# INTELLLM: LITTLE HINTS MAKE A BIG DIFFERENCE FOR LLM KV CACHE COMPRESSION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Large Language Models (LLMs) have demonstrated exceptional capabilities in integrating contextual knowledge, but their deployment is often constrained by the substantial computational resources required for long text sequences. To mitigate the inference time cost associated with attention mechanisms, LLMs utilize key-value embedding caching techniques (KV cache), which introduce significant storage pressure. In this paper, we propose IntelLLM, a novel and efficient approach to KV cache compression that strikes a balance between compression rate and performance. Drawing inspiration from sparse attention mechanism, we observe that only a small subset of tokens in lengthy texts capture the majority of attention weights. This sparsity, intrinsic to the attention mechanism, serves as the foundation for improving the KV compression ratio through a strategic eviction method. IntelLLM is composed of center of gravity eviction (CGE) strategy and remote gap localization (RGL) strategy. CGE is designed to address the potential loss of important semantic dependencies when evicting high-sparsity tokens, which prioritizes the retention of key tokens by shielding the center of gravity of attention during inference, thereby preserving critical information and optimizing the efficiency of attention computation. Additionally, RGL is proposed to leverage implicit positional features to maintain long-range dependencies, inspired by advancements in location encoding research. Our KV compression approach integrates seamlessly with existing LLMs, requiring minimal code modifications without the need for fine-tuning or model parameter changes. IntelLLM not only significantly reduces the storage requirements for KV cache but also consistently outperforms full KV models in long text processing tasks, while utilizing only 50% of the typical KV cache expenses.

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

#### 1 INTRODUCTION

Large Language Models (LLMs) have rapidly expanded from academic research to industrial appli-037 cations Wu et al. (2023)Zhao et al. (2024a) due to their exceptional capabilities in language understanding and generation, knowledge retention and retrieval, and multi-task learning Ge et al. (2024). However, as LLM applications evolve, the associated computational costs have become a critical 040 challenge. During autoregressive inference in LLMs, the quadratic complexity of the transformer 041 architecture and the large number of model parameters result in significant computational overhead 042 Zhou et al. (2024). To avoid recomputation, KV caches are preserved. However, the size of the 043 KV cache grows proportionally with the length of the inference sequence, eventually surpassing the 044 model's parameter size as the sequence accumulates, with the size of the KV cache emerging as 045 a key issue. This raises a significant challenge for LLM deployment: how to reduce the memory consumption of the KV cache without compromising model accuracy. 046

To tackle this challenge, Hooper et al. (2024b)Ye et al.Zhang et al. (b) have proposed various strategies to optimize KV caching and reduce memory consumption during inference. One prominent approach is cache compression Liu et al. (2024), which minimizes storage requirements by eliminating redundant information in the KV cache. Low-rank approximation techniques Zhang et al.
(a) have been employed to compress the original high-dimensional KV representations into lowerdimensional spaces. Additionally, sparse attention mechanisms Liu et al. (2023)Fu et al. have been introduced to selectively retain only the most critical key-value pairs, discarding less relevant information and reducing both computation and cache space requirements. However, these methods 054 present two main limitations: firstly, there is a trade-off between memory savings and accuracy, mak-055 ing it difficult to strike the right balance; secondly, the added complexity of these strategies often 056 leads to a conflict between increased computational overhead and memory optimization. Mean-057 while, the windowed attention approach Beltagy et al. (2020) has been proposed as a solution to 058 reduce KV caching costs during inference. While this approach effectively reduces memory consumption, it suffers from significant accuracy degradation when handling sequences that exceed the predefined memory window size. Comparable to the fixed window design, numerous studies have 060 explored reverse thinking, which involve adapting long texts to a preset window size. This can be 061 viewed as a hierarchical strategy that condenses the text, extracting only the core information to al-062 leviate KV cache pressure Song et al. (2024)Shao et al. (2024)He et al. (2024). These methods offer 063 valuable insights for our research, providing new perspectives on how to balance cache compression 064 with maintaining model performance. 065

Inspired by the above research insights, we propose the hypothesis that in long-text reasoning tasks, 066 LLMs may process information similarly to how humans read and comprehend lengthy articles by 067 relying primarily on key hints. In other words, the KV cache may only need to store the most 068 important tokens from the redundant text. Through statistical analysis, we identified a remarkably 069 high degree of sparsity in the attention layer. This finding validates our hypothesis that LLMs heavily 070 depend on key cues to achieve their superior long-text memory processing capabilities. It also 071 suggests that much of the redundant information in long texts may be superfluous for comprehension 072 and can be effectively optimized out of the KV cache. 073

