



SCOPE: Optimizing Key-Value Cache Compression in Long-context Generation

Jialong Wu^{*♡}, Zhenglin Wang^{*♡}, Linhai Zhang[◇], Yilong Lai[♡],
Yulan He^{◇♣}, Deyu Zhou^{♡†},

[♡] School of Computer Science and Engineering, Key Laboratory of Computer Network and Information Integration, Ministry of Education, Southeast University, China

[◇] Department of Informatics, King's College London, UK

[♣] The Alan Turing Institute, UK

{jialongwu, zhenglin, d.zhou}@seu.edu.cn

Abstract

Key-Value (KV) cache has become a bottleneck of LLMs for long-context generation. Despite the numerous efforts in this area, the optimization for the decoding phase is generally ignored. However, we believe such optimization is crucial, especially for long-output generation tasks based on the following two observations: (i) *Excessive compression* during the prefill phase which requires specific full context, impairs the comprehension of the reasoning task; (ii) *Deviation of heavy hitters*¹ occurs in the reasoning tasks with long outputs. Therefore, **SCOPE**, a simple yet efficient framework that separately performs KV cache optimization during the prefill and decoding phases, is introduced. Specifically, the KV cache during the prefill phase is preserved to maintain the essential information, while a novel strategy based on sliding is proposed to select essential heavy hitters for the decoding phase. Memory usage and memory transfer are further optimized using adaptive and discontinuous strategies. Extensive experiments on LONGGENBENCH show the effectiveness and generalization of SCOPE and its compatibility as a plug-in to other prefill-only KV compression methods.²

1 Introduction

Large Language Models (LLMs) (Dubey et al., 2024; Jiang et al., 2023; Yang et al., 2024a; Team et al., 2024; Achiam et al., 2023; Anthropic, 2024) have demonstrated powerful abilities for processing long-context tasks. When LLMs infer on these long-context tasks, the Key-Value (KV) cache occupies a larger amount of GPU memory and becomes a substantial bottleneck (Waddington et al., 2013;

* Equal Contribution.

† Corresponding Author.

¹According to Zhang et al. (2023), “heavy hitters” refer to the KV cache of pivotal tokens, a small subset of the entire KV cache, that effectively captures the critical information.

²The code is available in <https://github.com/Linking-ai/SCOPE>

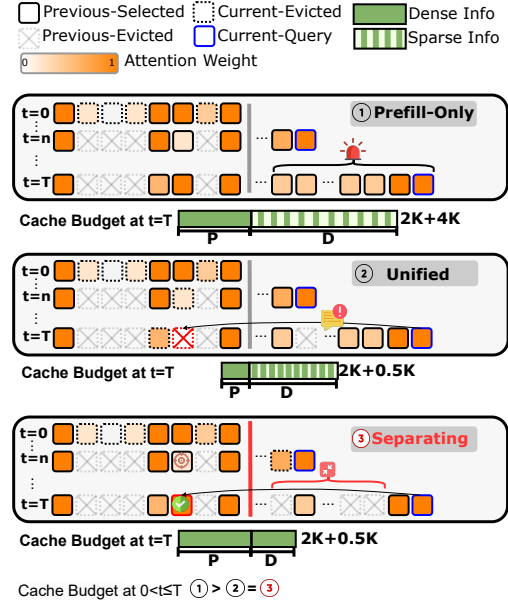


Figure 1: Illustration of three paradigms for compression during the decoding phase on a task with 4K input and 4K output. **Separating** the prefill and decoding phases facilitates the preservation of the essential information KV cache from the prefill phase while allowing for efficient allocation of the KV cache generated during the decoding phase.

Luohe et al., 2024; Yuan et al., 2024; Fu, 2024). For example, an RTX 3090 server struggles to handle the KV cache for a 64K context in LLaMA3.1-8B, which has a 128K context window. Therefore, compressing the KV cache while maintaining the performance is crucial.

LLM inference process involves the *prefill* phase and the *decoding* phase. For tasks with long inputs and short output (Kamradt, 2023; Bai et al., 2024) (e.g., long-form QA or sentence retrieve), effective compression of the KV cache during the prefill phase is crucial. However, for tasks with both long inputs and long outputs (Liu et al., 2024b,c) (e.g., lengthy text summarization and multi-question answering), KV cache compression holds equal im-

portance in both the prefill and decoding phases.

Previous methods fall into two categories: (1) The *Prefill-Only Compression* method compresses the KV cache only during the prefill phase while retaining all KV cache generated during the decoding phase. (2) The *Unified Compression* method treats both phases as a unified process. For *Prefill-Only Compression*, methods like SnapKV (Li et al., 2024) and PyramidKV (Cai et al., 2024), retaining all KV cache generated during the decoding phase, leading to linear cache growth with the output length and memory pressure 🔥, especially for long outputs, as shown in Figure 1. For *Unified Compression*, such as H₂O (Zhang et al., 2023) and PyramidInfer (Yang et al., 2024b), prioritizes retaining the KV cache generated during decoding while discarding the earlier KV cache influenced by recent tokens typically receiving higher attention weights (Zhao et al., 2021; Song et al., 2024). This poses substantial challenges for reasoning tasks that rely on understanding the fine-grained input content 📄. There has been no dedicated exploration of KV cache compression strategies for handling lengthy outputs.

