6.23	5.59	50.11	72.86	61.49	51.75
SnapKV	<u>99.90</u>	<u>82.90</u>	<u>91.40</u>	4.05	<u>8.45</u>	<u>6.25</u>	66.68	<u>74.48</u>	70.58	<u>59.61</u>
PyramidKV	100.00	85.07	92.53	1.91	6.09	4.00	64.55	76.14	<u>70.34</u>	59.32
CAKE	99.81	81.92	90.86	6.36	27.64	17.00	<u>64.57</u>	68.61	66.59	61.42



1482 Figure 12: Performance comparison of Llama3.1-8B-Instruct on NeedleBench 32K Multi-Needle
1483 Retrieval Task (EN) using different KV eviction methods under $B_{\text{total}} = 512L$.

1489 Table 13: Prompt prefilling and token decoding latency comparison.

Prefill+Decoding Length	Method	Overall Generation Time (s)	Prompt Prefill Time (s)	Decoding Time (s)	Decoding Time per Token (ms)
7168+1024	SnapKV	33.74	0.71	33.03	32.25
	PyramidKV	32.29	0.78	31.50	30.77
	CAKE	31.77	0.81	30.96	30.24
	Full Cache	34.54	0.70	33.83	33.04
15360+1024	SnapKV	34.78	1.66	33.12	32.34
	PyramidKV	33.02	1.73	31.29	30.56
	CAKE	32.96	1.77	31.19	30.46
	Full Cache	50.71	1.66	49.06	47.91
31744+1024	SnapKV	37.27	4.13	33.15	32.37
	PyramidKV	35.78	4.20	31.58	30.84
	CAKE	36.42	4.29	32.12	31.37
	Full Cache	86.42	4.14	82.28	80.35
64512+1024	SnapKV	44.29	11.33	32.96	32.19
	PyramidKV	43.11	11.48	31.63	30.89
	CAKE	43.10	11.47	31.63	30.89
	Full Cache	176.95	11.42	165.53	161.65
130048+1024	SnapKV	68.98	35.38	33.60	32.81
	PyramidKV	67.61	35.51	32.10	31.35
	CAKE	66.83	35.55	31.28	30.54
	Full Cache	364.15	35.86	328.28	320.59

1512 H ADDITIONAL EXPERIMENTS AND ANALYSIS ON EFFICIENCY 1513

1514 In this section, we present a detailed time breakdown during both prompt prefilling and token decoding
 1515 to better assess CAKE’s effectiveness across different inference stages compared to other KV
 1516 eviction methods, including SnapKV and PyramidKV with various allocation strategies. All meth-
 1517 ods are implemented as extensions of FlashAttention-2 (Dao, 2023) using Mistral-7B-Instruct-v0.3,
 1518 and their performance was benchmarked against a full cache baseline utilizing FlashAttention. To
 1519 measure computational time more precisely throughout the entire text generation phase, we have
 1520 separately evaluated the prompt prefilling time and the duration of autoregressive decoding. By fix-
 1521 ing the generation token count at 1024, the experiment simulates typical long-context processing
 1522 scenarios in LLMs, which are characterized by lengthy inputs and concise outputs.

1523 As shown in Table 13, while dramatically lowering latency in the decoding phase, CAKE, SnapKV,
 1524 and PyramidKV achieve comparable prompt prefilling times. For SnapKV and PyramidKV with
 1525 fixed allocation strategies, the prompt prefilling stage requires waiting for KV caching and per-
 1526 forming one eviction operation per layer, resulting in L eviction operations. However, compared to
 1527 the computation-intensive prefill stage, the eviction operation is relatively inexpensive and can be
 1528 considered negligible. Notably, while CAKE utilizes preference-guided cascading cache manage-
 1529 ment, it achieves similar timing to SnapKV and PyramidKV with fixed allocation strategies. This
 1530 is attributed to the parallelizable execution of KV cache eviction between layers (Algorithm 1, lines
 1531 9–14), which allows the dynamic update process to be incorporated into a single eviction operation.
 1532 Therefore, the total execution time ultimately equates to that of performing L eviction operations.