In this paper, we introduce IntelLLM, a lightweight framework tailored for long-text inference in 074 LLMs. IntelLLM incorporates center of gravity eviction (CGE) strategy and remote gap localization 075 (RGL) strategy, two novel techniques that effectively balance cache compression and model perfor-076 mance, all without fine-tuning. In Section 3, we analyze the shortcomings of the sliding window 077 mechanism in handling long-sequence inference tasks, focusing on the behavior of the attention layer and providing detailed insights into the attention mechanism's characteristics. Based on the 079 problem analysis and findings in Section3, Section4 introduces the KV cache eviction and update algorithm for IntelLLM, elaborating on the center of GE strategy and the RGL strategy to mitigate 081 performance loss. In Section5, we present the experimental validation, where IntelLLM is integrated with Llama-3-8B-instruct and Mistral-7B-inst-v0.2, and evaluated using LongBench. Without finetuning, IntelLLM achieves performance comparable to baseline models while significantly reducing 083 cache memory usage. Its lightweight design also ensures easy integration into any existing LLM 084 system. 085

086 087

880

090

#### 2 RELATED WORK

#### 2.1 LLM LONG TEXT REASONING

The significant computational and memory demands during LLM inference pose critical challenges 091 for deployment in resource-constrained environments, limiting the capacity of LLMs to efficiently 092 process long texts. Prior research has predominantly focused on two areas: long-text compression 093 Zhao et al. (2024b) and model compression strategies Ramesh et al. (2023). Long-text compres-094 sion techniques, such as scalable embeddings, prompt compression, and activation beacons, aim 095 to reduce storage overhead for extended sequences. Meanwhile, model compression strategies fo-096 cus on optimizing the attention mechanism, either by designing more efficient architectures or by compressing the key-value (KV) cache to alleviate memory pressure. In this paper, we address the 098 bottlenecks in LLM inference, with a particular focus on optimizing the KV cache. 099

100 2.2 QUANTIZATION

For parameter and computation intensive tasks in LLMs, quantization techniques have proven to be
 highly effective in reducing memory consumption and accelerating inference. Current research pri marily focuses on enhancing computation-heavy operations through low-bit integer quantization of
 model weights. For instance, Wu et al. (a) employs INT4 quantization for both weights and caches,
 significantly reducing GPU memory requirements. Additionally, advanced hyper-quantization meth ods such as KVQuant Hooper et al. (2024a) and WKVQuant Yue et al. (2024) have been introduced
 to further optimize memory usage. However, while these quantization techniques succeed in com-

pressing models from a parametric standpoint, they do not fully address the KV cache pressure
 inherent in long dialogue history inference tasks. This gap highlights the need for more comprehen sive strategies to tackle the specific challenges of long-context processing in LLMs.

111 112

113

114

#### 2.3 KV COMPRESSION / SPARSE ATTENTION

115 KV cache optimization has become a pivotal research area due to the substantial computational bur-116 den imposed by the quadratic complexity of the attention mechanism. Techniques such as clustering-117 based hashing Kitaev et al. (2020) and k-nearest neighbors (kNN) Nawrot et al. (2024) methods effectively reduce computational complexity to super-linear levels. Moreover, the introduction of 118 sparse Transformers improves computational efficiency through sparse attention mechanisms. Ap-119 proaches like MQA Shazeer (2019) and GQA Ainslie et al. (2023) for instance, optimize the atten-120 tion module by redesigning attention heads, while more recent methods focusing on KV cache reuse 121 Gim et al. (2024) and cross-layer cache sharing Brandon et al. (2024) further improve memory effi-122 ciency. However, these approaches often incur additional training overheads, posing challenges for 123 deployment in resource-constrained environments. This underscores the demand for more efficient 124 solutions that can minimize resource consumption without sacrificing performance.

125 126

#### 2.4 Offloading

127 128

145

129 To address the high memory demands of KV caches on GPUs, researchers have begun exploring 130 alternative memory strategies involving multi-cluster systems or CPUs. For instance, DistKV-LLM 131 enhances cloud service inference by distributing KV caches across multiple servers. Similarly, 132 Pan et al. (2024) Wu et al. (b) propose a method that offloads KV caches to CPUs, where only a 133 small portion of the cache is reloaded to the GPU, significantly lowering the memory requirements for high-performance computing (HPC) devices. Although these approaches optimize inference 134 by leveraging external storage, they also introduce additional memory resource requirements. In 135 contrast, the research presented in this paper focuses on alleviating the computational burden of 136 inference by efficiently utilizing limited memory resources, without relying on external memory 137 expansion. 138

While existing compression and offloading techniques partially alleviate KV cache pressure in LLM
inference, the trade-off between performance and resource consumption persists. In this paper, we
introduce an optimization strategy that eliminates reliance on external memory, enabling effective
KV cache compression without sacrificing model performance. This novel approach provides a
robust solution for long-text reasoning in LLMs, particularly in resource-constrained environments,
thereby enhancing the practical applicability of LLMs in real-world settings.