In this paper, we first unravel two essential observations that serve as the foundation for our exploration: (i) excessive compression during the prefill phase significantly affects the ability of LLM to reason through the query; (ii) heavy hitters deviate during the decoding phase in long-text generation, leading to skewed KV cache allocation. Building upon the insight, we introduce 🦋SCOPE, a simple yet efficient framework that Separately performs KV Cache Optimization during the Prefill and dEcoding phases. To our knowledge, we are the **first** to decouple the prefill and decoding phases to compress the KV cache independently. Specifically, we first maintain the KV cache generated during the prefill phase to ensure an understanding of long content. Then, we allocate heavy hitters using the sliding way in the decoding phase to optimize the memory of the KV cache. Building on the intuitive slide strategy, we further optimize *memory-usage* and *memory-transfer*, introducing adaptive strategy and discontinuous strategy.

To thoroughly validate our framework, we select LONGGENBENCH (Liu et al., 2024c) as the benchmark for our experiments over two mainstream LLMs. SCOPE can achieve comparable performance to the full KV cache when the overall compression rate is 35%. Additionally, our framework is seamlessly compatible with other compression

methods in the prefill phase.

The contributions of this work are as follows: 1). A simple yet efficient framework SCOPE is proposed to address the deviation of heavy hitters inspired by the observations and insights from an inference perspective. 2). Three strategies are developed to mitigate the deviation during the decoding phase. 3). Empirically, extensive experiments and analytical evaluations validate the effectiveness and generalizability of SCOPE.

2 Pilot Observation

2.1 KV Cache in Inference Perspective

Each request for an LLM involves two distinct phases (Zhou et al., 2024). The first phase, known as *prefill*, processes the complete input prompt to generate the initial output token. The second phase, termed *decoding*, iteratively produces the remaining output tokens, one at a time. We conduct pilot experiments through the lens of each phase in the inference process.

Prefill Phase: Existing work focusing on the prefill phase is grounded in the notion that attention is naturally sparse in typical tasks (Singhanian et al., 2024; Tang et al., 2024; Wu et al., 2024). For *PassageRetrieval-en* and *HotpotQA* within LongBench, a 20% compression ratio during the prefill phase still maintained performance nearly identical to that of the full cache, demonstrating the model’s ability to effectively retrieve and understand context even with significant compression, as shown in Figure 2a. However, when tasks require specific full context, such as reasoning tasks, attention is **not** always highly sparse (Chen et al., 2024), even if the output is short. As illustrated in Figure 2a, the same 20% compression rate during the prefill phase resulted in nearly 95% degradation in accuracy on the GSM8k+ task within LONGGENBENCH. Although sufficient performance is achieved on conventional tasks using KV cache compress during the prefill phase, the performance is notably poor on reasoning tasks when the compression ratio reaches a modest threshold, leaving room for targeted optimization through compression during the decoding phase.

Observations (i): For tasks that require specific fine-grained context, such as reasoning tasks, excessive compression during the prefill phase significantly compromises performance.

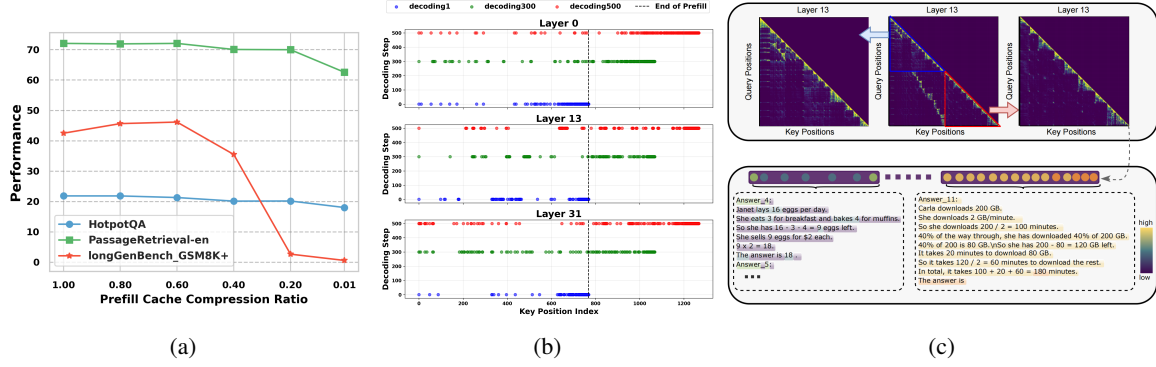


Figure 2: (a) Performances across various compression ratios during the prefill phase on three tasks under the full decoding cache condition. (b) Position distribution of the heavy hitters, selected by top 15% attention scores, at decoding steps 1, 300, and 500 across layers 0, 13, and 31. (c) Attention heatmaps for layer 13 of a GSM8k+ sample in LONGGENBENCH and details of the correspondence between attention scores and generated token positions. The complete case employed in the probing experiment is presented in Appendix 9.

Decoding Phase: We analyze the distribution of heavy hitters during the prefill and decoding phases as the decoding length increased in Figure 2b. Across all three layers, the retained heavy hitters predominantly originate from the KV cache generated during the decoding phase. This phenomenon has also been mentioned by several recent studies and can be attributed to the inherent properties of the attention mechanism, wherein tokens near the end often receive higher attention weights (Zhao et al., 2021; Song et al., 2024). This is particularly harmful for multi-question answering tasks, like LONGGENBENCH, as addressing such queries needs careful consideration of the question context. Previous prefill-only or unified compression strategies may overlook this distinction. In long-output tasks, as the output length increases, the deviation becomes more pronounced, making it imperative to preserve the heavy hitters identified during the prefill phase while providing appropriate management for those emerging in the decoding stage.