1533 1534 I ADDITIONAL ABLATION STUDIES 1535

1536 I.1 ABLATION STUDY ON METHOD COMPONENTS

1537 In this section, we conduct ablation studies on LongBench to analyze the contribution of each com-
 1538 ponent in our method. These studies highlight the importance of our design choices, focusing on
 1539 two key aspects: (1) the preference-prioritized adaptive allocation strategy, which combines atten-
 1540 tion dispersion (\mathcal{H}) for spatial influence and attention shift (\mathcal{V}) for temporal dynamics, and (2) the
 1541 attention-shift tolerant eviction indicator, which integrates mean attention score (Mean) for identi-
 1542 fying sustained important tokens and attention variance (Var) to capture dynamic attention changes.

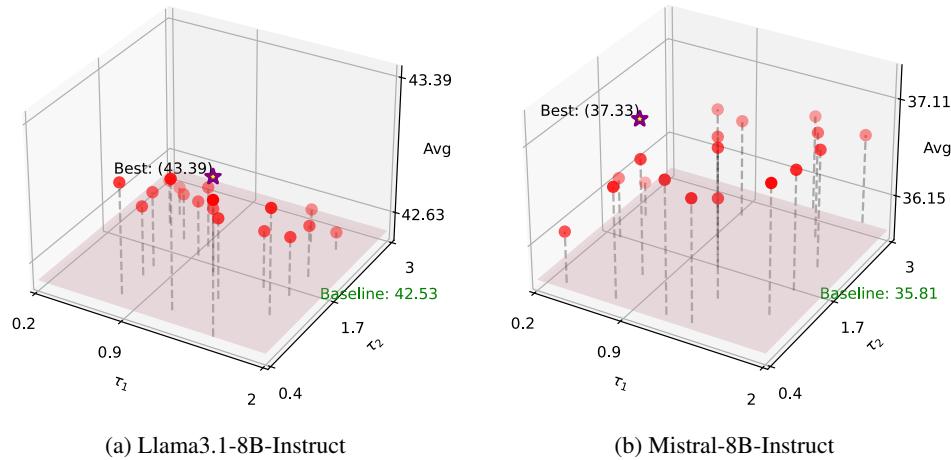
1543 As shown in Table 14, incorporating each additional component leads to a noticeable performance
 1544 improvement, demonstrating the synergistic effect of our method. The preference-prioritized adap-
 1545 tive allocation strategy improves performance from 28.14 to 28.82 when both \mathcal{H} and \mathcal{V} are included,
 1546 underscoring the importance of capturing spatial and temporal attention dynamics. This highlights
 1547 how our approach goes beyond traditional uniform strategies by tailoring cache allocation based
 1548 on complex layer-specific attention patterns. Furthermore, incorporating attention variance (Var) in
 1549 the eviction indicator significantly enhances performance by accounting for attention fluctuations,
 1550 further improving average performance to 29.29. This confirms that addressing both spatial distri-
 1551 bution and temporal evolution of attention is crucial for effective KV cache management, especially
 1552 in long-context scenarios where dynamic attention patterns heavily influence performance. Our
 1553 comprehensive approach successfully balances these factors, yielding superior results compared to
 1554 methods that overlook such nuanced attention characteristics.

1555 Table 14: Ablation study on different components of our method. The experiments are conducted
 1556 on Llama2-7b-Chat with total cache budget $B_{\text{total}} = 128L$.
 1557

Allocation strategy		Eviction indicator		Avg.
\mathcal{H}	\mathcal{V}	Mean	Var	
		✓		28.07
✓		✓		28.14
✓	✓	✓		28.82
✓	✓	✓	✓	29.29