161

### <sup>162</sup> 3 FINITE WINDOW HIGH COMPRESSIBILITY ASSUMPTION

163 164

#### 3.1 PROBLEM ANALYSIS

165 166 167

168

LLM is constructed based on casual transformer, which is divided into embedding layer, feedforward layer and self-attention layer. We focus on the self-attention layer in the decoding phase.

169 Extra-domain inference. During the pre-training phase, the model establishes a maximum text 170 length, making Softmax highly sensitive within this range. However, when covariates with se-171 quence lengths are introduced into the self-attentive mechanism with significantly exceeding the 172 pre-training length, they compromise the robustness of the attention distribution. This results in the 173 distributional moderation imbalance problem in Softmax, rendering LLM inference unpredictable. 174 As positional and semantic information cannot be easily decoupled, the pre-training window length becomes a constant challenge for extra-domain inference, especially without optimizing or fine-175 tuning the model architecture. We refer to this phenomenon as the problem of out-of-domain dis-176 tributional disequilibrium (ODD). 177

Intra-domain Inference. When the size of the sliding window is confined to the pre-training window, window attention mitigates the imbalance caused by relative positional differences in ODD. However, it still fails to reason effectively about long texts. From a semantic perspective, we contend that attention during window sliding relies entirely on short token sequences relative to the current time step. This premature discarding of long text information, coupled with the emphasis on short-term dependencies, directly contributes to the collapse of the LLM.

184 Consequently, we can draw two important insights:

Conclusion 1. Focusing solely on the attention layer, the stability of the LLM's ability to retain comprehension of long sequences is contingent upon the high sensitivity of Softmax in assigning weight streams.

Conclusion 2. The failure of sliding-window long-text inference arises from the limited robustness of the window's attention, indicating that the distribution of attention is heavily influenced by future markers. The query is continuously updated as the text expands, leading the model's attention to prioritize local information features at the expense of capturing contextual dependencies. This phenomenon is exacerbated by unpredictable noise in long texts, the prominence of locally relevant information, and conflicts between global features.

195 196

197

#### 3.2 LOCAL TOKEN MODELING

Building on the relationship between in-domain and out-of-domain inference as discussed in Section3.1, we divide the long context sequence into two components based on the current time step: local tokens  $L_{local}$  and remote tokens  $L_{remote}$ . Consequently, the long context can be denoted as  $L = L_{local} + L_{remote}$ . We then explore distinct compression strategies tailored to the unique features of these two components.

In the case of  $L_{local}$ , due to the autoregressive nature of LLM inference, the query at the current time step remains inaccessible to subsequent inference, while the query input at the next time step influences the reallocation of attention. This presents challenges in designing stable KV eviction strategies. To address this issue, we focus on uncovering solutions from the query cluster data itself.

**Research 1.** We explore the consistency characteristics of query clusters by analyzing tokens within the sequence  $L_{local}$ . For each head query in the same layer, we compute the cosine similarity and observe the intriguing vector distribution properties shown in Figure 1a. Notably, the closer the data points are to the diagonal split line, the redder their color becomes, indicating a high similarity between adjacent queries. This observation suggests that near-neighbor queries exhibit strong contextual similarity at the semantic level. Furthermore, this phenomenon aligns with the model's reliance on prior input information to generate appropriate responses.

**Theorem 1.** In longer inference streams, retaining a limited number of clusters of near-neighbor tokens helps stabilize normal inference for local short texts.

## 216 3.3 REMOTE TOKEN MODELING

For distant contexts, inspired by the success of sparse attention, we hypothesize that certain key tokens remain relevant for future reasoning in long-text inference. In other words, the LLM's comprehension and memory capabilities depend primarily on these key tokens, while irrelevant token clusters that appear in long texts serve only to support localized short-text reasoning and can be treated as inactive information suitable for eviction. If this assumption holds, it would significantly alleviate the memory pressure on the KV cache.

Research 2. We explore the relationship between key tokens and LLM capabilities in processing
 long texts through an analysis of attention scores. Focusing on the compressibility of distant con texts, we emphasize the importance of long-range dependencies within the attention mechanism.
 As shown in the Figure1b, we visualized attention scores from NarrativeQA's scene quizzes. A no table observation is the consistent alignment of red bands within the long-sequence attention scores,
 indicating that queries maintain coherence around specific tokens in lengthy historical texts. Addi tionally, the important tokens receiving high attention scores are highly sparse.

