

More Tokens, Lower Precision: Towards the Optimal Token-Precision Trade-off in KV Cache Compression

Anonymous ACL submission

Abstract

As large language models (LLMs) process increasing context windows, the memory usage of KV cache has become a critical bottleneck during inference. The mainstream KV compression methods, including KV pruning and KV quantization, primarily focus on either token or precision dimension separately. However, these works leaving the trade-off between these two orthogonal dimensions largely under-explored. In this paper, we leverage the Information Bottleneck principle to formulate KV cache compression within a unified theoretical framework. We demonstrate that a carefully managed token-precision trade-off can achieve an optimal point within the Information Bottleneck compared to standalone KV pruning or KV quantization. Experiments reveal that storing more tokens in the KV cache at lower precision—a strategy we term quantized pruning—can significantly enhance the long-context performance of LLMs. An in-depth analysis of this token-precision trade-off across key aspects shows that quantized pruning achieves substantial improvements in retrieval-related tasks and consistently performs well across varying input lengths. Furthermore, quantized pruning exhibits notable stability and effectiveness across different KV pruning methods, quantization strategies, and model scales. These findings offer valuable insights into optimizing KV cache compression through balanced token-precision trade-off strategies.

1 Introduction

Current long-context Large Language Models (LLMs) heavily depend on the Key-Value (KV) cache mechanism, which stores intermediate keys and values during text generation to avoid redundant computations (Waddington et al., 2013). As the input sequence length increases, the KV cache size expands proportionally, leading to substantial memory overhead (Achiam et al., 2023; Reid et al., 2024). Take Llama3-8B (Dubey et al., 2024) as an

example, storing the KV cache of a single sequence with 100k tokens necessitates a substantial memory overhead of 20GB — exceeding the model’s parameter memory of approximately 14GB. Moreover, since the decoding process is highly dependent on GPU memory bandwidth, the large KV cache also results in a significant increase in decoding time (Fu, 2024). As a result, efficiently compressing the KV cache has become a critical challenge in advancing LLM development (Pope et al., 2023).

To optimize KV cache memory usage, two primary approaches, focusing on token and precision dimensions, are widely adopted: KV pruning and KV quantization. For the token dimension, KV pruning methods reduce the KV cache footprint by discarding less salient tokens, thereby maintaining a fixed cache size (Xiao et al., 2023; Zhang et al., 2024b; Ren and Zhu, 2024; Li et al., 2024). For the precision dimension, KV quantization methods reduce memory usage by approximating KV cache with lower precisions, like 4-bit or even lower (Sheng et al., 2023; Liu et al., 2024c; Hooper et al., 2024; Yang et al., 2024b).

From the perspective of the Information Bottleneck (IB) principle, we propose a unified theoretical framework to describe KV pruning and KV cache quantization, and find an exclusive focus on optimizing either pruning or quantization in isolation may lead to a suboptimal trade-off between model performance and compression efficacy.

In this paper, we demonstrate that a strategy co-optimizing both pruning and quantization aspects can achieve an optimal operating point within the information-theoretic trade-off space than optimizing either aspect in isolation. Specifically, we uncovered an intriguing finding: *storing more tokens in the KV cache with lower precision*, a strategy we term *quantized pruning*, can significantly enhance the long-context performance of LLMs. For instance, storing $4\times$ tokens in 4-bit precision outperforms storing $1\times$ tokens in 16-bit precision

across various KV cache budget. Notably, in extremely low-resource scenarios, quantized pruning effectively preserves performance, whereas relying solely on KV pruning or quantization often leads to a significant performance collapse. Furthermore, we conduct in-depth analysis of the token-precision trade-off across critical aspects, including task type, input length, model scaling, quantization strategies, and layer-wise configurations. Extensive experiments reveal that quantized pruning achieves considerable gains in retrieval tasks and maintains robust performance across varying input lengths. Moreover, quantized pruning demonstrates strong stability across different KV pruning methods, quantization strategies, and model scales.

In summary, our contributions are as follows:

- We are the first to describe KV pruning and KV cache quantization from a unified theoretical framework, and find an optimal trade-off between model performance and compression efficacy leveraging token-precision trade-off.
- An important finding is revealed: storing more tokens with lower precision significantly outperforms standalone KV pruning or KV quantization methods under various KV cache budgets, highlighting the importance of balancing token count and precision in KV cache compression.
- Extensive experiments have been conducted on exploring token-precision trade-off across various critical aspects, including task types, input lengths, model scales, and quantization strategies, providing valuable insights for optimizing KV cache compression.

2 Related Work

KV Pruning KV pruning compresses the KV cache along the token dimension by selectively discarding less salient tokens to reduce memory usage. Mainstream methods typically identify salient tokens based on attention scores, as seen in (Liu et al., 2024b; Zhang et al., 2024b; Oren et al., 2024; Li et al., 2024). Other methods use alternative factors such as initial tokens (Xiao et al., 2023), variance (Ren and Zhu, 2024), special tokens (Ge et al., 2024) or the L2 norm (Devoto et al., 2024) to determine token importance. Recent studies delve deeper into optimizing the allocation of KV cache memory budgets. Some explore KV cache budget allocation strategies across layers (Cai et al., 2024; Yang et al., 2024a), while other studies explore head-level KV cache budget allocation (Feng et al.,

2024; Tang et al., 2024; Fu et al., 2024; Xiao et al., 2024).

KV Quantization KV quantization compresses KV cache from the precision dimension by storing KV cache using a reduced number of bits. FlexGen (Sheng et al., 2023) utilizes group-wise 4-bit quantization for both key and value cache. KIVI (Liu et al., 2024c) applies per-channel quantization on key cache and per-token quantization on value cache. KVQuant (Hooper et al., 2024) and CQ (Zhang et al., 2024a) use RoPE-related quantization, while KVQuant also preverges outliers without quantization. Atom (Zhao et al., 2024) and SKVQ (Duanmu et al., 2024) reorders the outlier channels for fine-grained group quantization with mixed-precision. QAA (Dong et al., 2024) and MiKV (Yang et al., 2024b), inspired by KV pruning methods, store discarded tokens from KV pruning methods using lower bit precision while retaining salient tokens in full precision. However, their memory usage scales proportionally with context length. In contrast, quantized pruning, fundamentally based on KV pruning, theoretically enables compression to a pre-defined memory budget.

Other KV Compression Methods Compressing KV cache from other dimensions typically requires modifying the model architecture, which usually necessitates additional training for adaptation. For the layer dimension, LCKV (Wu and Tu, 2024), CLA (Brandon et al., 2024) and MLKV (Zuhri et al., 2024) reduce memory usage by sharing the KV cache across adjacent layers. ShortGPT (Men et al., 2024) and DynamicSlicing (Dumitru et al., 2024) achieve compression by eliminating redundant layers. YOCO (Sun et al., 2024) changes the model structure and shares a single global KV cache across layers. For the head dimension, MQA (Shazeer, 2019) and GQA (Ainslie et al., 2023) share the KV cache within each head groups. DeepSeek-v2 (Liu et al., 2024a) employs dimension-reduction techniques to compress all heads into a single low-rank vector.

3 Preliminaries

During inference, KV cache is implemented within the self-attention module and operates in two distinct phases: i) the prefill phase, where the input prompt is used to generate KV cache for each transformer layer of LLMs; and ii) the decoding phase, where the model uses KV cache to generate the next token, and updates the KV cache with the new

token.