1566 **I.2 INFLUENCE OF TEMPERATURE PARAMETER**
 1567

1568 In this section, we examine the impact of temperature parameters τ_1 and τ_2 , which modulate the
 1569 influence of attention dispersion and shift in our preference-prioritized adaptive allocation strategy.
 1570 We evaluated various settings on LLama3.1-8B-Instruct and Mistral-7B-Instruct-v0.3 models, using
 1571 a cache budget of $B_{\text{total}} = 128L$. Our approach is compared against SnapKV, a state-of-the-art
 1572 KV eviction method employing uniform cache allocation, which serves as our baseline. As shown
 1573 in Figure 13, our method consistently outperforms the baseline across different configurations for
 1574 both models. Mistral-7B-Instruct-v0.3 achieves a peak performance of 37.33 (baseline: 35.81),
 1575 while LLama3.1-8B-Instruct reached 43.39 (baseline: 42.53), demonstrating the robustness and ef-
 1576 fectiveness of our strategy. While our approach is inherently robust, we found that simple, low-cost
 1577 adjustments to τ_1 and τ_2 can further enhance performance. These adjustments allow for better cap-
 1578 ture of model-specific attention patterns, optimizing cache allocation without expensive retraining
 1579 or extensive hyperparameter searches. This characteristic is particularly valuable in industrial de-
 1580 ployments where specific models need to be optimized under constrained memory budgets, offering
 1581 a practical solution for resource-efficient inference in production environments.



1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598 Figure 13: Performance comparison of different τ_1 and τ_2 configurations for Mistral-7B-Instruc-
 1599 t-v0.3 and LLama3.1-8B-Instruct models. The plots show the impact of varying τ_1 and τ_2 on model
 1600 performance, with the baseline (SnapKV) and best scores indicated.

1601
 1602 **J ADDITIONAL DETAILED ANALYSIS AND VISUALIZATION OF ATTENTION
 1603 DYNAMICS**

1606 In this section, we provide a more intuitive and comprehensive analysis of the differences in attention
 1607 patterns by presenting detailed visualizations of attention weights across various dimensions. To
 1608 illustrate the variations between different layers, we visualize the model’s attention weights across all
 1609 layers by aggregating the attention weights of each attention head. We also highlight the differences
 1610 in attention between different models by conducting experiments on both Llama3 and Mistral. To
 1611 demonstrate the attention differences across various input contexts, we visualize attention weights
 1612 for different task types on Llama3, including single-document QA, summarization, and code tasks.
 1613 Additionally, we visualize attention weights for the same task types but with varying contexts on
 1614 Mistral.

1615 **Attention Differences across Layers.** As shown in Figure 14 and Figure 15, attention patterns ex-
 1616 hibit significant differences across layers, as mentioned in the main text. Regardless of the model
 1617 or input context, we can observe varying degrees of attention dispersion and attention shift across
 1618 different layers. For example, we see layers with higher attention dispersion (e.g., layer 0 in Fig-
 1619 ure 14(b)) or lower dispersion (e.g., layer 1 in Figure 14(b)), as well as layers with higher shift (e.g.,
 layer 24 in Figure 15 (c)) or lower shift (e.g., layer 35 in Figure 15(c)). These observations demon-

1620
1621 strate the importance of carefully considering the differences in attention mechanisms and allocating
1622 appropriate cache resources accordingly.
1623

Attention Differences across Models. As shown in Figure 14 and Figure 15, these models have
1624 notable differences in attention patterns. For instance, Mistral often exhibits higher attention dis-
1625 persion at the beginning and end of its layer stack. In contrast, Llama3 frequently displays high
1626 attention dispersion primarily in the early layers. These distinct patterns underscore the importance
1627 of model-specific considerations when designing caching strategies, as different architectures may
1628 have varying attention characteristics across their layer stacks.
1629

Attention Differences across Contexts. As evident in Figure 14(a)-(c), there are striking differ-
1630 ences in attention patterns across different task types. Even within the same task type, while there
1631 may be some similarities in attention dispersion, attention shift patterns show marked differences.
1632 This can lead to different layer preferences for KV cache (comparing layers 10-16 in Figure 15(a)
1633 with layers 10-16 in Figure 15(b)). Therefore, whether the task types are the same or different,
1634 attention patterns are consistently dynamic and varied.
1635

Given the dynamic nature of attention across layers, models, and even contexts, adopting a uniform
1636 allocation strategy, while safe, does not effectively utilize memory. On the other hand, using a fixed
1637 pattern allocation strategy fails to generalize effectively across multiple models and task scenarios.
1638 In contrast, our CAKE thoroughly considers layer preferences for KV cache, formulating a reason-
1639 able cache size allocation scheme from a global perspective. It takes into account the substantial
1640 impact of attention dynamics, allowing it to dynamically adapt to different models and contexts.
1641 This approach ensures that CAKE can effectively respond to the varied and changing attention pat-
1642 terns observed across different scenarios.
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