To further validate the pervasive sparsity of these high-scoring tokens, we set the attention threshold at time step t to 1/t, considering only attention scores exceeding this threshold as important key locations of focus for the query. As depicted in the Figure1c, both individual attention heads and entire attention layers demonstrate significant sparsity, with over 90% of the attention being sparse. This observation reinforces our hypothesis that key tokens in distant long-text contexts are highly sparse and exert a lasting impact on the model's reasoning over extended sequences.

Theorem 2. For remote tokens, queries exhibit high attention to only a few key historical messages,
with highly scored keys displaying significant sparsity across both the attention layers and heads.
By retaining only the important key-value pairs, the long-text processing capabilities of the LLM
can be preserved without loss, enabling a high degree of compression for distant tokens.

241 242

243 244

245

246

247

248

249

250 251

253

254

260

261

262

264 265 266

267

#### 4 INTELLLM

Based on the theorems from Section3, this section introduces IntelLLM, a novel approach that eliminates the need for fine-tuning while significantly reducing the memory demands of the KV cache and improving the LLM's generalization capability for longer sequences. In Section4.1, we present the CGE strategies, designed based on two distinct characteristics of language task data. Then Section4.2 details the design of the RGL strategy, drawing inspiration from positional encoding research. Furthermore, Section4.3 outlines the windowing mechanism implemented in IntelLLM.



Figure 2: IntelLLM composed of CGE and RGL.

#### 4.1 CENTER OF GRAVITY EVICTION

In attention mechanisms, softmax functions are commonly employed to compute weights that guide the selection and focus of information. As discussed in Section3.1, softmax is highly sensitive to attentional matching within the pre-training window. However, when we compress the KV cache by evicting redundant KVs, the tokens that accumulate within this window become a center of weight concentration, due to their importance in long-text processing.

274 275

From a theoretical perspective, when a significant accumulation of weights occurs, we can assume that the center of gravity of attention is represented by  $x_k$ ,  $x_k$  represents the clustering of important KVs, with its attention score significantly exceeding that of other tokens. Based on the properties of softmax, we can simplify the output by considering the limit as:

 $SoftMax(x)_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$ 

283

 $SoftMax(x)_i \approx \begin{cases} 1, & \text{if } i = k\\ 0, & \text{if } i \neq k \end{cases}$ (2)

(1)

This implies that the softmax output will be heavily concentrated in the region corresponding to the
 maximum weight, i.e., the cumulatively important keys, effectively rendering the attention scores of
 other tokens close to zero.

When softmax becomes imbalanced, with large weights dominating the output, the attention mechanism may lose focus on other critical information. This imbalance can reduce the model's representation capacity, as it fails to fully leverage all available input data.

To restore the balance of the softmax function and mitigate the inference instability caused by distant context compression, we adopt a strategy of evicting the attention center of gravity, which prevents large weight dominance and further stabilizes the distribution of attention scores. Based on Equation 1, the following formula can be derived:

297 298

$$SoftMax(x)_{i} = \frac{e^{x_{i}-x_{k}}}{\sum_{j=1}^{n} e^{x_{j}-x_{k}}} = \frac{e^{x_{i}-x_{k}}}{1+\sum_{j\neq k} e^{x_{j}-x_{k}}}$$
(3)

when  $x_i = x_k$ ,  $e^{x_i - x_k} = e^0 = 1$ . As  $x_j \ll x_k$ ,  $e^{x_j - x_k} \to 0$ . The CGE strategy effectively mitigates the imbalance caused by dominant large weights by redirecting attention away from the center of gravity.

At the same time, based on the principle of weight accumulation, we focus the attention's center of gravity on two key regions. As shown in the Figure1b, a small number of tokens at the beginning consistently accumulate a significant portion of the attention scores, forming **the head gravity**. Additionally, as discussed in Section3.2, the high similarity of near-neighbor queries leads to geometric accumulation of attention scores. Thus, the second region of interest corresponds to **the tail gravity**, which concentrates on these similar near-neighbor query clusters.

308 309

#### 4.2 REMOTE GAP LOCALIZATION

During the attention computation process, both the query and the key originate from the input sequence. By calculating the attention scores between them, the model identifies similarities, allowing it to focus on the most relevant information at each time step. The attention scores are therefore positively correlated with the degree of focus. Typically, without explicit positional encoding (e.g., absolute or relative position), the attention mechanism is not inherently sensitive to positional distance and relies solely on query-key similarity to highlight pertinent content.