Observations (ii): During the decoding phase of long text generation, the use of the greedy algorithm may lead to a deviation in heavy hitters.

2.2 KV Cache Budget Reallocation

Building on the empirical observations from our pilot experiments, we derive the following insight:

Insight: It is crucial to allocate the budget of the KV cache during the prefill and decoding phases separately.

This insight inspires the design of SCOPE, which decouples compression into the prefill and decoding phases to effectively allocate the KV cache budget, preserving all KV cache generated during the prefill phase and enabling more effective reallocation of the KV cache budget. While numerous studies have explored the heavy hitters during the prefill phase, to our knowledge, no prior work has specifically addressed this aspect of the decoding phase. We dive deeper into the sparsity in the KV cache during the decoding to design strategy, selecting essential heavy hitters dynamically. To gain deeper insights, following prior works (Xiao et al., 2024b; Cai et al., 2024), we analyze the attention heatmaps, comparing the attention weights between the prefill and decoding phases, as shown in Figure 2c. The leftmost and rightmost plots represent the prefill and decoding phases, respectively. For tasks that require simultaneous reasoning for multiple questions, it is essential to recognize the position of the current prediction. This information can be captured by heavy hitters identified using a greedy algorithm, as illustrated in Figure 2c. Thus, it remains necessary to allocate a portion of the KV cache budget specifically for heavy hitters. Furthermore, owing to the autoregressive nature of LLMs, it remains essential to retain the recent tokens, which exhibit stronger correlations with current tokens.

3 Method

3.1 Revisiting KV Cache Compression

Initialization KV cache compression essentially involves adjusting the cache based on the given KV

cache budget, where we allocate a cache pool, denoted as Φ , consisting of Φ^p and Φ^d , which stores the KV cache generated during the prefill and decoding phases, respectively. The cache pool is updated at each step t , denoted as Φ_t . The widely recognized function for selecting heavy hitters based on the greedy algorithm is denoted as $\Psi_K(\mathbf{Att})$, which represents the selection of the Top- K KV caches from the given attention weights \mathbf{Att} .

Prefill Phase Given the input prompt tensor $\mathbf{P} \in \mathbb{R}^{M \times D}$, represented as $\mathbf{P} = \{\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_M\}$, where \mathbf{P}_i denote i -th token embeddings, and M represent the number of input tokens and D is the model’s hidden dimension. The key and value tensors are computed as follows:

$$\mathbf{K}_{\mathcal{P}}\mathbf{V}_{\mathcal{P}} = \mathbf{P}\mathbf{W}_K, \mathbf{P}\mathbf{W}_V, \quad (1)$$

where $\mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D \times D}$ are the weights matrices for the key and value projections, respectively. The KV pairs are denoted as $\mathbf{K}_{\mathcal{P}}\mathbf{V}_{\mathcal{P}}$. The attention weights $\mathbf{Att}_{\mathcal{P}}$ is calculated by \mathbf{P} and $\mathbf{K}_{\mathcal{P}}\mathbf{V}_{\mathcal{P}}$. The most effective and widely adopted approach, as established through early explorations (Zhang et al., 2023; Yang et al., 2024b; Li et al., 2024), two import hyperparameters α_1 and α_2 are introduced, where α_1 represents the length of *prefill essential history window* and α_2 represents the length of *prefill local window* during the prefill phase. The length of the total reserved KV cache is $\alpha_1 + \alpha_2$, which also corresponds to the size of the cache pool Φ^p during the prefill. For compression during the prefill phase is:

$$\mathbf{K}_0\mathbf{V}_0 = \Psi_{\alpha_1}(\mathbf{Att}_{\mathcal{P}}[: -\alpha_2]) \cdot \mathbf{K}_{\mathcal{P}}\mathbf{V}_{\mathcal{P}}[-\alpha_2 :], \quad (2)$$

where \cdot denotes concatenation and the function $\Psi_{\alpha_1}(\mathbf{Att}_{\mathcal{P}})$ selects the KV cache with the Top- α_1 attention weights from $\mathbf{Att}_{\mathcal{P}}[: -\alpha_2]$.

$\mathbf{K}_0\mathbf{V}_0$ is stored in Φ_0^p . Maintain an *essential history window* α_1 to retain KV with higher attention weights for the current query and a *local window* α_2 to reserve the KV of recently generated tokens, ensuring both contextual continuity and retention of attention. Notably, the compression is only executed once, at $t = 0$, marking the end of the prefill phase before transitioning into the decoding phase.

Decoding Phase During the decoding phase, the KV cache from the prefill phase is employed and updated to sequentially generate tokens. At each time step t , keys and values are computed only for

the new token tensor $\mathbf{X}_{t, t \in \{1, T\}}$ as follows:

$$\mathbf{K}_t\mathbf{V}_t = \mathbf{X}_t\mathbf{W}_K, \mathbf{X}_t\mathbf{W}_V, \quad (3)$$

$\mathbf{K}_t\mathbf{V}_t$ is concatenated with previously retained KV cache, which is stored in Φ , to obtain the current retained KV pairs. This is then computed with the current query \mathbf{X}_t to compute the attention \mathbf{Att}_t .