Prefill Phase. Let $X \in \mathbb{R}^{b \times l_{\text{prompt}} \times d}$ be the input tensor, where b is the batch size, l_{prompt} is the length of the input prompt, and d is the model hidden size. For clarity, we omit the layer index here. The key and value tensors can be computed by:

$$X_K = XW_K, X_V = XW_V \quad (1)$$

where $W_K, W_V \in \mathbb{R}^{d \times d}$ are the key and value layer weight. X_K, X_V are cached in the memory for utilization in the subsequent decoding phase.

Decoding Phase. Let $h \in \mathbb{R}^{b \times 1 \times d}$ be hidden state of current input token. $h_K = hW_K$ and $h_V = hW_V$ are the current key and value states. h_K and h_V are first employed to update the KV cache:

$$\begin{aligned} X_K &\leftarrow \text{Concat}(X_K, h_K), \\ X_V &\leftarrow \text{Concat}(X_V, h_V) \end{aligned} \quad (2)$$

then attention output h_O is calculated by:

$$h_O = \text{Softmax}(h_Q X_K^T) X_V \quad (3)$$

where $h_Q = hW_Q$ is the output of the query layer. For ease of illustration, we ignore the FFN module and other parts of the inference workflow that are not addressed in our approach.

KV Quantization The B-bit KV quantization process during the prefill phase can be expressed as follows: First, determine the minimum number z_i and the maximum number m_i in G_i , where G_i is a group of number in X_K or X_V . Using these numbers, compute the quantized result $Q(G_i)$ for each group according to the formula:

$$Q(G_i) = \left\lfloor \frac{G_i - z_i}{s_i} \right\rfloor, \quad s_i = \frac{m_i - z_i}{2^B - 1} \quad (4)$$

The notation $\lfloor \cdot \rfloor$ represents rounding to the nearest integer. The results from all groups are aggregated to obtain $Q(X_K)$ and $Q(X_V)$. During the decoding phase, the quantized $Q(X_K)$ and $Q(X_V)$ and the stored quantization parameters z_i and s_i are used to recover the original values. In the decoding phase, the dequantized result X'_K, X'_V are used to calculate the attention output. X'_K, X'_V are obtained through aggregated G'_i for each G_i . G'_i can be computed using:

$$G'_i = Q(G_i) \cdot s_x + z_x \quad (5)$$

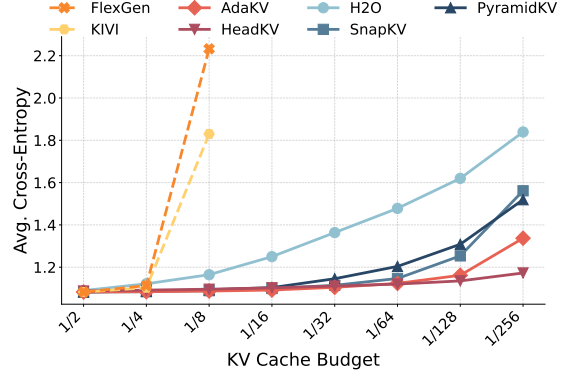


Figure 1: Cross-Entropy of Representative KV Compression methods under different KV Cache Budget on Gov_Report Dataset.

KV Pruning The goal of KV pruning is to find two submatrices $X_K^e, X_V^e \in \mathbb{R}^{b \times s \times d}$ from the full matrices X_K and X_V during the prefill phase, given a cache budget $s < n$, while maximizing performance preservation. During the decoding phase, LLMs with KV pruning only use X_K^e and X_V^e to update KV cache and generate new tokens.

4 Quantized Pruning

4.1 Information Bottleneck

From the perspective of the Information Bottleneck (IB) principle, we propose a unified theoretical framework to describe KV pruning and KV cache quantization. Previous work has leveraged the Information Bottleneck to formulate the objective of deep learning as an information-theoretic trade-off between representation compression and predictive information preservation (Tishby and Zaslavsky, 2015; Shwartz-Ziv and Tishby, 2017). In this work, we apply the perspective of IB on the KV cache. Given the input KV cache, denoted as X_k and X_v , and the model’s output Y generated using the KV cache, KV cache compression techniques, such as KV pruning and KV quantization, aim to minimize the following Lagrangian:

$$L = I(KV, KV_c) - \beta I(KV_c; Y) \quad (6)$$

where $I(KV; KV_c)$ represents the mutual information between the original KV cache (KV) and the compressed KV cache (KV_c), quantifying the information retained in the compressed representation KV_c . The term $I(KV_c; Y)$ quantifies the amount of information in KV_c that is pertinent to predicting the output distribution Y . β serves as a trade-off hyperparameter between representation

Pruning Method	<i>LongBench</i>											
	Token=128				Token=512				Token=2048			
	16-bit	8-bit	4-bit	2-bit	16-bit	8-bit	4-bit	2-bit	16-bit	8-bit	4-bit	2-bit
StreamingLLM	32.1	32.2	31.7	19.1	34.6	34.5	33.9	<u>20.7</u>	38.1	38.2	37.8	23.8
H2O	35.6	35.6	34.7	15.8	37.5	37.4	36.7	17.7	39.8	39.7	39.0	21.1
SnapKV	35.7	35.7	35.1	16.6	40.3	40.4	<u>39.7</u>	20.2	41.7	41.7	<u>41.0</u>	22.9
PyramidKV	37.4	37.3	36.4	<u>17.5</u>	40.3	40.3	39.6	20.9	<u>41.8</u>	<u>41.8</u>	41.3	<u>23.6</u>
Ada-KV	<u>39.3</u>	<u>39.2</u>	<u>37.4</u>	11.9	<u>40.9</u>	<u>40.8</u>	39.0	12.5	41.7	41.7	40.0	13.7
HeadKV	40.9	40.9	39.5	11.6	41.9	41.9	40.3	12.1	42.4	42.4	41.0	13.4

Pruning Method	<i>Needle-in-a-Haystack</i>											
	Token=128				Token=512				Token=2048			
	16-bit	8-bit	4-bit	2-bit	16-bit	8-bit	4-bit	2-bit	16-bit	8-bit	4-bit	2-bit
StreamingLLM	27.7	27.7	27.5	30.9	35.3	35.3	35.5	37.3	66.4	66.5	66.4	61.8
H2O	46.9	46.6	46.8	36.4	91.2	91.1	91.0	54.8	100	100	100	74.1
SnapKV	83.7	83.7	82.5	<u>55.9</u>	97.4	97.4	97.2	<u>66.3</u>	100	100	100	<u>78.1</u>
PyramidKV	98.9	98.9	98.8	67.5	100	100	100	78.6	100	100	100	79.6
Ada-KV	87.7	88.2	81.0	11.1	98.6	98.6	97.4	11.9	100	100	100	27.3
HeadKV	<u>98.6</u>	<u>98.5</u>	<u>98.4</u>	9.1	<u>100</u>	<u>99.9</u>	<u>99.7</u>	11.5	100	100	100	27.6

Table 1: Feasibility of quantized pruned tokens on LongBench and Needle-in-a-Haystack with Llama-3-8B-Instruct as backbone model. We use six different KV pruning methods to retain 128, 512 and 2048 tokens, and report the results of further quantization. KIVI is used as the default quantization method, except for Ada-KV and HeadKV, which adopt FlexGen. As they retain different tokens per head, making them difficult to be compatible with KIVI. The **best** results for each token-precision setting are in **bold**, the second best results are underlined.

compression and information preservation of the compressed KV cache.