316 We aligned the positional encoding with the compression window size and leveraged relative posi-317 tional differences to represent the relationship between the compressed historical long text and the 318 current query. However, our experimental results were less than satisfactory. So we hypothesize 319 that in the next time step of inference, simply reassigning positions is insufficient to capture the 320 attentional relationship between the current query and the compression window, even if the relative 321 positional distance remains within the pre-training length. As reasoning progresses, the information in the compression window spans much further than the current time step. Thus, the temporal struc-322 ture of the sequence at the semantic level remains intact, even when the KV cache is compressed 323 to fit within the pre-training window, and cannot be fully represented by a simple approximation of

<b>Input:</b> Last Query $Q \in \mathbb{R}^{n \times d}$ , Key Q	Cache $\hat{K} \in R^{m \times d}$ , Value Cache $\hat{V} \in R^{m \times d}$ , Head Gravity
Len $l_{head}$ , Tail Gravity Len $l_{head}$	ail
$A^0 \leftarrow QK^T / \sqrt{d}$	$\triangleright$ Attention score of $qk^2$
2 # Global dependency algorithms	
$A^1 \leftarrow softmax(A^0_{[:l_{bead}]})$	Masking head gravity of during normalization
4 $A^2 \leftarrow sum_v(A^1_{[:l_{tail}]})$	Masking tail gravity, sum along vertica
5 # Local fine-grained modeling algori	thms
6 $A^1 \leftarrow softmax(A^0_{[:l_{tail}]})$	Masking tail gravity of during normalizatio
$7 A^2 \leftarrow sum_v(A^1_{[:l_{head}]})$	▷ Masking head gravity, sum along vertica
s Indices $\leftarrow argsort(argtopk(A^2))$	Ranking indices of the top k position
9 $K_{comp}, V_{comp} \leftarrow gather(Indices)$	▷ Key clusters within the compression window
10 $\hat{K}, \leftarrow concat(K_{comp}, \hat{K}_{[-l_{tail}:]})$	
11 $\hat{V} \leftarrow concat(V_{comp}, \hat{V}_{[-l_{tail}:]})$	
12 $return\hat{K}, \hat{V}$	Integration of compression window and neighbor window

positional differences. This inherent positional relationship continues to influence the distribution of attention scores.

To address this, we propose the Remote Gap Localization (RGL) strategy, as shown in Figure2. The
 RGL strategy primarily addresses the issue of time span vanishing caused by high compression by
 assigning cache position values significantly larger than the length of the near-neighbor window to
 distant KVs within the compression window. Meanwhile, it preserves the relative position information of the near-neighbor window, ensuring stability in capturing long-range dependencies.

351 352

353

354

4.3 DESIGN

Leveraging the CGE and RGL strategies, IntelLLM incorporates two primary elements: the compression window and the nearest-neighbor window, as illustrated in Figure 2.

355 356 357

4.3.1 COMPRESSION WINDOW

In conjunction with the CGE strategy, we stabilize KV eviction by masking the influence of the center of gravity during the normalization process. However, through our experiments, we found that the combination of the masking order of the attention center of gravity and the normalization operation has different impacts on model performance. Based on these experimental observations, we categorize them into two types of characterization algorithms, providing explanations for each in relation to the semantic features of specific textual tasks.

364 Local Fine-grained Modeling Tasks. Compared to tasks requiring long-term global dependencies, 365 this task focuses more on modeling fine-grained local information in the current context. In natural 366 language processing, short-term dependencies are typically closely tied to the immediate context. 367 For example, in conversational tasks, information from neighboring time steps has a greater impact 368 on the model's predictions than distant contexts. Consequently, the model emphasizes the semantic 369 features of nearby tokens during inference. As shown in Algorithm1, we implement a local focus of attention by first masking the head gravity according to the short-term dependency strategy, followed 370 by removing the effect of the tail gravity after data normalization. 371

Global Dependency Tasks. Since this kind of task focuses more on requiring the model to rely
 on long-distance contextual information in the inference process, the compression strategy in this
 phase should focus more on the distant above. Therefore, as illustrated in Algorithm1, we first
 shield the influence of the tail gravity before the normalization operation, followed by masking
 the weaker influence of the head gravity during the selection of important keys. This approach
 effectively enhances the model's focus on distant information while emphasizing the consistency and completeness of global semantics.

## 4.3.2 NEAREST NEIGHBOR WINDOW

Based on the findings in Section 3.2, the nearest-neighbor window stabilizes the inference process for
 subsequent queries by retaining the tokens in close proximity. This approach ensures that the model
 maintains focus on relevant nearby information, supporting more accurate and stable predictions.

#### 5 EVALUATION

384 385 386

387

388

389 390 391 In this section, we demonstrate through validation experiments that IntelLLM achieves KV cache compression using a simple window combination. This approach effectively compensates for long text inference accuracy while delivering exceptional performance across a wide range of domain tasks.

