LOOKAHEADKV: FAST AND ACCURATE KV CACHE EVICTION BY GLIMPSING INTO THE FUTURE WITHOUT GENERATION

Anonymous authorsPaper under double-blind review

ABSTRACT

Transformer-based large language models (LLMs) rely on key-value (KV) caching to avoid redundant computation during autoregressive inference. While this mechanism greatly improves efficiency, the cache size grows linearly with the input sequence length, quickly becoming a bottleneck for long-context tasks. Existing solutions mitigate this problem by evicting prompt KV that are deemed unimportant, guided by estimated importance scores. Notably, a recent line of work improves eviction quality by "glimpsing into the future", in which a low-cost draft generator first produces a surrogate response that mimics the target model's true response, which is subsequently used to estimate the importance scores of cached KV. In this paper, we propose LookaheadKV, a lightweight eviction framework that leverages the strength of surrogate future response without the need for costly draft generation. LookaheadKV augments transformer layers with parameter-efficient modules trained to predict true importance scores with high accuracy. Our design ensures negligible runtime overhead comparable to existing inexpensive heuristics, while achieving accuracy superior to more costly approximation methods. Extensive experiments on long-context understanding benchmarks, across a wide range of models, demonstrate that our method not only outperform recent competitive baselines in long-context understanding tasks by 25%, but also reduces the eviction cost by up to 14.5x, leading to significantly faster time-to-first-token (TTFT).

1 Introduction

Long context length of Large Language Models (LLMs) is becoming increasingly critical for many emerging applications: processing long documents (Bai et al., 2024; Wang et al., 2024a; Hsieh et al., 2024), repository-level code understanding and generation (Luo et al., 2024; Liu et al., 2024; Jimenez et al., 2024), in-context learning (Li et al., 2025), extension to long multi-modal inputs such as video (Wang et al., 2024b), etc. However, a central challenge in enabling these applications is that the key-value (KV) cache size grows linearly in sequence length, which rapidly becomes a bottleneck for inference, restricting scalable deployment of such applications on both mobile devices and the cloud. For example, even for moderate-sized models, such as LLaMA3.1–70B (Dubey et al., 2024) in half-precision, storing a single 64K-token sequence already takes up 40GB of memory, while scaling to 512K tokens requires 160GB, exceeding the memory capacity of high-end consumer hardware.

A growing line of work addresses this challenge by identifying salient tokens to achieve effective KV cache eviction without loss of performance (Li et al., 2024; Cai et al., 2024; Galim et al., 2025; Wang et al., 2025; Zhang et al., 2023). Early methods often rely on simple heuristics, in which token importance is estimated based on the self-attention scores of the input tokens. SnapKV (Li et al., 2024), for instance, leverages the attention weights between the suffix of the input and the preceding context to estimate the importance of each prompt token. However, investigations in recent studies (SpecKV (Galim et al., 2025), LAQ (Wang et al., 2025)) reveal that leveraging the model's response, rather than the input suffix, can greatly improve the eviction quality. Furthermore, they show that

Source code to reproduce our results is available, released in supplementary.

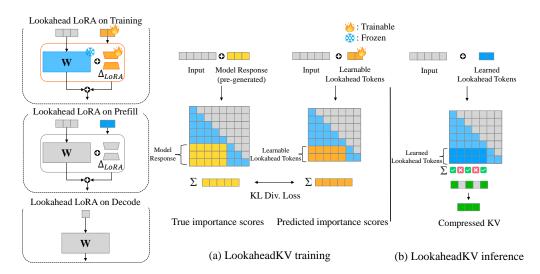


Figure 1: (a) Overview of LookaheadKV training (b) Overview of LookaheadKV inference.

a low-cost generated draft response, which closely approximates the true response, can serve as a powerful proxy for accurately estimating the importance scores. For example, SpecKV employs a smaller auxiliary model to produce draft tokens to approximate the target model's response, while Lookahead Q-Cache (LAQ) first applies a cheap KV eviction scheme to the target model, such as SnapKV, to obtain draft tokens, which in turn are used to approximate true importance scores.

While these draft-based methods substantially improve eviction quality, they often struggle with a fundamental trade-off between performance and efficiency, due to the need for costly draft token generation. Figure 2 presents the trade-off between accuracy and overhead of different approaches using the QASPER benchmark (Dasigi et al., 2021) and LLaMA3.1-8B-Instruct (Dubey et al., 2024) with a cache budget size of 128. While cheaper approaches like SnapKV are fast, inducing minimal overhead, they suffer a severe performance degradation under highly constrained budget settings. On the other hand, LAQ (Wang et al., 2025), a draft-based approach, shows impressive results even in extremely limited budget settings. However, it incurs a prohibitive computational overhead by generating an extra draft response, which limits its practicality in latency-sensitive applications such as mobile devices.

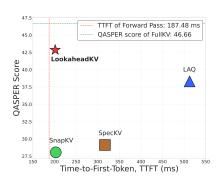


Figure 2: Accuracy-overhead Trade-off across KV cache eviction methods.

To overcome this limitation, we introduce LookaheadKV, a novel KV cache eviction method that augments LLMs with parameter-efficient modules, capable of accurately predicting future attention patterns, without the need for costly draft token generation. As shown in Figure 2, our method effectively overcomes the accuracy-overhead trade-off, achieving minimal performance loss with negligible overhead. LookaheadKV, as depicted in Figure 1, our method employs a set of learnable special tokens, together with *Lookahead LoRA* modules, novel low-rank adapters that selectively activate for the special tokens, to produce queries that can reliably estimate token-importance scores. By fine-tuning these modules to predict the true importance scores, LookaheadKV effectively minimizes the quality loss incurred by KV cache eviction with marginal inference overhead.

To rigorously assess the effectiveness of LookaheadKV, we evaluate it on a diverse set of long-context benchmarks (Bai et al., 2024; Hsieh et al., 2024; Ye et al., 2025) across multiple models of varying sizes (Dubey et al., 2024; Yang et al., 2025). Experimental results consistently demonstrate that LookaheadKV outperforms strong baselines across multiple budgets and context lengths while incurring significantly less eviction latency.

To summarize, our contributions are as follows:

- We propose LookaheadKV, the first KV cache eviction framework that employs learnable lookahead tokens and special LoRA modules to accurately predict the importance scores from the model's true response without generating costly approximate response.
- Through extensive experiments, we demonstrate that the proposed approach is effective and robust across different models and context lengths, and especially under low-budget settings, making our method particularly useful in resource-constrained environments.
- By conducting a rigorous analysis of eviction latency, both theoretically and empirically, we demonstrate that our method incurs negligible eviction overhead of less than 2.16% at 32K context length, while being $14.5\times$ faster than draft-based methods.

2 BACKGROUND

The primary objective of the KV cache eviction methods considered in this work, including our proposed approach, is to accurately estimate the importance score of individual key-value pairs of prompt tokens using attention weights, in order to guide the eviction process. In the following section, we formally define the problem of KV cache eviction and briefly discuss how prior methods have approached it.