With the Information Bottleneck, KV pruning and KV quantization can be understood as distinct yet complementary mechanisms for achieving this optimal compressed representation KV_c . KV pruning primarily influences the $I(KV; KV_c)$ term by determining the amount of information selected for transmission from the original KV . By discarding certain tokens, KV pruning reduces the complexity and raw information content that KV_c must encode from KV . KV quantization, conversely, primarily impacts the fidelity of this selected information within KV_c . It introduces noise or distortion due to reduced numerical precision, which can affect how much relevant information KV_c ultimately provides for predicting Y , thereby influencing $I(KV_c; Y)$.

However, an exclusive focus on optimizing either pruning or quantization in isolation may lead to a suboptimal trade-off between model performance and compression efficacy. To illustrate this, we evaluate standalone KV pruning and KV quantization methods, reporting their achieved KV cache compression ratios and the corresponding impact on model performance, as measured by cross-

entropy, more detailed can be seen in Appendix B. As demonstrated in Figure 1, both uni-dimensional approaches inherently encounter limitations in balancing information retention against the degree of compression; notably, cross-entropy escalates at an accelerated rate as the compression ratio increases.

Quantized Pruning Drawing upon the preceding analysis, we hypothesize that a unified strategy co-optimizing both pruning and quantization aspects can achieve a more favorable operating point within the information-theoretic trade-off space than optimizing either aspect in isolation. We propose Quantized Pruning based on this idea.

Specifically, Quantized Pruning first employs KV pruning methods to obtain the pruned key-value states \mathbf{X}_K^e and \mathbf{X}_V^e . Subsequently, these preserved states are quantized to $Q(\mathbf{X}_K^e)$ and $Q(\mathbf{X}_V^e)$ using KV quantization techniques during the prefill phase. In the decoding phase, the dequantized results derived from $Q(\mathbf{X}_K^e)$ and $Q(\mathbf{X}_V^e)$ are utilized to generate new tokens. The primary objective of Quantized Pruning is to maximize the model’s performance when utilizing the compressed KV cache, subject to a predefined KV Cache Budget B_{KV} for

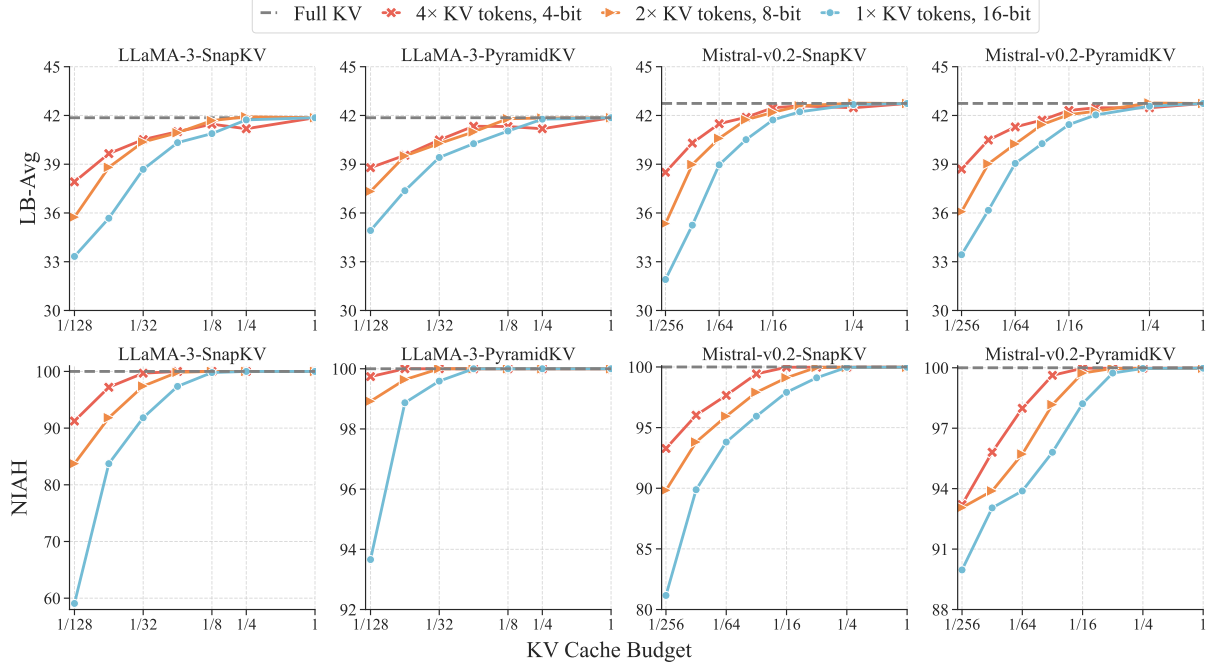


Figure 2: The token-precision trade-off under varying memory budgets on LongBench and NIAH. We report the results of SnapKV-based and PyramidKV-based quantized pruning on Llama-3 and Mistral-v0.2. We compare three configurations with approximately equivalent memory usage: 1) Using standalone KV pruning to retain $1\times$ tokens in 16-bit precision. 2) Quantized pruning by retaining $2\times$ tokens in 8-bit precision. 3) Quantized pruning by retaining $4\times$ tokens in 4-bit precision. Quantized pruning, which preserves more tokens at a lower precision, consistently outperforms standalone KV pruning methods across various budgets.

the KV cache.

$$\frac{\text{Memory}(Q(\mathbf{X}_K^e), Q(\mathbf{X}_V^e))}{\text{Memory}(\mathbf{X}_K, \mathbf{X}_V)} \leq B_{KV} \quad (7)$$

Given that the memory budget for Quantized Pruning is jointly determined by the number of preserved tokens and their average bit-width, Quantized Pruning is equivalent to identifying an optimal token-precision trade-off under fixed memory budget. This optimization aims to achieve a more favorable operating point within the Information Bottleneck.

5 Experimental Setup

Benchmarks We evaluate the performance of quantized pruning using the LongBench benchmark (Bai et al., 2024) and Needle-in-a-Haystack test (Kamradt, 2023). To better distinguish the performance impacts related to input lengths and layer-wise configurations (Sections 7.1 and C), we further utilize RULER (Hsieh et al., 2024), a dataset with different input length and diverse types of needles across 4 task categories. Detailed information on these datasets can be found in Appendix A.


LLMs In our primary evaluations, We employ state-of-the-art open-weight LLMs, specifically Llama-3-8B-Instruct (Dubey et al., 2024) and Mistral-7B-Instruct-v0.2 (Jiang et al., 2023). For scaling experiments in Section 7.2, we additionally assess the performance of Llama-3.2-1B, Llama-3.2-3B (Dubey et al., 2024) and Llama-3-70B (Dubey et al., 2024).