#### 5.1 EXPERIMENTAL SETTTING

To comprehensively evaluate IntelLLM's performance in long text reasoning tasks, we employ LongBench for long document benchmarking. To validate the rationale behind the IntelLLM architecture design, we employ the original model with full KV caching as a strong baseline. Since IntelLLM aims to balance performance and KV compression without requiring fine-tuning, we also include two windowed approaches LM-Infinite and StreamingLLM as additional baselines. All evaluation experiments are conducted on a single NVIDIA A100 80GB GPU server.

Table 1: Comparison of model performance based on Llama3. The baseline results, marked with \*, are reproduced from Xiao et al.. All baseline models utilized the pre-training window size as their context length, while IntelLLM is configured with a total window length of 4K ( $L_{comp} = L_{near} = 2K$ ). Refer to the AppendixA for corresponding data IDs and names.

I lama_3	Window	Single-Doc QA			Multi-Doc QA			Code	
Liama-5	nama-5 vinuovi -		1-2	1-3	2-1	2-2	2-3	3-1	3-2
Full*	8K	19.85	42.36	41.03	47.38	39.2	22.96	60.83	49.14
Infinite*	8K	19.39	42.80	40.44	43.77	37.89	18.33	60.12	48.62
Streaming*	8K	20.05	42.46	39.54	43.69	37.89	19.68	60.35	48.95
IntelLLM*	4K	22.16	40.88	48.06	44.3	35.2	22.74	58.39	53.76
T 1	**** 1	Few-shot Learning		Synthetic	Summarization				
	M/Indow				•				
Liama-3	Window	4-1	4-2	4-3	5-1	6-1	6-2	6-3	
Full*	Window 8K	4-1 74	4-2 90.5	4-3 42.3	5-1 8.5	6-1 29.94	6-2 21.45	6-3 27.51	
Full* Infinite*	Window 8K 8K	4-1 74 74	4-2 90.5 90.08	4-3 42.3 41.72	5-1 8.5 4.5	6-1 29.94 29.25	6-2 21.45 21.41	6-3 27.51 27.62	
Full* Infinite* Streaming*	Window 8K 8K 8K	4-1 74 74 73.5	4-2 90.5 90.08 90.08	4-3 42.3 41.72 41.55	5-1 8.5 4.5 5	6-1 29.94 29.25 29.17	6-2 21.45 21.41 21.33	6-3 27.51 27.62 27.56	

414 415 416

#### 417 418 5.2 MAIN RESULTS

419 We select the Llama-3-8B-Instruct model for a comprehensive evaluation on the LongBench dataset, 420 with the corresponding test results presented in Tables1. Our data analysis leads to the follow-421 ing conclusions: (1) Compared to the baseline streaming input model, IntelLLM's token eviction 422 strategy demonstrates superior performance, indicating that IntelLLM enhances the LLM's under-423 standing of long sequential contextual information through key hints, thus maintaining the stability of streaming inference; (2) While the two windowing mechanisms efficiently conserve KV cache 424 space, they struggle to effectively capture long-range text dependencies. IntelLLM mitigates the 425 performance degradation caused by compression through its adaptive tuning mechanism for the 426 compression windows, offering a partial solution to this challenge; (3) With 50% KV cache com-427 pression, IntelLLM achieves performance close to or even exceeding that of the original strong 428 baseline, validating the effectiveness of our strategy while significantly reducing memory costs and 429 enabling long sequence reasoning in LLMs with low-resource utilization. 430

The results of our evaluation of Mistral-7B-inst-v0.2 are shown in Table3, which also indicates that IntelLLM outperforms other approaches.

Table 2: Comparison of model performance based on Mistral. The baseline results, marked with \*, are reproduced from Xiao et al..IntelLLM is configured with a total window length of 4K ( $L_{comp} = L_{near} = 2K$ ).

Mistral	Window	1-1	1-2	2-1	4-1	4-2	4-3	5-1	6-1
Full*	32K	22.06	47.65	21.96	26.62	71	85.97	42.29	3.95
Infinite*	6K	18.44	39.05	22.27	26.65	70	85.22	41.6	2.08
Streaming*	6K	17.92	39.09	21.83	26.64	70	85.57	41.31	2.5
IntelLLM	4K	19.19	47.11	21.59	26.63	70.5	86.81	41.67	2.87

**Latency.** Based on Llama3, we measure the inference latency for 8K text sequences. In comparison to the full-cache model inference latency of 900.84 ms, the KV update algorithm of IntelLLM adds only an additional 2.37 ms. Considering the achieved 50% cache savings, this 2.63% increase in latency is entirely acceptable.