The main difference from previous KV compression methods lies in the distribution of Φ^p and Φ^d within the cache pool Φ . The *Prefill-Only Compression* method does not compress the KV cache generated during the decoding phase. Instead, it involves a linear growth of the KV cache with each newly generated token. Φ_p^t remains constant, and at each step t , it stores the originally preserved $\mathbf{K}\mathbf{V}_0$. Φ_d^t stores the KV cache at each time step t during the decoding phase, from $\mathbf{K}_1\mathbf{V}_1$ to $\mathbf{K}_T\mathbf{V}_T$, which leads to a significant increase of memory consumption as the length grows. The *Unified Compression* method in the decoding phase will apply the $\Psi_{\alpha_1}(\mathbf{Att}_t[: -\alpha_2])$ at each t to update cache pool Φ . As the number of generated tokens increases, the attention mechanism tends to assign higher weights to tokens at the end, meaning that the Top- α_1 KV caches returned by Ψ are all generated during the decoding phase, while those from the prefill phase are discarded. As t increases, Φ_t^p grows larger, while Φ_t^d becomes smaller. This results in more information being retained in Φ^d within Φ , while the information in Φ^p decreases, leading to potential essential information loss that may be needed in future decoding steps.

3.2 SCOPE

The primary goal of **SCOPE** is to mitigate the deviation of heavy hitters, thereby ensuring a more balanced allocation of Φ^p and Φ^d . Motivated by the findings in §2.1, where excessive compression during the prefill phase hinders performance on reasoning tasks, the cache pool Φ_p is constant at each t , i.e., we reserved all compressed KV Cache generated during the prefill phase. The operation on Φ^p in **SCOPE** is the same as that in the previously unified compression method.

It is necessary to leverage the sparsity of the KV cache generated during decoding to enable efficient allocation. Three strategies for the decoding phase are developed: **Slide**, **Adaptive**, and **Discontinuous**, all of which update only Φ_d . The adaptive strategy optimizes memory based on the slide strategy, and the discontinuous strategy optimizes computation on top of the adaptive one. We

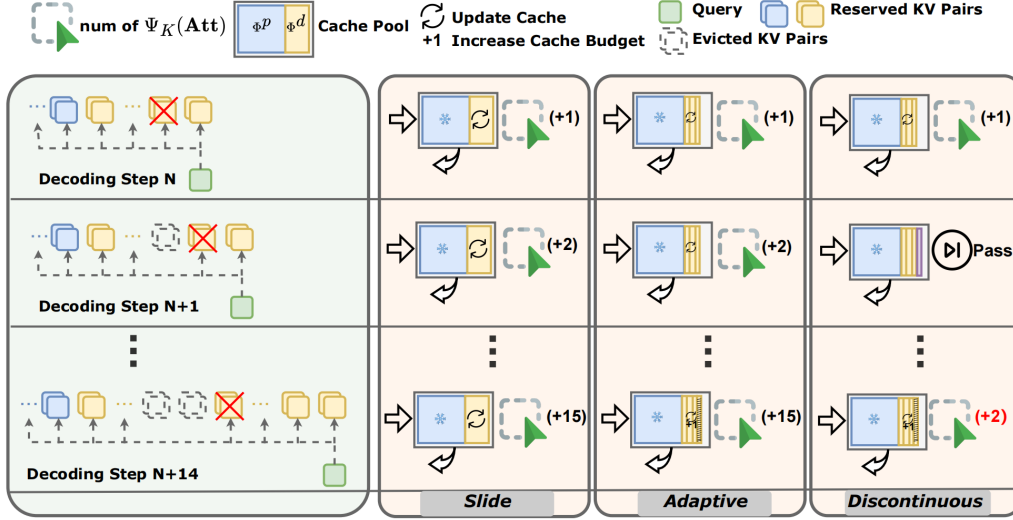


Figure 3: Illustration of three strategies of **SCOPE**. The prefilled cache pool Φ_t^p is constant at each t . The slide strategy updates the decoding cache pool at each decoding step while the size of the decoding cache pool is constant. The adaptive strategy incrementally increases the size of the decoding cache pool at regular intervals of $\frac{T-\beta_2}{\beta_1}$. The discontinuous strategy, built upon the adaptive strategy, executes $\Psi_K(\text{Att})$ at intervals of the same time period.

will introduce them one by one below in detail. For each strategy, Python-style pseudocode is provided in Figure 8 to facilitate comprehension of details.

Slide We compress the KV cache in the decoding phase by sliding the *decoding essential history window* β_1 and the *decoding local window* β_2 , where β_1 helps identity the position of the current prediction and β_2 stores global information strongly correlated with previous tokens other as discussed in §2.2.

The slide strategy starts from $t > M + \beta_1 + \beta_2$, applying the function $\Psi_{\beta_1}(\text{Att}_t[\alpha_1 + \alpha_2 : -\beta_2])$ to restrictively update only Φ_t^d while keeping Φ_t^p constant. This is achieved by limiting the selecting function Ψ to operate on Att_t starting from $\alpha_1 + \alpha_2$, thus excluding the attention weights from the prefill phase. It can be completely independent of the KV cache pool during the prefill phase Φ^p .