Setup Our experiments are designed to comprehensively investigate the optimal token-precision trade-off in KV cache compression. We quantify memory budget by reporting the compression ratio of the KV cache relative to the full, uncompressed KV cache. For KV pruning, we utilize PyramidKV (Cai et al., 2024) and SnapKV (Li et al., 2024) as representative state-of-the-art methods. To assess the feasibility of quantized pruning (Section 6), we also include StreamingLLM (Xiao et al., 2023), H2O (Zhang et al., 2024b), Ada-KV (Feng et al., 2024), and HeadKV (Fu et al., 2024). For KV quantization, we adopt KIVI (Liu et al., 2024c) as our primary method due to its stability and broad compatibility. Furthermore, in Section 7.3, we evaluate quantization strategies from FlexGen (Sheng et al., 2023) and KVQuant (Hooper et al., 2024) for

a comprehensive comparative analysis. Additional setup details are provided in Appendix B.

6 Optimal Token-Precision Trade-Off

In this section, we aim to find the optimal token-precision trade-off in KV cache compression. We first examine the feasibility of integrating KV pruning and quantization(Q1). Subsequently, we explore the best optimal allocation strategy between precision and token under varying memory budgets(Q2).


 Q1. Is it feasible to integrate KV pruning and quantization for a lower compression rate?

We first evaluate the feasibility of integrating KV pruning and quantization as a prerequisite for exploring the token-precision trade-off. We employ Llama-3-8B-Instruct and evaluate a range of KV pruning methods on the LongBench and NIAH. We report the results of quantizing the remaining tokens to different precision levels after applying KV pruning.

From Table 1, we observe that *it is feasible to quantize pruned KV cache for a lower compression rate*. For most KV pruning methods we evaluate, further quantizing the preserved tokens to as low as 4-bit precision results in minimal performance degradation, while quantizing to 8-bit precision shows negligible impact. However, reducing precision to 2-bit leads to a drastic performance decline across most KV pruning methods. This observation holds consistently across different KV pruning methods and varying numbers of preserved tokens.

Compared with precision, reducing the number of preserved tokens leads to more significant performance degradation. Specifically, when the number of preserved tokens is reduced to 1/4 (from 2048 to 512), all KV pruning methods experience a noticeable performance drop. In contrast, when the precision is reduced to 1/4 (from 16-bit to 4-bit), which has the same memory budget as token dimension, the performance degradation is relatively mild. This suggests that, under the same memory budget, tokens might have a more significant impact on the results compared to precision.

To conclude, KV pruning can be effectively integrated with KV quantization at a precision level of 4-bit without substantial performance degradation.

 Q2. What is the optimal allocation strategy between precision and token under varying memory budgets?

Observing that KV pruning and KV quantization can be effectively integrated, we further investigate that, given a fixed memory budget, how to balance the trade-off between number of preserved tokens and precision to achieve optimal performance. To this end, we evaluate the performance of quantized pruning using two leading KV pruning methods, SnapKV and PyramidKV, across different memory budgets on LongBench and NIAH.

As shown in Figure 2, we observe that *quantized pruning, which preserves more tokens at a lower precision, consistently outperforms standalone KV pruning methods across various budgets*. Besides, the 1/4 KV cache budget of quantized pruning on 8-bit 2x KV tokens outperforms 4-bit 4x tokens which represents the standalone KV quantization methods. This conclusion demonstrate that **quantized pruning can achieve a more favorable operating point within the Information Bottleneck than optimizing either aspect in isolation**.

For the NIAH task, the improvements from quantized pruning are particularly pronounced. This may be attributed to that quantized pruning can cover more tokens for retrieval under the same memory budget compared to standalone KV pruning.

In high-budget scenarios, the 8-bit strategy tends to deliver slightly better performance, which may due to the number of tokens at this budget is already quite large. In low-budget scenarios, such as 1/128 KV cache budget, storing more tokens at 4-bit precision yields superior results, highlighting the importance of token coverage when resources are constrained. Overall, using lower precision to preserve more tokens under a limited budget results in notable performance gains, compared to standalone KV pruning methods that use full precision to store fewer tokens.

7 Further Analysis

In this section, we further investigate series of key aspects regarding token-precision trade-off, including the impact of quantized pruning on various downstream task types and input lengths, model scaling effect and ablation on quantization strategies. We also explore fine-grained layer-wise quantized pruning in the Appendix C.

Models	Token	Bit	Task Types						
			SQA	MQA	SUMM	Fewshot	Syn.	Code	RULER-8k
Llama-3-8B-Instruct	512	16	28.2	31.9	23.5	67.6	37.7	57.6	67.5
	1024	8	29.6	33.1	24.3	67.9	37.4	58.0	74.9
	2048	4	30.7	32.5	25.3	68.8	37.2	57.6	82.2
Mistral-7B-Instruct-v0.2	512	16	33.7	27.3	24.3	65.6	41.75	54.0	53.1
	1024	8	34.2	29.0	25.6	66.4	43.73	54.8	62.1
	2048	4	35.2	28.14	26.6	66.9	43.08	55.4	73.6

Table 2: The token-precision trade-off in different task types. We report the results of 6 task types in LongBench and 8k subset of RULER. We use PyramidKV-based quantized pruning. The **best** results for each task type are in **bold**.

7.1 Impact on Task Types and Input Lengths

Task Types To further investigate the token-precision trade-off in different task types, we evaluate PyramidKV-based quantized pruning on six task types from LongBench and the 8K subset of the RULER dataset. We use 1/16 KV cache budget, and explore the token-precision trade-off under this fixed memory budget, as this setting exhibits minimal performance differences across three settings, facilitating a more direct comparison of performance across task types.

As illustrated in Table 2, we observe that the performance of quantized pruning is remarkably consistent across different task types. Specifically, lower precision, which retains more tokens in KV cache, leads to substantial performance improvements in the RULER task, which heavily relies on retrieving contents from the input. Tasks with high retrieval demands, such as Summarization and Single-Doc QA, also show noticeable gains with quantized pruning, particularly when $4\times$ tokens are preserved at 4-bit precision.

For tasks requiring more reasoning rather than intensive retrieval, such as Code Completion, Synthetic and Multi-Doc QA, the benefits of trading precision for more tokens are less pronounced. In these cases, storing fewer tokens with higher precision generally performs better. For example, using 1024 tokens in 8-bit precision achieves the highest score of 58 in Code task.

Input Lengths To evaluate the token-precision trade-off across various input lengths, we conduct experiments on subsets with different input length of the RULER dataset. Additionally, we analyze LongBench by grouping its data based on input length, more detailed information can be found in Appendix A. The results are shown in Figure 3.

Our observations are as follows: *quantized pruning*

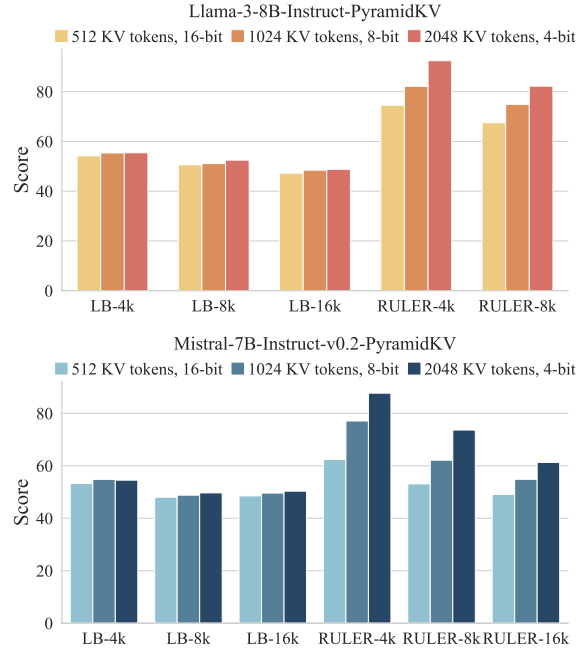


Figure 3: The token-precision trade-off in different input lengths. We report the results of LongBench and three subsets of RULER. We use PyramidKV-based quantized pruning.

ing consistently outperforms standalone KV eviction methods across various input lengths, regardless of the models and task types. Within the same dataset, scores decrease as input length increases; however, the relative differences among different compression methods remain similar across varying input lengths. Moreover, quantized pruning achieves significant performance improvements across all input lengths for retrieval demanded tasks like RULER.