**Ablation Study.** To validate the impact of IntelLLM's two policies on long text processing tasks, we select various task types and present the results of the ablation study for the Llama-3-8B-Instruct model in Table3. Given the inherent negative correlation between KV cache compression and performance, we use CGE and RGL as ablation terms to illustrate their effects on performance under limited KV caching conditions.

#### Table 3: Ablation with CGE and RGL.

IntelLLM	Head Gravity	Tail Gravity	RGL Gap	Result
	0	2K	4k	21.4
1-1	4	2K	4K	21.04
	6	2K	4k	22.16
IntelLLM	Head Gravity	Tail Gravity	RGL Gap	Result
IntelLLM	Head Gravity 4	Tail Gravity 2K	RGL Gap 4k	<b>Result</b> 30.94
IntelLLM 2-2	Head Gravity 4 4	Tail Gravity2K2K	RGL Gap 4k 6K	<b>Result</b> 30.94 31.83

**CGE Ablation.** CGE enhances the model's attentional focus by masking different attentional anchors, facilitating the management of both long and short dependencies across various task types. In short-term dependency tasks, we fix the positional interval between the tail center of gravity and the compression window, and assess the effectiveness of the CGE strategy by scaling the head center of gravity region. As shown in Table3, the weight distribution within the head center of gravity clusters has a significant impact on the model's performance in processing dialog tasks.

RGL Ablation. In Table3, we evaluate the correlation between RGL's positional distance interval settings and the model's inference performance. Given the constraints of the model's pre-training window, we limit the relative positional differences to fall within this range. The results show that using positional intervals to represent the semantic distance or time span between the nearest-neighbor window and the salient window proves to be an effective approach.

#### 6 CONCLUSION

In this paper, we tackle the memory challenges associated with KV caching in LLM deployments by introducing IntelLLM. Our extensive evaluation experiments demonstrate that IntelLLM significantly improves performance through the combined use of CGE and RGL strategies. This approach allows LLMs to achieve an optimal balance between KV cache compression and inference performance without requiring fine-tuning, all while avoiding substantial overhead in computational resources. As a result, IntelLLM effectively enhances the ability of LLMs to reason about long-text tasks.

## 486 REFERENCES

499

500

501

502

527

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit
   Sanghai. Gqa: Training generalized multi-query transformer models from multi-head check points. *arXiv preprint arXiv:2305.13245*, 2023.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
   *arXiv preprint arXiv:2004.05150*, 2020.
- William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan Kelly. Reducing transformer key-value cache size with cross-layer attention. *arXiv preprint arXiv:2405.12981*, 2024.
- Tianyu Fu, Xuefei Ning, Boju Chen, Tianqi Wu, Genghan Zhang, Guohao Dai, Huazhong Yang,
   and Yu Wang. Semsa: Semantic sparse attention is hidden in large language models.
  - Yingqiang Ge, Wenyue Hua, Kai Mei, Juntao Tan, Shuyuan Xu, Zelong Li, Yongfeng Zhang, et al. Openagi: When llm meets domain experts. *Advances in Neural Information Processing Systems*, 36, 2024.
- In Gim, Guojun Chen, Seung-seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. Prompt
   cache: Modular attention reuse for low-latency inference. *Proceedings of Machine Learning and Systems*, 6:325–338, 2024.
- Zhenyu He, Guhao Feng, Shengjie Luo, Kai Yang, Di He, Jingjing Xu, Zhi Zhang, Hongxia Yang, and Liwei Wang. Two stones hit one bird: Bilevel positional encoding for better length extrapolation. *arXiv preprint arXiv:2401.16421*, 2024.
- Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao,
   Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with
   kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024a.
- <sup>513</sup> Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao,
  <sup>514</sup> Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with
  <sup>515</sup> kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024b.
- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv* preprint arXiv:2001.04451, 2020.
- Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du,
  Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, et al. Cachegen: Kv cache compression and
  streaming for fast large language model serving. In *Proceedings of the ACM SIGCOMM 2024 Conference*, pp. 38–56, 2024.
- Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR, 2023.
- Piotr Nawrot, Adrian Łańcucki, Marcin Chochowski, David Tarjan, and Edoardo M Ponti. Dynamic memory compression: Retrofitting llms for accelerated inference. arXiv preprint arXiv:2403.09636, 2024.
- Xiurui Pan, Endian Li, Qiao Li, Shengwen Liang, Yizhou Shan, Ke Zhou, Yingwei Luo, Xiaolin
   Wang, and Jie Zhang. Instinfer: In-storage attention offloading for cost-effective long-context llm
   inference. *arXiv preprint arXiv:2409.04992*, 2024.
- Krithika Ramesh, Arnav Chavan, Shrey Pandit, and Sunayana Sitaram. A comparative study on the impact of model compression techniques on fairness in language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15762–15782, 2023.
- 539 Ninglu Shao, Shitao Xiao, Zheng Liu, and Peitian Zhang. Extensible embedding: A flexible multipler for llm's context length. *arXiv preprint arXiv:2402.11577*, 2024.