Adaptive We can optimize the β_1 of slide strategy to adaptively increase its size from a *memory-usage* perspective. When the length of the tokens generated during decoding is relatively short, a long *decoding essential history window* β_1 is unnecessary, it is unnecessary to place all these KV caches in the Φ_t^d . β_1 can be adaptively increased as needed. A function of time steps t and the maximum length T is proposed to adaptively adjust the length of the *decoding essential history window* β_1 , where $T \gg \beta_1 + \beta_2$. It starts with a

base size β_2 and grows linearly with time step t :

$$\hat{\beta}_1 = \frac{(t - \beta_2) \cdot \beta_1}{T - \beta_2} \text{ if } t > \beta_2, \quad (4)$$

The budget size of Φ_t^d also increases adaptively and is given by $\beta_2 + \frac{(t-\beta_2)\beta_1}{T-\beta_2}$ when $t < T$, which helps optimize memory usage, as the ratio $\frac{(t-\beta_2)}{(T-\beta_2)}$ is less than 1. As t reaches T , the size of Φ_t^d becomes $\beta_1 + \beta_2$. This adjustment aligns with the autoregressive token-by-token encoding characteristic of LLMs, ensuring more efficient use of resources. In addition, $\Psi_K(\text{Att})$ begins execution earlier than the sliding strategy. The adaptive strategy optimizes the budget of Φ_t^d and reduces unnecessary overhead while still retaining enough historical context for an effective generation without introducing additional hyperparameters.

Discontinuous We further optimize the adaptive strategy from a *memory-transfer* perspective to ensure by reducing the frequency of execution of $\Psi_K(\text{Att})$. The top- K selection operation $\Psi_K(\text{Att})$ would be executed a total of $T - \beta_2$ times using previous strategies, which potentially leads to frequent GPU I/O due to the update operation of Φ_d at each step. Motivated by the characteristic that consecutive queries tend to select similar keys (Zhao et al., 2024; Tang et al., 2024), we make the update operation, i.e., Top- K selection $\Psi_K(\text{Att})$ discontinuous. This strategy optimizes the times of execution frequency of selection

Table 1: Performance of our proposed SCOPE using three strategies and baselines on the LONGGENBENCH benchmark with LLaMA-3.1-8B-Instruct and Mistral-7B-Instruct-v0.3. The best results among all methods are in **bolded**. The prefill compression ratio averages around 60%.

Method	LONGGENBENCH-4K				LONGGENBENCH-8K			
	GSM8K+	MMLU+	CSQA+	Avg.	GSM8K++	MMLU++	CSQA++	Avg.
<i>LLaMA-3.1-8B-Instruct</i>								
Full Cache	53.26	54.40	71.67	59.78	44.44	51.01	64.92	53.46
	Decoding Compression Ratio=25.0%				Decoding Compression Ratio=12.5%			
StreamingLLM	10.78	34.15	43.50	29.48	26.98	41.26	65.33	44.53
H ₂ O	35.04	48.11	69.50	50.88	22.54	48.21	59.58	43.44
PyramidInfer	38.76	48.93	71.58	53.09	21.11	47.96	59.58	42.88
SCOPE (Slide)	46.51	46.62	72.50	56.21	30.24	51.64	65.75	49.21
SCOPE (Adaptive)	43.10	50.25	71.50	54.95	27.86	51.23	62.50	47.19
SCOPE (Discontinuous)	42.02	49.75	72.67	54.81	24.37	50.00	59.33	44.57
	Decoding Compression Ratio=12.5%				Decoding Compression Ratio=6.25%			
StreamingLLM	11.94	35.97	41.42	29.78	20.56	41.79	64.92	42.42
H ₂ O	26.59	45.97	65.25	45.94	21.27	45.94	55.08	40.77
PyramidInfer	28.29	46.41	61.42	45.38	19.84	45.50	55.08	40.14
SCOPE (Slide)	42.56	50.94	73.50	55.67	26.59	49.59	65.08	47.09
SCOPE (Adaptive)	37.29	50.19	74.00	53.83	30.56	49.94	65.92	48.80
SCOPE (Discontinuous)	39.85	50.06	72.92	54.27	28.41	50.63	67.92	48.99
<i>Mistral-7B-Instruct-v0.3</i>								
Full Cache	11.01	28.30	64.33	34.55	9.37	20.35	51.75	27.15
	Decoding Compression Ratio=25.0%				Decoding Compression Ratio=12.5%			
StreamingLLM	6.90	22.83	65.25	31.66	2.62	17.48	46.75	22.28
H ₂ O	7.91	26.48	60.42	31.60	5.71	16.20	40.17	20.69
PyramidInfer	10.00	24.15	62.92	32.36	5.56	16.48	40.17	20.73
SCOPE (Slide)	7.67	21.51	58.58	29.26	5.95	16.95	45.50	22.80
SCOPE (Adaptive)	11.47	29.06	64.50	35.01	9.76	20.35	51.75	27.29
SCOPE (Discontinuous)	11.55	29.06	64.50	35.04	9.84	20.35	51.75	27.31
	Decoding Compression Ratio=12.5%				Decoding Compression Ratio=6.25%			
StreamingLLM	6.51	19.69	57.92	28.04	3.57	17.14	46.83	22.51
H ₂ O	7.13	21.07	49.83	26.01	5.63	16.07	33.25	18.32
PyramidInfer	7.13	20.94	51.75	26.61	5.63	16.76	33.25	18.55
SCOPE (Slide)	8.84	18.68	51.75	26.42	5.71	17.08	46.17	22.99
SCOPE (Adaptive)	10.93	29.06	64.50	34.83	7.30	19.94	51.75	26.33
SCOPE (Discontinuous)	10.39	29.06	64.50	34.65	8.33	20.19	51.75	26.76

operation $\Psi_K(\text{Att})$, with ζ occurring once every interval of $\frac{T-\beta_2}{\beta_1}$, whereas previous strategies execute at each step t . This interval is consistent with the growth of $\hat{\beta}_1$ in the adaptive strategy. The frequency can be reduced to $\frac{T-\beta_2}{\beta_1} = \beta_1$ using this strategy, thereby alleviating the memory I/O pressure caused by frequent updates to Φ_d .