7.2 Scaling Effect on Quantized Pruning

To investigate the impact of model scaling on quantized pruning, we conducted experiments on four models from the Llama series: Llama3-70B,

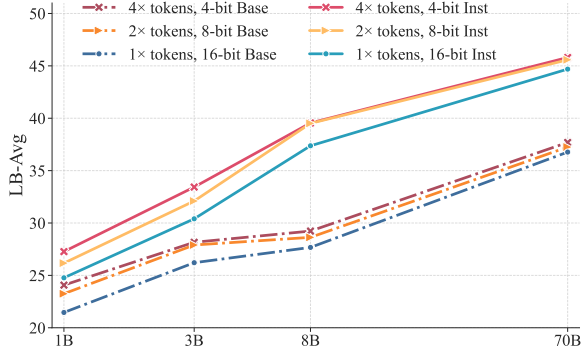


Figure 4: Scaling effect on Llama family models, with PyramidKV-based quantized pruning. All models are under 1/64 KV cache budget.

Llama3-8B, Llama3.2-3B and Llama3.2-1B. Since the Llama3 series does not include 1B and 3B versions, we used the Llama3.2 series as substitutes. However, it is important to note that the performance of Llama3.2-3B-Base is quite similar to Llama3-8B-Base. For both the Base models and Instruct models, we evaluated their performance on LongBench under 1/64 KV cache budgets. To further validate the conclusion, we also experimented with a 1/16 KV cache budget, and the results are presented in the Appendix C.

As shown in Figure 4, we observe that quantized pruning consistently achieves better performance across all scaling levels. Notably, the performance gap between quantized pruning and standalone KV pruning methods remains relatively stable across different model scales. For Base models, although the performance improvement from scaling is smaller compared to Instruct models, quantized pruning still provides a noticeable performance boost.

These findings highlighting the robustness and effectiveness of quantized pruning across model scaling.

7.3 Ablation on Quantization Strategies

While there has been extensive research on strategies for KV quantization, it remains unclear whether existing quantization strategies remain effective when integrated with KV pruning methods. In this section, we aim to investigate the impact of KV quantization strategies on quantized pruning. We also conducted an analysis of the effect of quantization group size on the quantized pruning, with further results available in Appendix C.

Quantization methods We explore the methods in FlexGen (Sheng et al., 2023), KIVI (Liu et al.,

2024c), and KVQuant (Hooper et al., 2024). To elaborate, for the FlexGen methods, KV quantization is applied to both the key and value caches along the token dimension, grouping every 64 elements without filtering outlier numbers. We modify the FlexGen by (1) filtering 1% of outlier numbers in both the key and value cache, as mentioned in KVQuant (2) quantizing the key along the channel dimension, as in KIVI and (3) combining (1) and (2). These correspond to the results labeled as FlexGen+Outlier 1%, KIVI, and KIVI+Outlier 1% in the Figure 10.

We can observe that none of the quantization strategies show significant performance degradation when combined with KV pruning methods, demonstrating the relative stability of quantized pruning. The KIVI method consistently outperforms FlexGen across various models and KV pruning methods. The improvement is particularly pronounced for PyramidKV on the Mistral model, underscoring the significance of quantizing key states along the channel dimension. Filtering 1% of outlier numbers proves effective for the FlexGen strategy but yields limited improvements for KIVI. It shows some benefit on Llama models but results in negligible gains on the Mistral models.

Overall, KIVI demonstrates strong performance when integrated with KV pruning methods, while other KV quantization strategies also maintain good results, highlighting the stability of quantized pruning.

8 Conclusion

We investigate a series of critical yet unexplored questions regarding the effectiveness and feasibility of token-precision trade-off in KV cache compression. Through comprehensive experiments, we demonstrate that storing more tokens in the KV cache with lower precision can significantly enhance the long-context performance of LLMs, and demonstrating robust performance across diverse input lengths, downstream tasks, with particularly significant gains in retrieval tasks. Moreover, we show that quantized pruning demonstrates strong feasibility across different KV pruning methods, quantization strategies, and model scales. Our analysis sheds light on the token-precision trade-off of KV cache memory optimization, offering valuable insights into designing more efficient compression strategies. We hope this work deepens the understanding of KV cache compression and inspires future research.

Limitations

While our work demonstrates the effectiveness of KV compression through trade-offs between token and precision dimensions, other potential dimensions, such as head and layer, remain unexplored. Investigating the feasibility of combining these dimensions with token and precision for a more substantial compression potential represents an avenue for future research.

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*.
- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit Sanghai. 2023. [GQA: Training generalized multi-query transformer models from multi-head checkpoints](#). In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- 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, multi-task 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.
- William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan Kelly. 2024. [Reducing transformer key-value cache size with cross-layer attention](#). *Preprint*, arXiv:2405.12981.
- Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie Hu, and Wen Xiao. 2024. [Pyramidkv: Dynamic kv cache compression based on pyramidal information funneling](#). *Preprint*, arXiv:2406.02069.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011*.
- Alessio Devoto, Yu Zhao, Simone Scardapane, and Pasquale Minervini. 2024. [A simple and effective \$l_2\$ norm-based strategy for kv cache compression](#). *Preprint*, arXiv:2406.11430.
- Shichen Dong, Wen Cheng, Jiayu Qin, and Wei Wang. 2024. [Qaq: Quality adaptive quantization for llm kv cache](#). *Preprint*, arXiv:2403.04643.

- Haojie Duanmu, Zhihang Yuan, Xiuhong Li, Jiangfei Duan, Xingcheng Zhang, and Dahua Lin. 2024. [Skvq: Sliding-window key and value cache quantization for large language models](#). *arXiv preprint arXiv:2405.06219*.
- 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*.
- Razvan-Gabriel Dumitru, Paul-Ioan Clotan, Vikas Yadav, Darius Peteleaza, and Mihai Surdeanu. 2024. Change is the only constant: Dynamic llm slicing based on layer redundancy. *arXiv preprint arXiv:2411.03513*.
- Alexander R Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir R Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. *arXiv preprint arXiv:1906.01749*.
- 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](#). *Preprint*, arXiv:2407.11550.
- Yao Fu. 2024. [Challenges in deploying long-context transformers: A theoretical peak performance analysis](#). *Preprint*, arXiv:2405.08944.
- Yu Fu, Zefan Cai, Abdelkadir Asi, Wayne Xiong, Yue Dong, and Wen Xiao. 2024. Not all heads matter: A head-level kv cache compression method with integrated retrieval and reasoning. *arXiv preprint arXiv:2410.19258*.
- Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2024. [Model tells you what to discard: Adaptive KV cache compression for LLMs](#). In *The Twelfth International Conference on Learning Representations*.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. *arXiv preprint arXiv:1911.12237*.
- Daya Guo, Canwen Xu, Nan Duan, Jian Yin, and Julian McAuley. 2023. Longcoder: A long-range pre-trained language model for code completion. In *International Conference on Machine Learning*, pages 12098–12107. PMLR.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W. Mahoney, Yakun Sophia Shao, Kurt Keutzer, and Amir Gholami. 2024. [Kvquant: Towards 10 million context length llm inference with kv cache quantization](#). *Preprint*, arXiv:2401.18079.