KV cache eviction using importance scores. Let $X = \{x_1, ..., x_{n_{\text{in}}}\}$ be an input token sequence (e.g., a user instruction, part of a code snippet, etc.) and $Y = \{y_1, ..., y_{n_{\text{out}}}\}$ the model's generated response to X. For a given layer and attention head in an LLM, the attention scores of the complete sequence are given by:

$$\mathbf{Q} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \mathbf{W}_q \qquad \mathbf{K} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \mathbf{W}_k \qquad \mathbf{A} = \operatorname{Softmax} \left(\frac{\mathbf{Q} \mathbf{K}^{\top}}{\sqrt{d}} \right), \tag{1}$$

where $\mathbf{X} = [\mathbf{x}_1,...,\mathbf{x}_{n_{\text{in}}}]^{\top} \in \mathbb{R}^{n_{\text{in}} \times d}$ and $\mathbf{Y} = [\mathbf{y}_1,...,\mathbf{y}_{n_{\text{out}}}]^{\top} \in \mathbb{R}^{n_{\text{out}} \times d}$ are the hidden states of the input prompt and model-generated response, respectively. For better readability, we omit the layer and head index. We define the ground-truth importance scores $\mathbf{s}_{\text{GT}} = [s_1,...,s_{n_{\text{in}}}]$ of the KV cache as the average cross-attention scores between the queries of \mathbf{Y} and the keys of \mathbf{X} , i.e., $s_j = \frac{1}{n_{\text{out}}} \sum_{i=n_{\text{in}}+n_{\text{out}}}^{n_{\text{in}}+n_{\text{out}}} \mathbf{A}_{i,j}$. Intuitively, these scores quantify the relative contribution of each prompt token's key-value pair to the model's response generation. Based on these scores, the pruned KV cache can be obtained by retaining a subset of (e.g., TopK) important KV pairs to minimize the attention output perturbation, such that:

$$Attn(x, KV_{orig}) \approx Attn(x, KV_{GT}),$$
 (2)

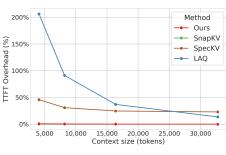
where KV_{orig} and KV_{GT} are the original and evicted KV cache using the ground-truth importance scores, respectively.

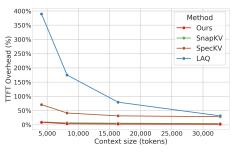
However, since the model's true future response is unknown during the prefill phase, such scores cannot be computed directly. Consequently, prior methods resorted to constructing a surrogate response sequence $\tilde{\mathbf{Y}} = [\tilde{y}_1,...,\tilde{y}_{n_{\mathrm{window}}}]^{\top} \in \mathbb{R}^{n_{\mathrm{window}} \times d}$ to approximate the model's (partial) future response and predict the attention pattern:

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{Y}} \end{bmatrix} \mathbf{W}_q$$
 $\tilde{\mathbf{K}} = \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{Y}} \end{bmatrix} \mathbf{W}_k$ $\tilde{\mathbf{A}} = \operatorname{Softmax} \left(\frac{\tilde{\mathbf{Q}} \tilde{\mathbf{K}}^{\top}}{\sqrt{d}} \right),$ (3)

resulting in the estimated importance score vector $\mathbf{s}_{\text{approx}} = [\tilde{s}_1, ..., \tilde{s}_{n_{\text{in}}}]$, whose entries are computed as $\tilde{s}_j = \frac{1}{n_{\text{window}}} \sum_{i=n_{\text{in}}+1}^{n_{\text{in}}+n_{\text{window}}} \tilde{\mathbf{A}}_{i,j}$. In short, these methods aim to obtain the estimated score vector whose ranking is similar to that of the ground-truth, such that the overlap between the retained KV pairs and KV_{GT} is high. Various approaches have been suggested to approximate the future response for effective KV cache eviction.

SnapKV. SnapKV (Li et al., 2024) proposes to use the suffix of input prompt to compute the estimate of the true future importance scores. Because SnapKV requires only marginal extra computation to perform eviction, as it uses attention weights that are already computed during the prefill forward pass, it has widely been adopted as a cheap and effective heuristic for KV cache eviction.





(a) Theoretical latency overhead

(b) Actual latency overhead

Figure 3: Time-to-First-Token (TTFT) latency overhead ratio across context lengths. Similar to SnapKV, LookaheadKV introduces negligible TTFT overhead across all tested context lengths; draft-based methods (LAQ, SpecKV) incur substantial latency, especially for shorter contexts.

SpecKV and LAQ. Recently, several methods have proposed to use a low-cost generator to generate a (partial) approximate response first, and subsequently use it to estimate the true future importance scores. Notably, SpecKV (Galim et al., 2025) employs a smaller LLM to first generate a draft response, while Lookahead Q-Cache (LAQ) (Wang et al., 2025) first applies SnapKV to the target model to generate a draft response, which is in turn used to approximate the future importance.

These draft-based methods have consistently shown superior performance compared to cheaper heuristics (Li et al., 2024), demonstrating the effectiveness of employing surrogate future response, i.e., by "glimpsing into the future". However, the extra draft generation step still incurs substantial additional compute, resulting in significant increase in latency, as shown in Figure 3. In summary, existing methods face a clear trade-off: inexpensive heuristics are fast but less accurate, whereas draft-based techniques improve performance at the cost of increased inference time.

3 Proposed Method: LookaheadKV

To overcome the challenge of fast and accurate importance prediction, we introduce LookaheadKV, a framework that augments the LLM with a set of lightweight learnable modules which are optimized to predict ground-truth importance scores and guide the eviction process. LookaheadKV achieves the best of both worlds: 1) it eliminates the need for generating a draft response for each query, resulting in significantly faster KV cache eviction, and 2) it employs learned special tokens that serve as approximate future response for importance estimation, leveraging the strength of draft-based methods. The following section (and Figure 1) presents the detailed workflow of LookaheadKV.

3.1 MAIN COMPONENTS

Learnable Lookahead Tokens. LookaheadKV initially performs KV cache eviction using a set of learnable special tokens during the pre-fill phase, and subsequently decodes auto-regressively with the retained KV cache. Specifically, for given input prompt tokens X, LookaheadKV appends a sequence of trainable "soft" lookahead tokens $L = \{l_1, ..., l_{n_{\text{lookahead}}}\}$ whose queries in the attention layers are used to estimate the attention pattern of the true model response. In essence, these tokens are trained to compress the attention information of the true response to serve as the "observation window" in the eviction phase. These tokens are used during the pre-fill stage only for eviction, and adds zero additional overhead for decoding.

Lookahead LoRA. To enhance the quality of estimation, we introduce *Lookahead LoRA*, a novel low-rank adapter module that only activates for the lookahead tokens. Lookahead LoRA allows these tokens to learn richer representations to compress the information of the true future attention pattern, while ensuring that the outputs of normal input tokens are unchanged, preserving the original model behavior. Moreover, since the original model weights remain unaltered, LookaheadKV modules can be selectively enabled or disabled depending on the particular requirements of a given application, thereby broadening the method's applicability.

Combining the modules together, LookaheadKV computes the queries and keys of the complete sequence as follows:

$$\mathbf{Q}_{\mathrm{LKV}} = \begin{bmatrix} \mathbf{X} \\ \mathbf{L} \end{bmatrix} \mathbf{W}_{q} + \begin{bmatrix} \mathbf{0} \\ \mathbf{L} \end{bmatrix} \Delta \mathbf{W}_{q} \qquad \mathbf{K}_{\mathrm{LKV}} = \begin{bmatrix} \mathbf{X} \\ \mathbf{L} \end{bmatrix} \mathbf{W}_{k} + \begin{bmatrix} \mathbf{0} \\ \mathbf{L} \end{bmatrix} \Delta \mathbf{W}_{k}$$
(4)

where $\mathbf{L} \in \mathbb{R}^{n_{\mathrm{lookahead}} \times d}$ denotes the hidden states of the lookahead embeddings, and $\Delta \mathbf{W}_q$, $\Delta \mathbf{W}_k$ are the Lookahead LoRA modules for query and key projections. Accordingly, the attention pattern $\mathbf{A}_{\mathrm{LKV}} = \mathrm{Softmax}(\frac{\mathbf{Q}_{\mathrm{LKV}} \, \mathbf{K}_{\mathrm{LKV}}^{\top}}{\sqrt{d}})$, is used to estimate the importance score $\tilde{s}_j = \frac{1}{n_{\mathrm{lookahead}}} \sum_{i=n_{\mathrm{in}}+1}^{n_{\mathrm{in}}+n_{\mathrm{lookahead}}} \mathbf{A}_{\mathrm{LKV} \, i,j}$, which is in turn used for effective KV cache eviction.