4 Experiments

4.1 Datasets

We develop two **open-sourced** datasets, LONGGENBENCH-4K ({subtask}+) and LONGGENBENCH-8K ({subtask}++), where multiple reasoning tasks must be handled

simultaneously³, each containing three sub-tasks synthesized from GSM8K (Cobbe et al., 2021), MMLU (Hendrycks et al., 2021), and CSQA (Hendrycks et al., 2021). These subtasks are designed to address long-input challenges with output lengths of 4K and 8K, respectively.⁴ To validate the effectiveness of SCOPE on general long-output tasks, we select the En.Sum task from ∞ BENCH (Zhang et al., 2024), with an average output length of 1.1K. For the detailed statistics of datasets and additional details corresponding to each subtask, refer to Appendix B.

³Prompt template is provided in Appendix B.

⁴The selected examples have output lengths of 4K and 8K, ensuring no premature cessation of the response.

4.2 Baselines

To validate the effectiveness of SCOPE, we compare it with **Full Cache** and representative unified compression methods, including **StreamingLLM** (Xiao et al., 2024b), which keeps the KV of early and recent tokens; **H₂O** (Zhang et al., 2023), which balances recent and Heavy Hitter (H₂) tokens based on cumulative attention scores; and **PyramidInfer** (Yang et al., 2024b), which reduces the cache in deeper layers using sparse attention patterns. To validate modularity, we apply SCOPE in combination with **SnapKV** (Li et al., 2024) and **PyramidKV** (Cai et al., 2024) during the decoding phase, as detailed in §4.4.

4.3 Implementation Details

We build SCOPE using two open-sourced LLMs, specifically LLaMA-3.1-8B-Instruct and Mistral-7B-Instruct-v0.3. Based on the preliminary experiments (Figure 2a), we set the size of Φ_p , i.e., $\alpha_1 + \alpha_2$ to 2048 for LongGenBench-4K and 4096 for LongGenBench-8K, corresponding to approximately 60% of the input length. α_2 is set to 8 following previous works (Cai et al., 2024; Li et al., 2024). $\beta_1 + \beta_2$ are set to 512 and 1024 in two configurations, corresponding to different compression ratios for outputs of 4K and 8K. β_2 is set to 256 to accommodate the CoT length in answers, avoiding performance loss from overly short sequences. For a fair comparison, the total budget of KV cache during both the prefill and decoding phases is consistent across all methods. More details can be found in the Appendix C.

4.4 Results

Comparison with Baselines Table 1 presents a comprehensive analysis of our proposed SCOPE and baselines. SCOPE (with three strategies) achieves the best results under both decoding compression methods, and the discontinuous strategy, optimized for memory-usage and memory-transfer, delivers outstanding performance. On the challenge GSM8K+/GSM8K++ tasks, SCOPE highlights the importance of preserving the KV cache generated during the prefill process, while other compression methods lead to marked performance degradation. This ensures that the understanding of the problem statement remains intact, achieving comparable performance to the full cache without compromising comprehension. StreamingLLM poses challenges on LONGGENBENCH, where

Table 2: The plug-in experiment results of LLaMA3.1-8B on the GSM8K+ task from LONGGENBENCH-4K. The results comparable to the full cache are in **bold**.

Decoding Phase Strategy	Prefill Phase		
	Full Cache	SnapKV	PyramidKV
Full Cache	53.26	27.75	27.75
	Decoding Compress Ratio=25.0%		
Slide	52.17	26.90	26.90
Adaptive	43.88	27.60	27.60
Discontinuous	47.21	27.67	27.67
	Decoding Compress Ratio=12.5%		
Slide	49.69	22.56	22.56
Adaptive	44.03	27.60	27.60
Discontinuous	46.98	27.67	27.67

vital information may lie within the middle of the input, consistent with the findings in prior study (Zhang et al., 2023). This inevitably results in the loss of crucial information if only the first few tokens and local tokens are preserved. Performance between PyramidInfer and H₂O shows no notable difference, indicating that the layer-wise sparsity feature is not prominent for tasks with long outputs.

Plug-in to Prefill-Only Methods Table 2 shows the results of seamlessly integrating our decoding phase compression strategy with prefill-only compression methods. Some strategies even outperform the full cache results, despite compressing 35%⁵ of the KV cache. This validates the sparsity of the KV cache generated during the decoding phase in multi-QA tasks and demonstrates the effectiveness of our proposed strategies. PyramidKV (Cai et al., 2024), a variant of SnapKV, adjusts the budget allocation across layers without observing improvements in the preliminary experiments, consistent with the empirical finding (§4.4).

Actually, the retained KV cache during the prefill phase can be regarded as “*attention sinks*”, which bears a resemblance to the principle of StreamingLLM. We extend this concept to broader, more realistic long-output scenarios.

5 Analysis and Discussion

5.1 Mitigating the Loss of Essential H₂

The unified compression method, such as H₂O, suffers from the loss of crucial KV cache generated during prefill, which is essential to understanding

⁵The average input-output length is 7.4K in the GSM8K+ task. With budgets Φ_p of 2K and Φ_d of 0.5K, the total reserved KV cache size is 2.5K, leading to a full compression ratio of about 35%.