- 540 Noam Shazeer. Fast transformer decoding: One write-head is all you need. arXiv preprint 541 arXiv:1911.02150, 2019. 542
- Woomin Song, Seunghyuk Oh, Sangwoo Mo, Jaehyung Kim, Sukmin Yun, Jung-Woo Ha, and 543 Jinwoo Shin. Hierarchical context merging: Better long context understanding for pre-trained 544 llms. arXiv preprint arXiv:2404.10308, 2024.
- 546 Jianbo Wu, Jie Ren, Shuangyan Yang, Konstantinos Parasyris, Giorgis Georgakoudis, Ignacio La-547 guna, and Dong Li. Lm-offload: Performance model-guided generative inference of large lan-548 guage models with parallelism control. a.
- Jianbo Wu, Jie Ren, Shuangyan Yang, Konstantinos Parasyris, Giorgis Georgakoudis, Ignacio La-550 guna, and Dong Li. Lm-offload: Performance model-guided generative inference of large lan-551 guage models with parallelism control. b. 552
- 553 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, 554 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multiagent conversation framework. arXiv preprint arXiv:2308.08155, 2023. 555
- 556 Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, and Maosong Sun. Infilm: Training-free long-context extrapolation for llms with an efficient 558 context memory. In First Workshop on Long-Context Foundation Models@ ICML 2024. 559
- Lu Ye, Ze Tao, Yong Huang, and Yang Li. Chunkattention: Efficient attention on kv cache with 560 chunking sharing and batching. 561
- Yuxuan Yue, Zhihang Yuan, Haojie Duanmu, Sifan Zhou, Jianlong Wu, and Liqiang Nie. Wkvquant: 563 Quantizing weight and key/value cache for large language models gains more. arXiv preprint 564 arXiv:2402.12065, 2024. 565
- Michael Zhang, Aaryan Singhal, Benjamin Frederick Spector, Simran Arora, and Christopher Re. 566 Low-rank linearization of large language models. In Workshop on Efficient Systems for Founda-567 tion Models II@ ICML2024, a. 568
- Yuxin Zhang, Yuxuan Du, Gen Luo, Yunshan Zhong, Zhenyu Zhang, Shiwei Liu, and Rongrong Ji. 570 Cam: Cache merging for memory-efficient llms inference. In Forty-first International Conference on Machine Learning, b.
- Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm 573 agents are experiential learners. In Proceedings of the AAAI Conference on Artificial Intelligence, 574 volume 38, pp. 19632–19642, 2024a. 575
  - Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. Dense text retrieval based on pretrained language models: A survey. ACM Transactions on Information Systems, 42(4):1–60, 2024b.
    - Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, et al. A survey on efficient inference for large language models. arXiv preprint arXiv:2404.14294, 2024.

592

588 589

549

569

571

572

576

577

578

579

580

593

#### 594 Appendix А 595

596

597	Table 4: Dataset ID Mapping Table						
598							
599	Dataset	ID	Source	Metric			
600	Single-Doc QA		<b>X</b> 1				
601	NarrativeQA	1-1	Literature,Film	F1			
602	Qasper MultiFieldOA	1-2	Science				
603	Multi De sum ent QA	1-3	Multi-field	F1			
604	HotpotOA	2.1	Wilcipadia	F1			
605	2WikiMultihonOA	$2^{-1}$	Wikipedia	F1			
606	MuSiQue	2-3	Wikipedia	F1			
607	Code Completion	20	() hipedia				
608	LCC	3-1	Github	Edit Sim			
609	RepoBench-P	3-2	Github repository	Edit Sim			
610	Few-shot Learning		¥ ¥				
611	TREC	4-1	Web question	Accuracy			
612	TriviaQA	4-2	Wikipedia, Web	F1			
613	SAMSum	4-3	Dialogue	Rouge-L			
614	Synthetic Task						
615	PassageCount	5-1	Wikipedia	Accuracy			
616	Summarization		<b>a</b>	<b>D T</b>			
617	GovReport	6-1	Government report	Rouge-L			
618	QMSum	6-2	Meeting	Rouge-L			
619	Multinews	0-3	News	Rouge-L			
620							
621							
622							
623							
624							
625							
626							
627							
628							
620							
629							
621							
620							
632							
624							
625							
635							
607							
637							
638							
640							
040							
041							
042							
043							
044							
040							
040							
047							

#### Table 4: Dataset ID Mapping Table