3.2 LOOKAHEADKV TRAINING

We train LookaheadKV modules to compress the attention pattern of the true future response, using the model-generated responses as target. Specifically, given a data pair (X,Y), one iteration of LookaheadKV training consists of the following steps:

- 1. **GT Forward Pass.** For each layer l=1,...,L and head h=1,...,H, the ground-truth importance scores $\mathbf{s}_{\mathrm{GT}}^{l,h}$ between the input prompt X and model-generated response Y are computed.
- 2. **Lookahead Forward Pass.** Similarly, for each layer l and head h, we obtain the importance score estimates $\mathbf{s}_{\mathsf{LKV}}^{l,h}$ between the input prompt X and the lookahead tokens L.
- 3. **Loss Computation.** We first normalize all score vectors so that they sum to 1, and compute the average KL divergence loss between the GT and LookaheadKV importance scores across all heads and layers:

$$\mathcal{L} = \frac{1}{|L|} \frac{1}{|H|} \sum_{l}^{L} \sum_{h}^{H} \text{KL} \left(\text{Norm}(\mathbf{s}_{\text{GT}}^{l,h}) \parallel \text{Norm}(\mathbf{s}_{\text{LKV}}^{l,h}) \right).$$
 (5)

The loss is backpropagated to update the weights of the lookahead embeddings and Lookahead LoRA modules, while all other LLM layers remain frozen. The pseudo-code for LookaheadKV training and eviction is given in Algorithm 1 and Algorithm 2.

Training Objective. We want to ultimately optimize the similarity of the ranking between the two attention score vectors, such that we obtain TopK indices identical to those from ground-truth importance scores. Inspired from works on distilling attention scores (Wang et al., 2020; Izacard & Grave, 2021), we minimize the KL divergence between these normalized attention scores. As our attentions scores are normalized, this KL divergence is equivalent to the popular ListNet (Cao et al., 2007) ranking loss, with ϕ of ListNet as identity instead of exp.

Lookahead LoRA Overhead. LookaheadKV consistently achieves strong performance with negligible eviction overhead. In principle, one can apply Lookahead LoRA to only a subset of the linear layers to tradeoff accuracy and latency. However, even when Lookahead LoRA is applied to every linear layer, there is a mere 1.6% increase in the latency compared to not using Lookahead LoRA at all (see Table 3 for ablation results). Further, appending a small number of lookahead tokens (≤ 32 in our experiments) during prefill for importance estimation similarly incurs minimal cost. Consequently, we train LookaheadKV with LoRA modules applied to all linear layers.

To avoid materializing the full attention score matrix, we use FlashAttention (Dao et al., 2022) in the forward pass, coupled with eager attention for importance score computation and loss backpropagation, as detailed in Section C.

4 EXPERIMENTS

4.1 Training

Dataset. To encourage the model to learn from diverse attention patterns, we curate training samples of varying lengths and sources, comprising of both instruction-following datasets as well as pretraining texts. We collect 50K samples from the long_sft subset of the ChatQA2 (Xu et al., 2025) dataset,

20K samples from the Tulu (Lambert et al., 2025) instruction-following dataset, 7K samples from the Stack (Kocetkov et al., 2023), and 9K few-shot completion data samples that we create based on the training splits of the MetaMath, ARC, and HellaSwag datasets, originally curated in Pal et al. (2024). For instruction-following data, we remove the last assistant response, and use the target model to obtain the (X,Y) pairs of input prompt and model response. For pretraining documents, we first truncate the text at random positions to obtain X, and use the target model to complete the sequence to obtain Y. We limit the maximum input sequence length to 16K, and generate all training responses using greedy decoding and max generation length of 512.

Training Details. We apply LookaheadKV on two widely used open-source architectures, LLaMA (Dubey et al., 2024) and Qwen (Yang et al., 2025), covering three model sizes each: LLaMA3.2-1B, LLaMA3.2-3B, LLaMA3.1-8B, Qwen3-1.7B, Qwen3-4B, and Qwen3-8B. For all models, we set the lookahead size $n_{\rm lookahead}=32$, and apply LoRA to all projection and feed-forward modules ($\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o, \mathbf{W}_{up}, \mathbf{W}_{down}$, and \mathbf{W}_{gate}) with rank r=8 and scaling factor $\alpha=32$. This configuration introduces less than 0.5% additional trainable parameters across all models, as summarized in Table 1. Full hyperparameter settings are provided in Table 6.

Table 1: Additional trainable parameters introduced by LookaheadKV.

Model	Trainable Params					
1110401	Params	% of Model				
LLaMA3.2-1B	5.4M	0.44				
LLaMA3.2-3B	11.9M	0.37				
LLaMA3.1-8B	20.6M	0.26				
Qwen3-1.7B	8.5M	0.49				
Qwen3-4B	16.2M	0.40				
Qwen3-8B	21.5M	0.26				

4.2 EVALUATION SETUP

We evaluate our method on two popular long-context benchmarks: LongBench (Bai et al., 2024) and RULER (Hsieh et al., 2024). LongBench is a multi-task benchmark that comprehensively assesses long-context understanding across diverse tasks, such as question answering, summarization, few-shot learning, and code completion. We report results on the 16 English tasks, and use the average score as the main metric. RULER is another multi-task synthetic benchmark, primarily comprising 13 Needle-in-a-Haystack-style subtasks. Each sample can be constructed at varying sequence lengths, allowing systematic evaluation of scaling behavior. Similar to LongBench, we use average score as the main metric, and report the results at 4K, 8K, 16K and 32K context lengths.

Baselines. We compare our method against popular KV-cache eviction methods: 1) **SnapKV** (Li et al., 2024), 2) **PyramidKV** (Cai et al., 2024), and 3) **StreamingLLM** (Xiao et al., 2024). Additionally, we include stronger, more recent baselines that involve costly approximate future response generation, including 4) **Lookahead Q-Cache** (LAQ) (Wang et al., 2025), and for 8B-scale models, 5) **SpecKV** (Galim et al., 2025). In all experiments, Llama3.2-1B-Instruct and Qwen3-1.7B are used as draft models for Llama3.1-8B-Instruct and Qwen3-8B, respectively. We follow the standard eviction configuration settings for all baseline methods, which we detail in Section F

4.3 PERFORMANCE RESULTS

LongBench evaluation. Figure 4 shows the average LongBench scores of LookaheadKV and baselines, across cache budget settings ranging from 64 to 2048. Our method consistently demonstrates superior performance across all models and all budgets tested, demonstrating the effectiveness and robustness of our approach. Overall, results show that expensive draft-based methods, e.g., LAQ and SpecKV, outperform cheaper baselines, corroborating that employing approximate future response for importance estimation is effective. Nevertheless, our method significantly outperforms the draft-based approaches, especially at lower budget settings, highlighting that learning to estimate future importance is crucial for performance preservation. Due to space limitation, we report the results of 1B-scale models in Section E.

RULER evaluation. We report the RULER evaluation results of all methods with a fixed budget of 128 in Figure 4 (1B-scale results are provided in Section E). LookaheadKV consistently outperforms other baseline approaches here as well, maintaining strong performance across all evaluated context lengths. Further, note that while we limit the maximum training sequence length of LookaheadKV to 16K, our method generalizes to longer context length of 32K. We conduct additional experiments on the impact of training context length in Section 5.3.