692	Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shan-	Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong,	746
693	tanu Acharya, Dima Rekesh, Fei Jia, and Boris Gins-	Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and	747
694	burg. 2024. RULER: What’s the real context size of	Xia Hu. 2024c. KIVI: A tuning-free asymmetric 2bit	748
695	your long-context language models? In <i>First Confer-</i>	quantization for KV cache . In <i>Proceedings of the</i>	749
696	<i>ence on Language Modeling</i> .	<i>41st International Conference on Machine Learning</i> ,	750
		volume 235 of <i>Proceedings of Machine Learning</i>	751
697	Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng	<i>Research</i> , pages 32332–32344. PMLR.	752
698	Ji, and Lu Wang. 2021. Efficient attentions for		
699	long document summarization. <i>arXiv preprint</i>	Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang,	753
700	<i>arXiv:2104.02112</i> .	Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng	754
		Chen. 2024. Shortgpt: Layers in large language	755
701	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	models are more redundant than you expect. <i>arXiv</i>	756
702	sch, Chris Bamford, Devendra Singh Chaplot, Diego	<i>preprint arXiv:2403.03853</i> .	757
703	de las Casas, Florian Bressand, Gianna Lengyel, Guil-		
704	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi,	758
705	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	and Roy Schwartz. 2024. Transformers are multi-	759
706	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	state rnns . <i>Preprint</i> , arXiv:2401.06104.	760
707	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,		
708	arXiv:2310.06825.	Reiner Pope, Sholto Douglas, Aakanksha Chowdhery,	761
		Jacob Devlin, James Bradbury, Jonathan Heek, Kefan	762
709	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	Xiao, Shivani Agrawal, and Jeff Dean. 2023. Effi-	763
710	Zettlemoyer. 2017. Triviaqa: A large scale distantly	ciently scaling transformer inference . In <i>Proceedings</i>	764
711	supervised challenge dataset for reading comprehen-	<i>of Machine Learning and Systems</i> , volume 5, pages	765
712	sion. <i>arXiv preprint arXiv:1705.03551</i> .	606–624. Curran.	766
713	Greg Kamradt. 2023. Needle in a haystack - pressure	Machel Reid, Nikolay Savinov, Denis Teplyashin,	767
714	testing llms. https://github.com/gkamradt/	Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste	768
715	LLMTest_NeedleInAHaystack .	Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Fi-	769
		rat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Un-	770
716	Tom��� Ko��isk��y, Jonathan Schwarz, Phil Blunsom, Chris	locking multimodal understanding across millions of	771
717	Dyer, Karl Moritz Hermann, G��bor Melis, and Ed-	tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	772
718	ward Grefenstette. 2018. The narrativeqa reading		
719	comprehension challenge. <i>Transactions of the Asso-</i>	Siyu Ren and Kenny Q. Zhu. 2024. On the effi-	773
720	<i>ciation for Computational Linguistics</i> , 6:317–328.	cacy of eviction policy for key-value constrained	774
		generative language model inference . <i>Preprint</i> ,	775
721	Xin Li and Dan Roth. 2002. Learning question clas-	arXiv:2402.06262.	776
722	sifiers. In <i>COLING 2002: The 19th International</i>		
723	<i>Conference on Computational Linguistics</i> .	Noam Shazeer. 2019. Fast transformer decod-	777
		ing: One write-head is all you need . <i>Preprint</i> ,	778
724	Yuhong Li, Yingbing Huang, Bowen Yang, Bharat	arXiv:1911.02150.	779
725	Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai,		
726	Patrick Lewis, and Deming Chen. 2024. Snapkv:	Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuo-	780
727	Llm knows what you are looking for before genera-	han Li, Max Ryabinin, Daniel Y. Fu, Zhiqiang Xie,	781
728	tion . <i>Preprint</i> , arXiv:2404.14469.	Beidi Chen, Clark Barrett, Joseph E. Gonzalez, Percy	782
		Liang, Christopher R��, Ion Stoica, and Ce Zhang.	783
729	Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang,	2023. Flexgen: High-throughput generative infer-	784
730	Bo Liu, Chenggang Zhao, Chengqi Deng, Chong	ence of large language models with a single gpu .	785
731	Ruan, Damai Dai, Daya Guo, et al. 2024a.	<i>Preprint</i> , arXiv:2303.06865.	786
732	Deepseek-v2: A strong, economical, and efficient		
733	mixture-of-experts language model. <i>arXiv preprint</i>	Ravid Shwartz-Ziv and Naftali Tishby. 2017. Opening	787
734	<i>arXiv:2405.04434</i> .	the black box of deep neural networks via informa-	788
		tion . <i>Preprint</i> , arXiv:1703.00810.	789
735	Tianyang Liu, Canwen Xu, and Julian McAuley.		
736	2023. Repobench: Benchmarking repository-level	Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wenhui	790
737	code auto-completion systems. <i>arXiv preprint</i>	Wang, Shuming Ma, Quanlu Zhang, Jianyong Wang,	791
738	<i>arXiv:2306.03091</i> .	and Furu Wei. 2024. You only cache once: Decoder-	792
		decoder architectures for language models . <i>ArXiv</i> ,	793
739	Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao	abs/2405.05254.	794
740	Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyril-		
741	lidis, and Anshumali Shrivastava. 2024b. Scis-	Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Shikuan	795
742	sorhands: Exploiting the persistence of importance	Hong, Yiwu Yao, and Gongyi Wang. 2024. Razo-	796
743	hypothesis for llm kv cache compression at test time.	rattention: Efficient kv cache compression through	797
744	<i>Advances in Neural Information Processing Systems</i> ,	retrieval heads. <i>arXiv preprint arXiv:2407.15891</i> .	798
745	36.		