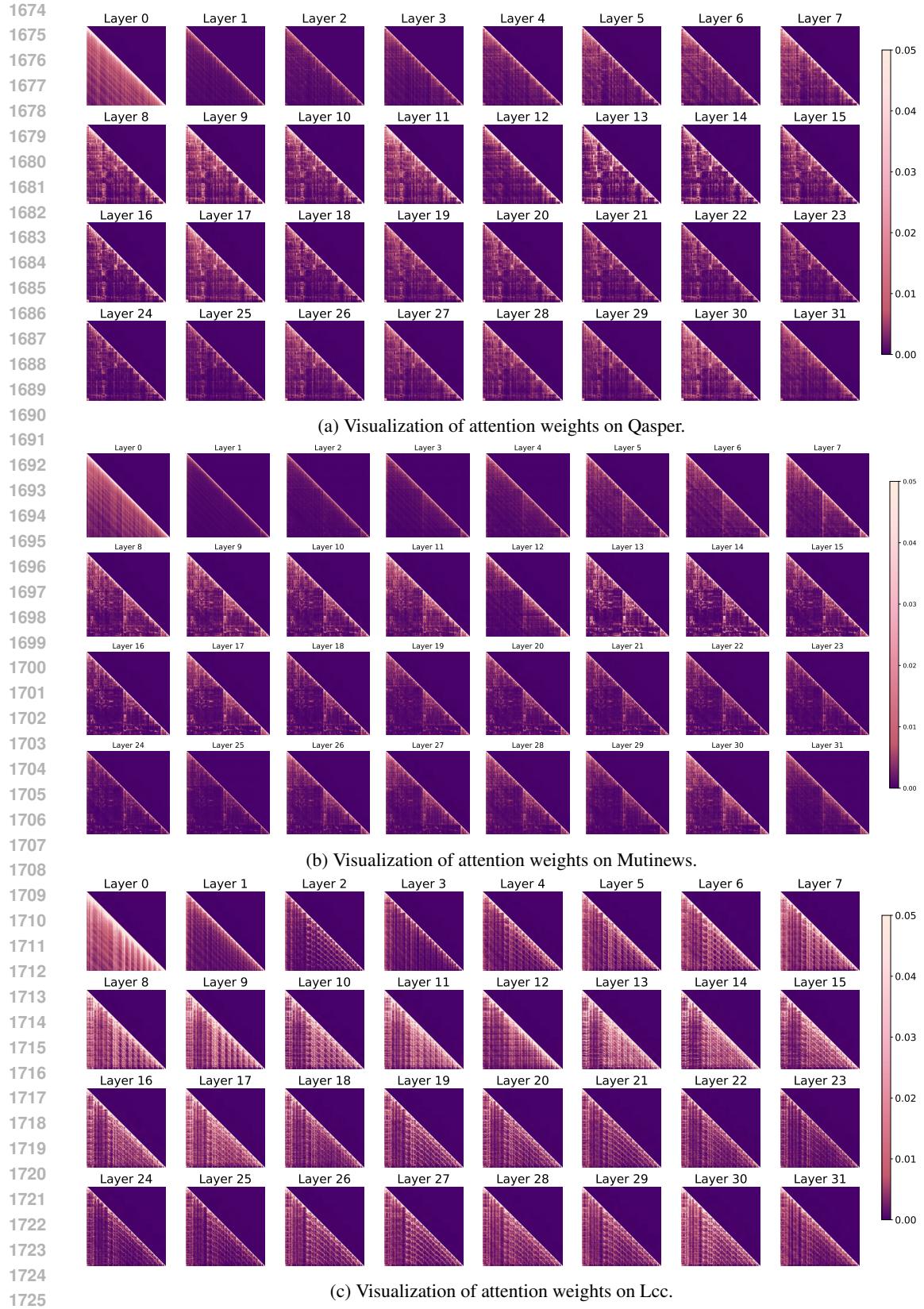
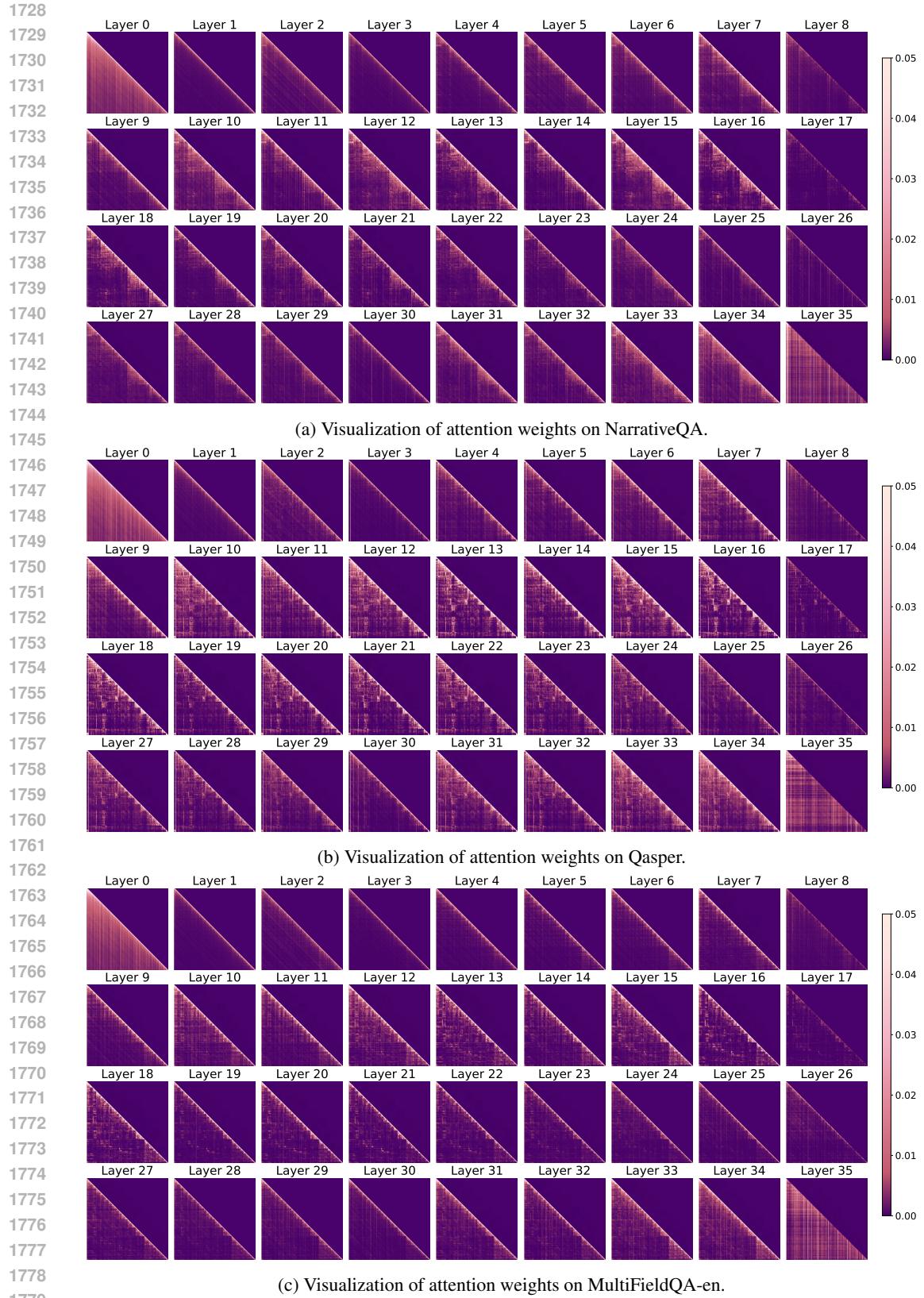


Figure 14: Visualization results of attention weights for Llama3.1-8B-Instruct across datasets from different tasks.



1780 Figure 15: Visualization results of attention weights for Mistral-7B-Instruct-v0.3 across datasets
1781 from Single-document QA task.