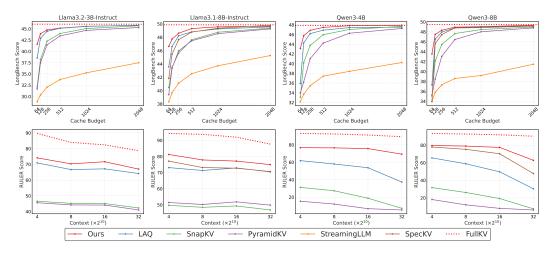


Figure 4: Top row: Average LongBench results across multiple budgets and models. Bottom row: Average RULER results across varying context lengths with a fixed budget of 128. Across all tested models, budgets and context lengths, our method consistently demonstrates superior performance.

Long-Form output Evaluation. We further evaluate LookaheadKV on the HTML to TSV task from LongProc (Ye et al., 2025), which involves extracting structured information from long HTML documents and converting it into TSV format. This benchmark tests not only the model's ability to process long-context inputs, but also its capacity to generate long-form outputs. We assess LookaheadKV and baseline methods under two input—output settings: 12K—0.5K and 23K—2K tokens, both at a fixed cache budget ratio of 30%.

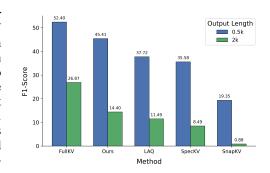


Figure 5 presents the results on the HTML to TSV task. Across both sequence-length config-

Figure 5: HTML-to-TSV evaluation results.

urations, LookaheadKV consistently outperforms prior approaches. We hypothesize that LookaheaKV, which learns to predict the attention pattern of the entire future response, is particularly superior in long-form generation tasks compared to draft-based methods that rely only on partial future response as the observation window.

5 ANALYSIS

5.1 EFFICIENCY COMPARISON

To assess the efficiency of our method against the baselines, we measure the Time-To-First-Token (TTFT) across multiple context lengths using their official implementations, with the exception of LAQ which we re-implement since it does not have an official implementation. Furthermore, since the latency of a method can vary significantly depending on the implementation, we conduct rigorous analysis and derive the theoretical latency for each method, based on the analytical model proposed in Davies et al. (2025). We discuss further details in Section B.

Table 2 presents the results of the TTFT analysis for 8K and 32K context lengths (see Table 5 for 4K and 16K results). Overall, we observe that draft-based methods incur significant overhead, either due to increased computation (SpecKV) or memory traffic (LAQ). On the contrary, LookaheadKV requires marginal additional cost across all tested context lengths, achieving 14.5 times faster eviction overhead compared to LAQ at 32K sequence length.

Table 2: Theoretical and Practical Analysis across various context lengths and methods.

			Theoretica		Empirical Cost			
Context Length	Method	Compute (TFLOPs)	Memory Traffic (GB)	TTFT (ms)	TTFT Overhead (ms)	TTFT (ms)	TTFT Overhead (ms)	
	Forward Pass Only	136	13	257	N/A	291	N/A	
	LookaheadKV (ours)	137	13	258	1.03	302	11	
8192	SnapKV	136	13	257	0.01	311	20	
	SpecKV	159	81	337	79.53	411	121	
	LAQ	137	445	492	234.59	800	509	
	Forward Pass Only	928	13	1754	N/A	1760	N/A	
	LookaheadKV (ours)	929	13	1755	1.74	1798	38	
32768	SnapKV	928	13	1754	0.01	1838	78	
	SpecKV	1115	106	2156	402.80	2263	503	
	LAQ	930	406	1993	239.26	2314	554	

5.2 ABLATION ON TRAINABLE MODULES

We study the impact of lookahead size $n_{lookahead}$ and LoRA placement through a 2D ablation across four lookahead sizes (4, 8, 16, 32) and three configurations: emb-only (No LoRA applied), QV (LoRA applied to Q and V), and all (LoRA applied to all linear layers). The results show that both larger lookahead sizes and broader LoRA coverage consistently improve average LongBench performance. Importantly, these gains come with negligible inference overhead, as the additional parameters account for only a tiny fraction of pre-filling compute and memory. Based on this analysis, we set $n_{lookahead} = 32$ and apply LoRA to all linear modules in our main experiments.

Table 3: 2D ablation across lookahead sizes and trainable modules, on LLaMA 1B. Average Long-Bench scores and TTFT overhead are reported.

	n_{lo}	okahead = 4	$n_{ m lc}$	$n_{\mathrm{lookahead}} = 8$		okahead = 16	$n_{\mathrm{lookahead}} = 32$		
Module	score	overhead(%)	score	overhead(%)	score	overhead(%)	score	overhead(%)	
emb-only	25.5	3.4	25.7	3.8	26.4	3.4	26.4	4.2	
QV	26.5	3.7	26.4	4.1	26.9	4.0	26.9	4.4	
all	26.6	4.2	27.0	4.2	27.0	4.7	27.1	5.0	

5.3 ROBUSTNESS TO TRAINING CONTEXT LENGTH

Transformer-based language models trained with fixed context lengths often struggle to generalize beyond their training window. Similarly, one may raise concern about the context length generalization of our method. To examine this effect, we apply LookaheadKV training to LLaMA-3B with limited training context lengths of 2K, 4K, and 8K, and evaluate on RULER (Figure 6). We observe that while longer training context lengths yield better performance as expected, training on shorter contexts still remains effective with relatively minor degradation in performance, demonstrating that our method generalizes robustly to unseen sequence lengths.

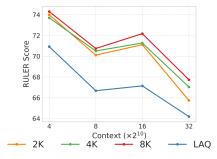


Figure 6: RULER evaluation on LookaheadKV trained with shorter contexts.

6 RELATED WORK

KV Cache Eviction. Early analyses revealed that attention scores tend to be sparse (Zhang et al., 2023), implying that only a small subset of KV entries substantially contributes to the attention output. Subsequent work showed that the importance of these tokens remains stable throughout generation, i.e., tokens deemed important early on tend to stay important (Liu et al., 2023). These observations motivated a range of eviction methods aimed at discarding unimportant KV entries while preserving model performance. A representative method is H2O (Heavy-Hitter Oracle) (Zhang et al., 2023), which proposes an eviction policy that considers the historical importance of tokens based on attention weights. NACL (Chen et al., 2024) performs eviction in a chunk-wise fashion, computing token importance locally within each chunk.

Prefill KV Cache Eviction. Another line of work, which we discuss extensively in our paper, focuses on eviction of prefill KV-cache. SnapKV (Li et al., 2024) introduced the notion of an "observation window" consisting of the suffix of the input prompt, which is used to predict important tokens to keep for subsequent response generation. Further, SpecKV (Galim et al., 2025) proposed to generate an approximate response with a smaller model and use the resulting tokens as a more reliable observation window for future importance prediction. Similarly, Lookahead Q-Cache (Wang et al., 2025) first applies a cheap eviction method, such as SnapKV, to obtain a partial low-cost draft response, then re-evicts KV entries based on the importance scores derived from the draft. KV-zip (Kim et al., 2025) adopts a query-agnostic strategy by inserting a repeated prompt and measuring which KV entries are essential for accurately reconstructing the input. Orthogonal to these approaches, several works proposed to allocate non-uniform budgets for each layer (Cai et al., 2024) and head (Feng et al., 2024) to further improve performance.