799	Naftali Tishby and Noga Zaslavsky. 2015. Deep learning and the information bottleneck principle . In <i>2015 IEEE Information Theory Workshop (ITW)</i> , pages 1–5.	853
800		854
801		855
802		856
803	Daniel Waddington, Juan Colmenares, Jilong Kuang, and Fengguang Song. 2013. Kv-cache: A scalable high-performance web-object cache for manycore . In <i>2013 IEEE/ACM 6th International Conference on Utility and Cloud Computing</i> , pages 123–130.	857
804		858
805		
806		
807		
808	Haoyi Wu and Kewei Tu. 2024. Layer-condensed kv cache for efficient inference of large language models . <i>Preprint</i> , arXiv:2405.10637.	859
809		860
810		861
811	Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, Junxian Guo, Shang Yang, Haotian Tang, Yao Fu, and Song Han. 2024. Duoattention: Efficient long-context llm inference with retrieval and streaming heads. <i>arXiv preprint arXiv:2410.10819</i> .	862
812		863
813		
814		
815		
816	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2023. Efficient streaming language models with attention sinks. <i>arXiv</i> .	864
817		
818		
819	Dongjie Yang, XiaoDong Han, Yan Gao, Yao Hu, Shilin Zhang, and Hai Zhao. 2024a. Pyramidinfer: Pyramid kv cache compression for high-throughput llm inference . <i>Preprint</i> , arXiv:2405.12532.	865
820		866
821		867
822		868
823		869
824		870
825		871
826		872
827		873
828		874
829		875
830		876
831		877
832		878
833		879
834		880
835		881
836		882
837		883
838		884
839		885
840		886
841		887
842		888
843		889
844		890
845		
846		
847		
848		
849		
850		
851		
852		
	Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, et al. 2021. Qmsum: A new benchmark for query-based multi-domain meeting summarization. <i>arXiv preprint arXiv:2104.05938</i> .	
	Zayd Muhammad Kawakibi Zuhri, Muhammad Farid Adilazuarda, Ayu Purwarianti, and Alham Fikri Aji. 2024. Mlkv: Multi-layer key-value heads for memory efficient transformer decoding . <i>Preprint</i> , arXiv:2406.09297.	
	A Datasets	
	LongBench LongBench (Bai et al., 2024) includes 17 datasets covering 6 categories of tasks, which can be divided into single-document QA (Dasigi et al., 2021; Kočiský et al., 2018), multi-document QA (Yang et al., 2018; Ho et al., 2020), summarization (Huang et al., 2021; Fab-bri et al., 2019; Zhong et al., 2021), few-shot learning (Gliwa et al., 2019; Joshi et al., 2017; Li and Roth, 2002), synthetic, and code generation (Guo et al., 2023; Liu et al., 2023). LongBench features an average input length ranging from 1,235 to 18,409 tokens. For inputs exceeding the model’s context window length(8k for Llama-3-8B-Instruct (Dubey et al., 2024), we split the data and only take the beginning and end segments of the input to fill the context window length. Additionally, we reserve sufficient space for newly generated tokens based on the specific type of sub-dataset. For evaluate the impact of input lengths, we select datasets with sufficient data to cover three input length ranges: (<4k, 4k 8k, and >8k). These datasets include MultiFieldQA-en, 2Wiki-MultihopQA, GovReport, TREC, TriviaQA, SAM-Sum, and RepoBench-P, representing a variety of task types. We refer to the three subsets as LB-4k, LB-8k, and LB-16k, respectively.	
	NIAH Needle-in-a-Haystack(NIAH) (Kamradt, 2023) is a challenging pressure test designed to assess the ability of models to accurate identify and retrieve relevant information from lengthy context. NIAH randomly inserts key information into an arbitrary position within a long essay. In our setup, we use PaulGrahamEssays as the haystack and the sentence "The best thing to do in San Francisco is eat a sandwich and sit in Dolores Park on a sunny day." as the needle, which is the default setting of NIAH. We vary the essay length from 1,000 tokens up to the models’ context window limits, increasing by 100 tokens per step for Llama-series models and 400 tokens per step for Mistral. The	

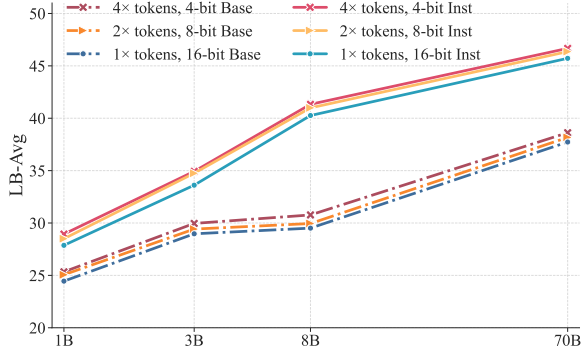


Figure 5: Scaling effect on Llama family models, with PyramidKV-based quantized pruning. All models are under 1/16 KV cache budget.

results are reported as the average score across all tests.

RULER RULER (Hsieh et al., 2024) generates synthetic examples to evaluate long-context language models with configurable sequence lengths and varying task complexities. It includes four task categories: Retrieval, Multi-hop Tracing, Aggregation, and Question Answering. The dataset comprises six subsets with input lengths of 4K, 8K, 16K, 32K, 64K and 128K tokens. In our experiments, we use the 4K, 8K and 16K subsets to test the models within their context window limits.

B Experiment Setup

Memory Budgets We report the ratio of compressed KV cache and the full KV cache for memory budge. The full KV cache for Llama-3 is 8k KV tokens in 16-bit on both LongBench and NIAH, while for Mistral-v0.2 is 16k KV tokens on LongBench and 32k KV tokens on NIAH in 16-bit.

KV eviction methods We retain the last 32 tokens for StreamingLLM (Xiao et al., 2023), H2O (Zhang et al., 2024b), and SnapKV (Li et al., 2024), while keeping 8 tokens for PyramidKV (Cai et al., 2024), Ada-KV (Feng et al., 2024) and HeadKV (Fu et al., 2024), as recommended in the corresponding paper (Cai et al., 2024; Fu et al., 2024). For other settings, we adopt the default configurations from their papers.

KV quantization We utilize HQQQuantized-Cache from Huggingface and adjust the group dimensions of keys and values to implement grouped quantization strategies from FlexGen (Sheng et al., 2023) and KIVI (Liu et al., 2024c). We use 64 as the default group size which is suggested in

Model	Method	Group Size		
		32	64	128
Llama-3	SnapKV	40.4	39.6	38.9
	PyramidKV	40.3	39.6	38.9
Mistral-v0.2	SnapKV	40.4	40.3	40.1
	PyramidKV	40.3	40.5	40.0

Table 3: The impact of group size for quantized pruning on LongBench.

FlexGen (Sheng et al., 2023). In the experiments involving outlier filtering, we exclude numbers in the KV cache with a absolute value exceeding 6 from quantization, which roughly corresponds to the top 1% of outliers based on our validation set analysis.

C More results in Experiments

Group Size We analyze the impact of group size during KV quantization. We employ SnapKV and PyramidKV to retain 512 tokens with 4-bit KIVI quantization and report the performance variations when the group sizes were set to 32, 64 and 128. As shown in Table 3, smaller group sizes lead to performance improvements at the cost of higher memory usage. Reducing the group size from 128 to 64 resulted in a notable improvement, but further decreasing it from 64 to 32 yielded minimal gains for the Mistral model. Therefore, we set the default quantization group size to 64 to balance performance and memory usage in our experiments.

Scaling Effect on Quantized Pruning We validate our findings by measuring performance at a budget of 1/16, with the results shown in Figure 5. This corroborates the conclusion reported in the main text using a budget of 1/64. When the KV cache budget is relative small to 1/64, the performance improvement brought by quantized pruning is higher compared to 1/16 KV cache budget, which aligns with the conclusions we observed earlier in Q2.