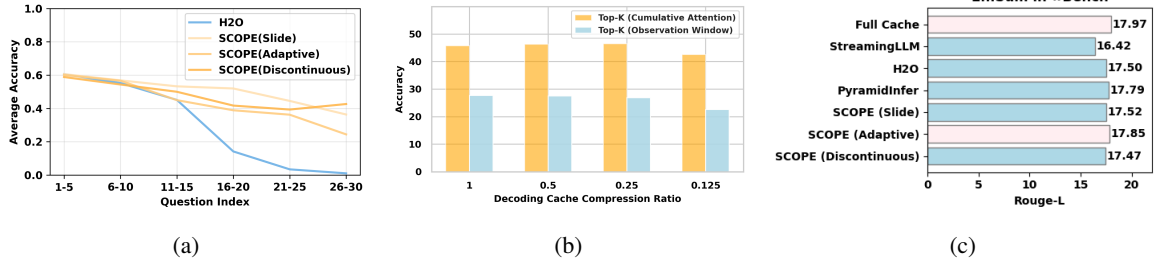


Figure 4: (a) Accuracy distribution of different question positions. (b) Accuracy across different cache compression ratios during the decoding phase using two Tok- K selection algorithms while the KV cache during the prefill phase is compressed to a 60% ratio. Top- K (Observation Window) (Li et al., 2024; Cai et al., 2024), is computed within a fixed-size local window of recent key-value pairs while Top- K (Cumulative Attention) (Zhang et al., 2023; Yang et al., 2024b), attention is computed globally across all key-value pairs. (c) Results on En.Sum task from ∞ BENCH, with the condition $\beta_1 + \beta_2 = 512$.

the context due to the deviation of heavy hitters. In Figure 4a, we show the relationship between prediction position and performance. The performance of H2O drops markedly in later predictions, while all three of our strategies mitigate this decline, validating the effectiveness of preserving the prefill KV cache.

5.2 Influence on $\beta_1 + \beta_2$ and $\Psi_K(\text{Att})$

KV cache budget during the decoding phase $\beta_1 + \beta_2$ and selection algorithm $\Psi_K(\text{Att})$ are the key hyperparameters within the SCOPE framework. The budget $\beta_1 + \beta_2$, i.e., the compression ratio is scaled using two mainstream top- K selection algorithms as illustrated in Figure 4b. Unlike the prefill phase, where performance on the GSM8k+ task significantly drops as the compression ratio increases, compressing to 25% during the decoding phase only results in a 15% performance decline. It validated that compression in both phases is better than solely focusing on extreme compression during prefill and the necessity of optimizing the KV cache **separately** for the prefill and decoding phases. Using the Top- K selection strategy based on cumulative attention yields better results than the Top- K selection strategy based on the observation window. For tasks like LONGGENBENCH, predictions still require reviewing and capturing the corresponding question, making a short observation window insufficient.

5.3 Efficiency on Memory Usage and Transfer

Our adaptive and discontinuous strategies building on slide strategy improve memory efficiency, as explored in Table 3. Compared to the full cache and prefill-only compression methods, both our method and the unified compression approach effectively reduce memory usage pressure by storing less KV

Table 3: Efficiency analysis on Peak KV memory and latency (Lat.) for LLaMA3.1-8B with a prefill compression ratio of 60% and a decoding compression ratio of 12.5%.

Method	Peak KV Mem.	Tokens/s
Full Cache	15.6(100%)	36.57
SnapKV	12.5(80.1%)	38.28
PyramidKV	12.5(80.1%)	36.90
StreamingLLM		22.02
H2O	5.8(37.1%)	21.78
PyramidInfer		22.38
SCOPE (Slide)		18.28
SCOPE (Adaptive)	5.8(37.1%)	18.28
SCOPE (Discontinuous)		25.92

cache overall. Our adaptive strategy further optimizes performance by dynamically adjusting the budget. However, this introduces frequent updates to the stored KV cache pool, leading to increased I/O transfer and latency. The optimized strategy effectively mitigates this issue by executing computations discontinuously.

5.4 Generalization of SCOPE

Results of our proposed SCOPE and baselines on ∞ BENCH are shown in Figure 4c. Adaptive strategy demonstrates the **closest** performance alignment with the full cache setting, effectively validating its generalized capability design. It is effective not only for multi-QA tasks but also for summarization tasks, demonstrating that traditional tasks may also be suited to the separation of prefill-decoding KV cache budget allocation.

6 Related Work

KV Cache Compression KV cache compression methods focus on leveraging the sparsity in attention to address memory bottlenecks,

complementing other efficient techniques (Kwon et al., 2023; Dao, 2024; Wang et al., 2024; Liu et al., 2024d). While recent work has optimized the prefill phase by adjusting the compression budget (Yang et al., 2024b; Feng et al., 2024; Cai et al., 2024), phase-specific optimization remains unexplored. Our approach tailors KV cache compression to the distinct characteristics of each phase, offering a novel perspective.

Long-context Tasks Recent advancements in LLMs have focused on enhancing the capabilities for long-context tasks. Previous evaluations of long-context tasks have mainly concentrated on tasks with long inputs, and numerous benchmarks have been proposed, such as Needle-in-a-Haystack (NIAH) (Kamradt, 2023), LongBench (Bai et al., 2024) and ∞ BENCH (Zhang et al., 2024) for comprehensive understanding tasks, where the output is generally short for most sub-tasks. Most research on KV cache compression has been conducted within the context of these benchmarks, where the focus has been primarily on optimizing the prefill phase. In this work, we leverage LONGGENBENCH, which focuses on long-input and long-output tasks (Liu et al., 2024c), to optimize KV cache compression in scenarios where the output can be as long as 8K tokens.