Prompt Tuning for Task Adaptation. Another line of work closely related to ours is parameter-efficient finetuning through learned prompts. Prompt Tuning (Lester et al., 2021) inserts a sequence of continuous, learnable embeddings into the frozen LLM for downstream task adaptation, while Prefix-Tuning (Li & Liang, 2021) extends this idea by pre-pending learned vectors across multiple layers. Further, P-Tuning v2 (Liu et al., 2022) demonstrated that prompt-based adaptation scales well across a wide range of model sizes. Unlike conventional prompt-tuning methods that aim to improve task performance, our work leverages learned prompts to predict internal model statistics, thereby enhancing computational efficiency rather than accuracy.

Training objectives similar to ours have been used in distillation (Wang et al., 2020), or in ranking/retrieval (Cao et al., 2007; Izacard & Grave, 2021). Some contemporaneous works (Greenewald et al., 2025; Peng et al., 2025; Samragh et al., 2025) also propose LoRA modules that selectively activate only for some tokens.

7 CONCLUSION AND LIMITATION

We introduce LookaheadKV, a trainable prefill-time KV cache eviction framework that accurately predicts token importance without relying on draft generation. The method augments a frozen LLM with a small set of learnable lookahead tokens and Lookahead LoRA modules that activate only on these tokens. Trained to match ground-truth importance distributions across layers and heads, LookaheadKV achieves performance superior to costly draft generation-based approaches while adding negligible inference overhead. Empirically, across LLaMA and Qwen model families and multiple long-context benchmarks, our approach consistently outperforms training-free heuristics and draft-based baselines, especially in low-budget regimes and long-from output tasks, while introducing less than 0.5% additional parameters and incurring only a marginal increase in prefill latency.

Due to limited compute resources, we were unable to conduct experiments on larger-sized models. Experiments on longer contexts should also be explored, but our analysis indicates that LookaheadKV training generalizes to longer context lengths. We greedily generate all responses for training, but the interaction between the decoding parameters, e.g., temperature, and LookaheadKV performance could be further explored. Lastly, LookaheadKV focuses on the prefill KV cache eviction; extending LookaheadKV to also perform decoding-stage eviction remains a future work.

8 REPRODUCIBILITY STATEMENT

Our source code is released in supplementary to reproduce our results, and pseudo-code is also provided in Section A. Section G provides links to datasets and evaluation benchmarks used, and Section 4.1 describes the pre-processing steps on the data.

REFERENCES

- 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. LongBench: A bilingual, multitask benchmark for long context understanding. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3119–3137, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.172. URL https://aclanthology.org/2024.acl-long.172/.
- Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Yucheng Li, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Junjie Hu, et al. Pyramidky: Dynamic ky cache compression based on pyramidal information funneling. *ArXiv preprint*, abs/2406.02069, 2024. URL https://arxiv.org/abs/2406.02069.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th International Conference on Machine Learning*, ICML '07, pp. 129–136, New York, NY, USA, 2007. Association for Computing Machinery. ISBN 9781595937933. doi: 10.1145/1273496.1273513. URL https://doi.org/10.1145/1273496.1273513.
- Yilong Chen, Guoxia Wang, Junyuan Shang, Shiyao Cui, Zhenyu Zhang, Tingwen Liu, Shuohuan Wang, Yu Sun, Dianhai Yu, and Hua Wu. NACL: A general and effective KV cache eviction framework for LLM at inference time. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7913–7926, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.428. URL https://aclanthology.org/2024.acl-long.428/.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/67d57c32e20fd0a7a302cb81d36e40d5-Abstract-Conference.html.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. A dataset of information-seeking questions and answers anchored in research papers. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4599–4610, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.365. URL https://aclanthology.org/2021.naacl-main.365/.
- Michael Davies, Neal Crago, Karthikeyan Sankaralingam, and Christos Kozyrakis. Efficient Ilm inference: Bandwidth, compute, synchronization, and capacity are all you need, 2025. URL https://arxiv.org/abs/2507.14397.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.

- Yuan Feng, Junlin Lv, Yukun Cao, Xike Xie, and S Kevin Zhou. Ada-kv: Optimizing kv cache eviction by adaptive budget allocation for efficient llm inference. *ArXiv preprint*, abs/2407.11550, 2024. URL https://arxiv.org/abs/2407.11550.
 - Kevin Galim, Ethan Ewer, Wonjun Kang, Minjae Lee, Hyung Il Koo, and Kangwook Lee. Draft-based approximate inference for llms. *ArXiv preprint*, abs/2506.08373, 2025. URL https://arxiv.org/abs/2506.08373.
 - Kristjan Greenewald, Luis Lastras, Thomas Parnell, Vraj Shah, Lucian Popa, Giulio Zizzo, Chulaka Gunasekara, Ambrish Rawat, and David Cox. Activated lora: Fine-tuned llms for intrinsics, 2025. URL https://arxiv.org/abs/2504.12397.
 - Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and Boris Ginsburg. RULER: What's the real context size of your long-context language models? In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=kIoBbc76Sy.
 - Gautier Izacard and Edouard Grave. Distilling knowledge from reader to retriever for question answering. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id=NTEz-6wysdb.
 - Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R. Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
 - Jang-Hyun Kim, Jinuk Kim, Sangwoo Kwon, Jae W. Lee, Sangdoo Yun, and Hyun Oh Song. KVzip: Query-agnostic KV cache compression with context reconstruction. In *ES-FoMo III: 3rd Work-shop on Efficient Systems for Foundation Models*, 2025. URL https://openreview.net/forum?id=gcqzyyF654.
 - Denis Kocetkov, Raymond Li, Loubna Ben allal, Jia LI, Chenghao Mou, Yacine Jernite, Margaret Mitchell, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro Von Werra, and Harm de Vries. The stack: 3 TB of permissively licensed source code. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=pxpbTdUEpD.
 - Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxi Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training. In *Second Conference on Language Modeling*, 2025. URL https://openreview.net/forum?id=iluGbfHHpH.
 - Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3045–3059, Online and Punta Cana, Dominican Republic, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.243. URL https://aclanthology.org/2021.emnlp-main.243.
 - Cheng Li. Llm-analysis: Latency and memory analysis of transformer models for training and inference. https://github.com/cli99/llm-analysis, 2023.
 - Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context LLMs struggle with long in-context learning. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL https://openreview.net/forum?id=Cw2xlg0e46.
 - Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*

Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 4582–4597, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353. URL https://aclanthology.org/2021.acl-long.353.

- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. Snapkv: LLM knows what you are looking for before generation. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/28ab418242603e0f7323e54185d19bde-Abstract-Conference.html.
- Tianyang Liu, Canwen Xu, and Julian J. McAuley. Repobench: Benchmarking repository-level code auto-completion systems. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=pPjZIOuQuF.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 61–68, Dublin, Ireland, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-short.8. URL https://aclanthology.org/2022.acl-short.8.
- Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance hypothesis for LLM KV cache compression at test time. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/a452a7c6c463e4ae8fbdc614c6e983e6-Abstract-Conference.html.
- Qinyu Luo, Yining Ye, Shihao Liang, Zhong Zhang, Yujia Qin, Yaxi Lu, Yesai Wu, Xin Cong, Yankai Lin, Yingli Zhang, Xiaoyin Che, Zhiyuan Liu, and Maosong Sun. RepoAgent: An LLM-powered open-source framework for repository-level code documentation generation. In Delia Irazu Hernandez Farias, Tom Hope, and Manling Li (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 436–464, Miami, Florida, USA, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-demo.46. URL https://aclanthology.org/2024.emnlp-demo.46/.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *ArXiv preprint*, abs/2402.13228, 2024. URL https://arxiv.org/abs/2402.13228.
- Yuqi Peng, Lingtao Zheng, Yufeng Yang, Yi Huang, Mingfu Yan, Jianzhuang Liu, and Shifeng Chen. Tara: Token-aware lora for composable personalization in diffusion models, 2025. URL https://arxiv.org/abs/2508.08812.
- Mohammad Samragh, Arnav Kundu, David Harrison, Kumari Nishu, Devang Naik, Minsik Cho, and Mehrdad Farajtabar. Your llm knows the future: Uncovering its multi-token prediction potential, 2025. URL https://arxiv.org/abs/2507.11851.
- Minzheng Wang, Longze Chen, Fu Cheng, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan Xu, Lei Zhang, Run Luo, Yunshui Li, Min Yang, Fei Huang, and Yongbin Li. Leave no document behind: Benchmarking long-context LLMs with extended multi-doc QA. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 5627–5646, Miami, Florida, USA, 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.322. URL https://aclanthology.org/2024.emnlp-main.322/.

- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 5776–5788. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. Videoagent: Long-form video understanding with large language model as agent. In *European Conference on Computer Vision*, pp. 58–76. Springer, 2024b.
- Yixuan Wang, Shiyu Ji, Yijun Liu, Yuzhuang Xu, Yang Xu, Qingfu Zhu, and Wanxiang Che. Lookahead q-cache: Achieving more consistent kv cache eviction via pseudo query. *ArXiv preprint*, abs/2505.20334, 2025. URL https://arxiv.org/abs/2505.20334.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=NG7sS51zVF.
- Peng Xu, Wei Ping, Xianchao Wu, Chejian Xu, Zihan Liu, Mohammad Shoeybi, and Bryan Catanzaro. Chatqa 2: Bridging the gap to proprietary llms in long context and RAG capabilities. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=cPD2hU35x3.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Xi Ye, Fangcong Yin, Yinghui He, Joie Zhang, Howard Yen, Tianyu Gao, Greg Durrett, and Danqi Chen. Longproc: Benchmarking long-context language models on long procedural generation. In Second Conference on Language Modeling, 2025. URL https://openreview.net/forum?id=ruWC5LIMSo.
- Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark W. Barrett, Zhangyang Wang, and Beidi Chen. H2O: heavy-hitter oracle for efficient generative inference of large language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/6ceefa7b15572587b78ecfcebb2827f8-Abstract-Conference.html.

A PSEUDO-CODE

702

703 704

705

706 707

731 732

755

The pseudocode for LookaheadKV training and eviction is described in Algorithm 1 and Algorithm 2, respectively.

Algorithm 1 LookaheadKV Training

```
708
                               Require: dataset \mathcal{D} of input-response pairs
709
                                   1: scores ← []
                                                                                                                                                                                                                                                                                                             710

⊳ score estimates using LookaheadKV

                                   2: estimates ← []
711
                                   3: for each training sample (X, Y) in dataset D do
712
                                                           for each layer l do
                                                                                                                                                                                                                                                                                                                                                              ⊳ GT pass
                                   4:
713
                                   5:
                                                                        for each head h in layer l do
714
                                   6:
                                                                                      S \leftarrow GT importance score for head (l, h)
                                   7:
                                                                                      scores.append(S)
715
                                                                         end for
                                   8:
716
                                                           end for
                                   9:
717
                               10:
                                                           for each layer l do

    b lookahead pass
    b lookahead pass
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
    c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
   c
718
                                                                        for each head h in layer l do
                               11:
719
                                                                                      \hat{S} \leftarrow \text{importance scores using lookahead embeddings for head } (l, h)
                               12:
720
                                                                                      estimates.append(\hat{S})
                               13:
721
                               14:
                                                                        end for
722
                                                           end for
                               15:
723
                               16:
                                                           L \leftarrow 0
                                                                                                                                                                                                                                                                                                                                            724
                                                           for all (S, \hat{S}) in scores, estimates do
                               17:
725
                                                                          L \leftarrow L + \mathrm{KL}(\mathrm{Norm}(S) \parallel \mathrm{Norm}(\hat{S}))
                               18:
726
                               19:
                                                           end for
727
                                                           L \leftarrow \frac{L}{|\text{scores}|}
                               20:
728
                               21:
                                                           L.backward()
729
                               22: end for
730
```

Algorithm 2 LookaheadKV Eviction

```
733
               Require: Input prompt X = (x_1, \ldots, x_{n_{in}})
734
                Require: cache budget k
735
                Require: learned lookahead tokens L = (l_1, \ldots, l_{n_{\text{lookahead}}})
736
                 1: \hat{X} \leftarrow (x_1, \dots, x_{n_{\text{in}}}, l_1, \dots, l_{n_{\text{lookahead}}})
                                                                                                                      > append learned lookahead tokens to input
737
                 2: Perform a forward pass with \hat{X} to populate KV cache
738
                 3: For each layer l and head h we now have
739
                                 K_{l,h} \in \mathbb{R}^{(n_{\text{in}} + n_{\text{lookahead}}) \times d}
740
                                 V_{l,h} \in \mathbb{R}^{(n_{\text{in}} + n_{\text{lookahead}}) \times d}
741
                 4: for each layer l do
742
                              \mathbf{for} \ \mathrm{each} \ \mathrm{head} \ h \ \mathrm{in} \ \mathrm{layer} \ l \ \mathbf{do}
                 5:
743
                                     Q_{l,h}^{\text{lookahead}} \leftarrow L_{l,h}^{\text{lookahead}} \ W_{l,h}^q \ + \ L_{l,h}^{\text{lookahead}} \ \Delta W_{l,h}^q
                                                                                                                                              > queries of lookahead tokens
                 6:
744
                                     \mathbf{A}_{l,h} \leftarrow \operatorname{Softmax}\!\left(\frac{Q_{l,h}^{\operatorname{lookahead}} K_{l,h}^{\top}}{\sqrt{d}}\right)
745
                 7:
746
                                    \mathbf{s}_{l,h} \leftarrow \frac{1}{n_{\text{lookahead}}} \sum_{i=1}^{n_{\text{lookahead}}} \mathbf{A}_{l,h}[i,:]
                 8:
                                                                                                                                                   \triangleright score vector of length n_{\rm in}
747
                 9:
                                    \mathcal{I}_{l,h} \leftarrow \text{TopK}(\mathbf{s}_{l,h}, k)
                                                                                                                                                             \triangleright select Top-k indices
748
                                     K_{l,h}^{\text{pruned}} \leftarrow K_{l,h}[\mathcal{I}_{l,h}]
                10:
749
                                     V_{l,h}^{\text{pruned}} \leftarrow V_{l,h}[\mathcal{I}_{l,h}]
750
                11:
                                     Cache (K_{l,h}^{\text{pruned}}, V_{l,h}^{\text{pruned}})
751
                12:
752
                13:
                              end for
753
                14: end for
               15: return Pruned KV cache \{(K_{l,h}^{\text{pruned}}, V_{l,h}^{\text{pruned}})\}
754
```

B THEORETICAL ESTIMATION DETAILS

This section details our methodology for theoretically estimation the Time-to-First-Token (TTFT) latency for various KV cache eviction algorithms. Our analysis is based on the analytical model for FLOPs and memory traffic proposed by Davies et al. (2025). To align configurations of theoretical estimates with them of actual measurements, we simulate the execution of LLaMA3.1-8B on a single NVIDIA H100 80GB GPU with a batch size of 1, assuming all weights and activations are in half-precision. We set KV cache budget size of 128, lookahead size as 32, and window size as 32. We only consider tensor operations which are dominant parts of the computations. To provide estimates that closely reflect real-world performance, our calculations incorporate practical hardware utilization by assuming a flops efficiency of 0.7 and a memory efficiency of 0.9, as described in Li (2023).