Layer-Wise Quantized Pruning Inspired by the observation that different layers may have varying requirements for the number of tokens in PyramidKV (Cai et al., 2024) and PyramidInfer (Yang et al., 2024a), we further investigate whether the demands for precision and preserved tokens are consistent across layers. To explore this, we use

Token-Precision Settings	Throughput (tokens/s)			KV Cache Memory per Batch Item	LongBench Avg	NIAH Avg
	batch=64	batch=96	batch=128			
128 tokens, 16-bit	634.9	1460.1	1658.8	~18 MB	35.24	81.16
256 tokens, 8-bit	575.9	1250.8	1389.5	~19.1 MB	38.98	89.83
512 tokens, 4-bit	554.8	1317.6	1482.8	~20.2 MB	40.3	93.28

Table 4: Speed Evaluation on Mistral-7B-v0.2 with SnapKV-based pruning and KIVI-based quantization. We report the decoding speed of different methods on single NVIDIA A100 GPU.

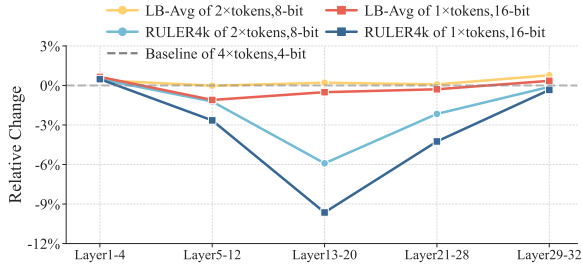


Figure 6: The results of Layer-Wise Quantized Pruning on Llama-3-8B-Instruct, with SnapKV as pruning method. KV cache budget=1/16.

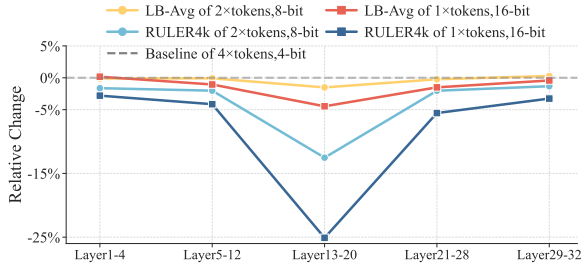


Figure 7: The results of Layer-Wise Quantized Pruning on Mistral-7B-v0.2-Instruct, with SnapKV as pruning method. KV cache budget=1/64.

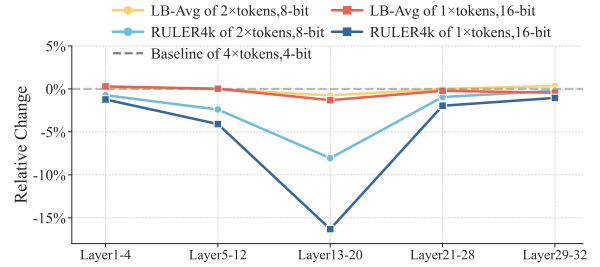


Figure 8: The results of Layer-Wise Quantized Pruning on Mistral-7B-v0.2-Instruct, with SnapKV as pruning method. KV cache budget=1/16.

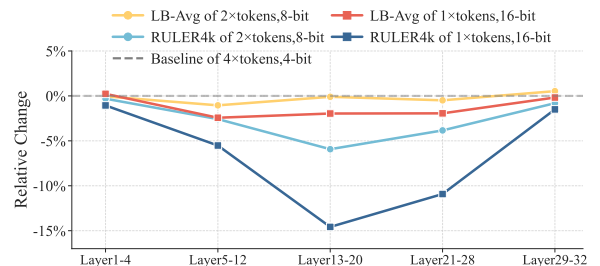


Figure 9: The results of layer-wise quantized pruning on Llama-3-8B-Instruct, with SnapKV as pruning method. We use $4 \times$ KV token 4-bit as baseline and report the relative change. Configurations are modified every 4 layers for the initial and final layers, while intermediate layers are reconfigured every 8 layers.

the best-performing configuration from previous experiments, 4-bit precision with $4 \times$ tokens, as the baseline and compare it against layer-wise configurations adopting 8-bit precision with $2 \times$ tokens and 16-bit precision with $1 \times$ tokens. Using SnapKV as the KV pruning method, we present the results for Llama-3 under 1/64 KV cache budget in the Figure 9. The x-axis shows the modified layers range, while the y-axis shows the relative change to the baseline (4-bit precision with $4 \times$ tokens) on LongBench and RULER-4k. Furthermore, we present the results for Llama-3 evaluated with a 1/64 KV cache budget and Mistral-v0.2 evaluated with 1/64 and 1/16 KV cache budgets.

It is evident that for most layers, transitioning from $4 \times$ tokens with 4-bit precision to higher precision and fewer tokens results in a performance decline under constrained KV cache budgets. Specif-

ically, the shift to 8-bit shows a relatively minor performance drop, whereas moving to 16-bit with fewer preserved tokens leads to a more significant decrease. These layers-wise trade-off conclusions are consistent with our experiments before.

Notably, modifying intermediate layers causes a drastic performance decline, while adjustments made at the initial and final layers result in comparatively smaller performance reductions. This effect is especially pronounced in retrieval-related tasks such as RULER-4k, where significant performance differences are observed. On LongBench, changes are less significant, with a notable performance drop only observed at 16-bit precision. These findings highlight that, under the same memory budget,

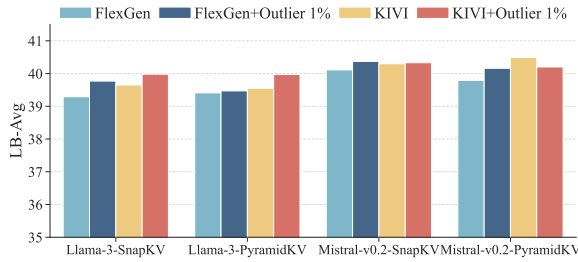


Figure 10: Ablation of quantization strategies on quantized pruning, remaining 512 KV tokens in 4-bit.

these configurations is comparable (due to similar memory bandwidth requirements), the choice of token-precision trade-off has a substantial impact on model quality.

preserving more tokens with lower precision in intermediate layers is crucial for the performance, while the token-precision trade-off in the initial and final layers exerts a more balanced influence.

For Mistral-v0.2 in Figure 7 and Figure 8, we can see the layer-wise results are similar to Llama-3. Modifying intermediate layers causes a drastic performance decline, while adjustments made at the initial and final layers result in comparatively smaller performance reductions. This effect is especially pronounced in retrieval-related tasks such as RULER-4k, where significant performance differences are observed. On LongBench, changes are less significant, with a notable performance drop only observed at 16-bit precision.

Speed Evaluation Both KV Pruning and KV Quantization primarily aim to reduce the memory footprint of the KV cache. During the autoregressive decoding phase, performance is often bandwidth-bound, meaning the latency is highly correlated with the amount of data that needs to be read from and written to memory (i.e., the size of the KV cache). Therefore, methods achieving similar memory compression levels are expected to yield similar throughput.

To provide concrete data on this, we conducted additional experiments measuring throughput alongside accuracy for three configurations using Mistral-7B-v0.2 (with SnapKV pruning and KIVI-based quantization). The configurations maintain nearly identical memory usage while varying the token-precision balance:

As shown in the table 4, the KV cache memory consumption and the resulting throughput are very similar across these three settings. However, the task performance (LongBench Avg, NIAH Avg) shows marked differences, with the quantized pruning approach (512 tokens with 4-bit precision) achieving significantly higher accuracy. This demonstrates that while the speed impact of