7 Conclusion

In this paper, we propose SCOPE, a framework that optimizes KV cache usage for long-context generation in LLMs. We observe that excessive compression during the prefill phase harms reasoning capabilities while the deviation of heavy hitters during decoding. To resolve these issues, SCOPE preserves essential KV cache during the prefill phase and employs a sliding strategy to efficiently manage the KV cache generated during decoding. Additionally, we introduce adaptive and discontinuous strategies to further optimize memory usage and transfer. Our extensive experiments demonstrate that SCOPE achieves near-full KV cache performance with only 35% of the original memory while remaining compatible with existing prefill compression methods.

Limitations

SCOPE separates the prefill and decoding phases for long-text generation tasks, while a Top- K algorithm is used to select the *heavy hitters* in both

the prefill and decoding phases. We discuss the following limitations:

Prefill Phase We employ the widely recognized top- K algorithm during the prefill phase, and future work could explore chunking or other techniques (Song et al., 2024; Xu et al., 2024) to further enhance the estimation of previous tokens. As discussed in §4.4, the retained KV cache during the prefill phase can be regarded as an “*attention sinks*”. Enhancing the quality of this overall “*attention sinks*” is a potential direction for future research. Moreover, our phase-level approach is orthogonal to other KV reuse methods (Xiao et al., 2024a; Lee et al., 2024; Liu et al., 2024a) and could be integrated with these techniques to further optimize memory management and computation efficiency.

Decoding Phase The execution of Top- K at each decoding step is time-costly due to the frequent GPU I/O. Though we optimize the operation frequency in the discontinuous strategy, we can also reduce the I/O size to lower latency. Specifically, by leveraging the PD-separated framework, optimizing I/O for just Φ_d would be more efficient, as the size of Φ_p is constant, while we currently update the entire Φ .

Modality Although SCOPE has shown advantages for long-output tasks in the text modality, there is potential for our method to be applied to long-output tasks in **vision**, such as multi-image generation, where the KV cache required for storing each image is substantial.

Dataset Our experiments demonstrate the effectiveness of SCOPE on **two** well-established benchmarks: LONGGENBENCH and ∞ BENCH. In both benchmarks, our strategies consistently outperform the baseline, highlighting the generalization of SCOPE. While these results are robust, we also expect to evaluate SCOPE on more diverse and challenging benchmarks in the future, further validating its scalability and broader applicability.

Acknowledgement

The authors would like to thank the anonymous reviewers for their insightful comments. This work is funded by the National Natural Science Foundation of China (Grant No.62176053). This work is supported by the Big Data Computing Center of Southeast University.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. [arXiv preprint arXiv:2303.08774](#).
- Anthropic. 2024. [The claude 3 model family: Opus, sonnet, haiku](#). Accessed: 2024-07-09.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. [LongBench: A bilingual, multitask benchmark for long context understanding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137, Bangkok, Thailand. Association for Computational Linguistics.
- Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie Hu, et al. 2024. Pyramidkv: Dynamic kv cache compression based on pyramidal information funneling. [arXiv preprint arXiv:2406.02069](#).
- Zhuoming Chen, Ranajoy Sadhukhan, Zihao Ye, Yang Zhou, Jianyu Zhang, Niklas Nolte, Yuandong Tian, Matthijs Douze, Leon Bottou, Zhihao Jia, et al. 2024. Magicpig: Lsh sampling for efficient llm generation. [arXiv preprint arXiv:2410.16179](#).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. [arXiv preprint arXiv:2110.14168](#).
- Tri Dao. 2024. [Flashattention-2: Faster attention with better parallelism and work partitioning](#). In *The Twelfth International Conference on Learning Representations*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. [arXiv preprint arXiv:2407.21783](#).
- Yuan Feng, Junlin Lv, Yukun Cao, Xike Xie, and S Kevin Zhou. 2024. Ada-kv: Optimizing kv cache eviction by adaptive budget allocation for efficient llm inference. [arXiv preprint arXiv:2407.11550](#).
- Yao Fu. 2024. Challenges in deploying long-context transformers: A theoretical peak performance analysis. [arXiv preprint arXiv:2405.08944](#).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. [arXiv preprint arXiv:2310.06825](#).
- Greg Kamradt. 2023. [Llms need needle in a haystack: Test-pressure testing llms](#). Accessed: 2024-11-20.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.
- Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. 2024. {InfiniGen}: Efficient generative inference of large language models with dynamic {KV} cache management. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, pages 155–172.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024. [SnapKV: LLM knows what you are looking for before generation](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Guangda Liu, Chengwei Li, Jieru Zhao, Chenqi Zhang, and Minyi Guo. 2024a. Clusterkv: Manipulating llm kv cache in semantic space for recallable compression. [arXiv preprint arXiv:2412.03213](#).
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024b. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Xiang Liu, Peijie Dong, Xuming Hu, and Xiaowen Chu. 2024c. [LongGenBench: Long-context generation benchmark](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 865–883, Miami, Florida, USA. Association for Computational Linguistics.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. 2024d. [KIVI: A tuning-free asymmetric 2bit quantization for KV cache](#). In *Forty-first International Conference on Machine Learning*.
- Shi Luohe, Hongyi Zhang, Yao Yao, Zuchao Li, and hai zhao. 2024. [Keep the cost down: A review on methods to optimize LLM’s KV-cache consumption](#). In *First Conference on Language Modeling*.