To isolate the specific overhead introduced by each eviction algorithm, we first establish a baseline by calculating the theoretical latency of a single forward pass. The TTFT overhead for each eviction method is then determined by subtracting this baseline forward pass latency from the method's total estimated TTFT. Notably, we do not add memory IO overhead incurred by KV cache unlike other eviction methods, since we only aim to calculate the computational overhead of a single forward pass operation. For LAQ, the total latency is calculated by summing the costs of its three costituent steps—the first eviction, low-cost generation of pseudo response, and the second eviction. Similarly, the total latency of SpecKV is estimated by aggregating the latencies of its draft prefill, draft decode, and target model eviction phases. A comprehensive implementation of the code to derive theoretical estimates of all baselines is available in the Suppplementary Materials.

C IMPLEMENTATION OPTIMIZATION

Efficient attention implementations such as FlashAttention (Dao et al., 2022) do not materialize the full attention score matrix, but is required in our setting to compute importance scores and enable gradient backpropagation. A possible solution is to compute the complete attention matrix using native PyTorch (i.e., eager attention), but this quickly leads to an out-of-memory error as the matrix size grows quadratically with the sequence length, which is incompatible with our training setting (upto 16K sequence length). Fortunately, for our objective, we only require the cross-attention scores between the generated response and the entire input sequence, and the response length is typically much shorter than the input prompt.

Leveraging this observation, we adopt the following approach: for the attention layers' forward computation, we use flash attention, while for the importance score computation and loss backpropagation, we employ eager attention to only compute the partial attention score matrix with the queries of model response and keys of the entire sequence. This reduces the memory requirement of eager attention from $\mathcal{O}((|X|+|Y|)^2)$ to $\mathcal{O}(|X|\cdot|Y|+|Y|^2)$, where |X| and |Y| denote the lengths of the input prompt and model response, respectively, with $|X|\gg |Y|$.

D NEED FOR DATA GENERATION

One of the requirements of LookaheadKV training is that the target model's generated responses must be available as training data. However, generating these responses from the model can sometimes be costly, e.g., when applying LookaheadKV across multiple models. Hence, to assess whether this requirement of can be relaxed, we evaluate an alternative setting where training uses the responses from the source datasets instead of model-generated outputs.

We observe in Figure 7 that this substitution leads to a relatively minor drop in average LongBench performance in lower-budget regimes. We hypothesize that if the attention distribution of the model-generated responses and that of the source dataset responses are moderately similar, our method can still successfully learn to accurately predict the importance scores. Overall, these results suggest that, in scenarios where training data generation is impractical, using source responses provides a viable and effective alternative.

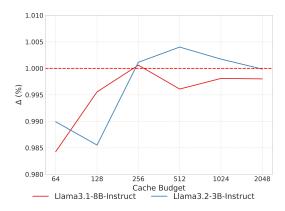


Figure 7: Performance ratio of training using model-generated data vs. source data.

E ADDITIONAL RESULTS

We provide a comprehensive experimental results excluded from the main text due the page limitation.

E.1 RESULTS ON LONGBENCH

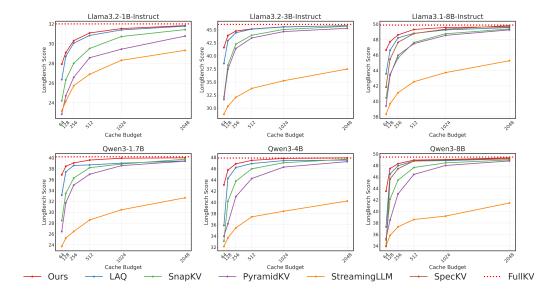


Figure 8: Full Longbench results across multiple cache budgets.

Table 4: Performance comparison of different methods across various LLMs on LongBench.

Fullix	_	LLMs	Single	e-Docume	nt QA	Mu	lti-Document Ç	QA .	S	ummarizati	on	F	ew-shot Lea	rning	Synti	netic	Co	ode	Avg.
StreamingLLM 24.95 21.50 32.56 50.67 42.89 24.31 18.49 21.25 18.19 40.50 85.57 38.28 7.50 99.50 59.03 49.72 39.08 59.00 59.03 49.72 39.08 59.00		LLMS	NrtQA	Qasper	MF-en	HotpotQA	2WikiMQA	Musique	GovReport	QMSum	MultiNews	TREC	TriviaQA	SAMSum	PCount	Pre	Lcc	RB-P	Avg.
Streaming LLM 24.95 21.50 32.56 50.67 42.89 24.31 18.49 21.25 18.19 40.50 85.57 38.28 7.50 99.50 50.91 49.72 98.68		FullKV	31.63	46.66	56.93	58.10	48.50	31.57				72.50	91.65	43.79	6.64	99.50	65.12	58.78	49.88
SampkV 29.13 28.06 51.23 56.79 45.30 27.81 19.99 23.03 19.73 46.00 89.72 40.44 7.50 99.50 59.08 51.06 43.06 58.06 59.00 56.76 46.11 28.13 19.86 22.81 20.03 44.50 88.14 39.73 75.0 99.50 59.84 51.06 43.06 58.06 58.06 59.00 24.24 24.20 21.59 60.00 59.00 41.04 72.5 59.50 61.11 61.38 45.45 58.06 59.00 22.51 24.25 24.20 21.59 60.00 20.02 42.14 8.83 99.50 61.11 61.38 45.45 67.00																			
PyramidKV 27.70 28.86 52.00 56.76 46.11 28.13 19.86 22.81 20.03 44.50 88.41 39.73 7.50 99.50 59.84 51.96 43.35 psckV 29.22 29.12 54.05 56.54 46.30 29.90 22.65 23.18 21.25 52.00 90.02 42.14 8.83 99.50 61.15 55.38 46.61 29.00 10.00 65.00 99.02 42.14 8.83 99.50 61.17 55.29 47.72 19.00 10.00																			
Figure F																			
Speck V 29.22 29.12 54.05 56.54 46.30 29.90 22.65 23.18 21.25 52.00 90.02 42.14 8.83 99.50 61.11 61.38 45.45																			
LookaheadKV 31,32 42,85 56,78 57,04 47,44 30,82 25,18 24,33 23,09 65,50 92,24 42,06 75,0 99,50 61,75 55,29 47,72																			
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	#																		
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	Ĕ	LookaneadKv	31.32	42.65	30.78	37.04	47.44	30.62				05.50	92.24	42.90	7.30	99.30	01./5	33.29	47.72
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	E.							** **											
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	ф																		
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	~																		
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	33																		
Cookaneadky S196 44.01 50.80 57.99 47.41 51.46 27.25 24.50 24.59 69.00 92.55 42.59 67.50 93.50 63.15 55.73 43.75 58.75 43.75 48.75	E																		
StreamingLLM 27.23 30.80 36.64 50.59 43.26 23.45 23.75 21.67 25.49 63.50 88.84 42.56 7.50 93.50 63.15 55.73 43.73	3																		
StreamingLLM 27:23 30.80 36.64 50.59 43.26 23.45 25.73 21.67 25.84 69.50 92.04 42.78 70.8 99.50 64.57 58.46 48.76 PyramidkV 30.79 44.91 56.65 58.13 48.17 30.56 26.65 24.53 25.88 68.00 91.78 42.20 6.83 99.50 64.47 57.77 48.55		2500Kuncuuze (5100	74101	20.00	31.33	17.11	22110				0,100	72100	1200	7.00	100100	02.01	37.02	10101
SampKV 29.64 44.60 57.30 57.62 48.31 31.18 27.57 24.17 25.84 69.50 92.04 42.78 7.08 99.50 64.57 58.46 48.76		Strooming I I M	27.22	20.90	26.64	50.50	42.26	22.45				62.50	00 04	12.56	7.50	02.50	62.15	55 72	42.72
PyramidKV 30.79 44.91 56.65 58.13 48.17 30.56 26.65 24.53 25.88 68.00 91.78 42.20 6.83 99.50 64.41 57.77 48.55 EACH Charles and the control of the control o																			
L\(\text{AC}\) 31.63 & 45.63 & 55.02 & 57.70 & 50.27 & 31.28 & 28.82 & 25.10 & 26.18 & 72.50 & 92.33 & 43.31 & 6.50 & 100.00 & 62.75 & 59.04 & 49.25 \\ \text{SpecKV}\) 31.59 & 45.44 & 57.98 & 57.5 & 49.16 & 31.95 & 28.67 & 24.95 & 25.77 & 67.50 & 92.23 & 43.94 & 6.60 & 99.50 & 65.21 & 62.30 & 49.36 \\ \text{LookaheadKV}\) 31.14 & 46.04 & 57.77 & 58.22 & 48.43 & 30.72 & 30.75 & 25.31 & 26.66 & 72.50 & 91.92 & 43.39 & 7.08 & 100.00 & 64.87 & 58.36 & 49.57 \\ \text{FullKV}\) 26.04 & 47.76 & 53.33 & 59.23 & 43.37 & 36.05 & 33.66 & 24.05 & 24.79 & 71.50 & 90.21 & 44.43 & 2.00 & 100.00 & 69.39 & 65.27 & 49.46 \\ \text{StreamingLLM}\) 17.65 & 26.69 & 28.40 & 41.05 & 33.46 & 20.82 & 15.72 & 19.15 & 15.14 & 43.00 & 82.57 & 38.44 & 1.50 & 70.00 & 62.86 & 56.69 & 35.82 \\ \text{SpapKV}\) 91.44 & 32.65 & 45.99 & 54.81 & 38.95 & 26.59 & 17.66 & 20.83 & 16.04 & 49.50 & 87.10 & 38.90 & 3.50 & 99.50 & 64.62 & 88.29 & 42.13 \\ \text{PyramidKV}\) 15.57 & 30.19 & 41.84 & 46.01 & 35.73 & 19.57 & 16.51 & 19.67 & 14.86 & 47.00 & 83.51 & 35.56 & 2.50 & 92.00 & 62.14 & 53.07 & 38.48 \\ \text{LookaheadKV}\) 20.04 & 42.15 & 53.55 & 57.89 & 42.84 & 36.74 & 21.33 & 22.25 & 18.34 & 64.50 & 89.55 & 40.93 & 30.01 & 100.00 & 66.74 & 61.70 & 46.52 \\ \text{SpecKV}\) 23.03 & 37.14 & 53.58 & 56.77 & 42.24 & 31.82 & 21.33 & 22.86 & 19.04 & 60.00 & 88.31 & 41.50 & 3.50 & 100.00 & 66.74 & 61.70 & 46.52 \\ \text{StreamingLLM}\) 18.18 & 28.53 & 28.82 & 42.81 & 33.58 & 21.34 & 18.63 & 19.05 & 40.00 & 88.31 & 41.05 & 3.50 & 100.00 & 66.55 & 59.41 & 37.32 \\ \text{StreamingLMV}\) 23.03 & 38.32 & 51.04 & 57.36 & 40.67 & 32.82 & 21.13 & 18.89 & 18.97 & 59.50 & 89.46 & 41.06 & 2.00 & 100.00 & 65.50 & 59.41 & 37.32 \\ \text{StreamingLMV}\) 23.03 & 38.32 & 51.04 & 57.36 & 40.67 & 32.82 & 21.51 & 21.89 & 18.97 & 59.50 & 89.46 & 41.06 & 2.00 & 100.00 & 65.50 & 59.41 & 37.32 \\ \text{StreamingLMV}\) 24.26 & 46.13 & 52.48 & 58.84 & 42.99 & 34.53 & 36.62 & 24.22 & 23.38 & 20.38 & 70.00 & 89.05 & 42.47 & 3.00 & 100.00 & 66.51 & 57.34 & 40.00 \\ \text{StreamingLMV}																			
SpeckV 31.59 45.44 57.98 57.51 49.16 31.95 28.67 24.95 25.77 67.50 91.23 43.94 6.00 99.50 65.21 62.30 49.35																			
FullKV 26.04 47.76 53.33 59.23 43.37 36.05 33.66 24.05 24.79 71.50 90.21 44.43 2.00 100.00 69.39 65.57 49.46 StreamingLLM 17.65 26.69 28.40 41.05 33.346 20.82 15.72 13.03 15.14 43.00 82.57 38.44 1.50 70.00 62.86 56.69 35.82 SmpKV 19.14 32.65 45.99 54.81 38.95 26.59 17.66 19.96 19.60 49.50 87.10 38.90 35.0 99.50 64.62 38.29 42.13 PyramidKV 19.14 32.65 45.99 41.84 40.10 35.73 99.71 17.66 19.96 19.96 47.00 83.51 35.50 2.50 92.00 62.42 38.29 42.13 LookaheadKV 26.06 44.30 53.24 58.78 42.29 35.89 22.92 22.92 13.66.50 88.95 41.64 35.00 99.50 65.52 62.88 47.40 StreamingLLM 18.18 28.53 28.52 42.81 33.58 21.34 18.63 19.20 17.67 48.00 88.51 41.50 35.00 100.00 66.82 61.05 45.62 FyramidKV 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 88.20 41.06 20.0 100.00 65.55 59.41 37.32 FyramidKV 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 88.20 41.06 20.0 100.00 65.76 61.88 45.45 Expression LLM 18.18 28.53 28.52 42.81 33.58 21.34 18.63 19.20 17.76 48.00 88.51 41.00 69.00 65.50 59.41 37.32 FyramidKV 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 88.20 38.98 3.50 100.00 65.76 61.88 45.45 Expression LLM 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 88.20 38.98 3.50 100.00 65.75 59.14 37.32 Expression LLM 18.18 28.53 28.52 42.42 34.50 24.53 23.64 21.25 68.00 88.13 43.12 3.00 100.00 65.76 61.88 45.45 Expression LLM 21.25 32.82 31.44 45.94 34.38 23.34 23.53 23.64 21.25 68.00 88.13 43.12 3.00 100.00 65.76 61.88 45.45 Expression LLM 21.25 32.82 31.44 45.94 34.38 23.34 23.53 23.64 21.25 68.00 88.13 43.12 3.00 100.00 67.83 64.07 47.25 Expression LLM 21.25 32.82 31.44 45.94 34.38 23.34 23.34 23.34 22.35 23.50 20.55 23.50 62.00 88.71 41.18 0.50 44.00 68.39 63.65 39.19 Expression LLM 21.25 32.82 31.44 45.94 34.38 23.34 23.34 23.34 23.36 20.38 70.00 89.55 43.13 2.00 100.00 67.83 64.74 42.74 23.75 20.55 20.55 23.50 62.00 89.55 43.13 2.00 100.00 67.85 64.75 48.00 Expression LLM 21.25 32.82 31.44 45.94 34.38 23.34 23.34 23.36 20.38 23.60 89.55 43.13 2.00 100.00 67.85 64.75 48.00																			
StreamingLLM 17.65 26.69 28.40 41.05 33.46 20.82 15.72 19.15 15.14 43.00 82.57 38.44 1.50 70.00 62.86 56.69 35.82		LookaheadKV	31.14	46.04	57.77	58.22	48.43	30.72	30.75	25.31	26.66	72.50	91.92	43.39	7.08	100.00	64.87	58.36	49.57
StreamingLLM 17.65 26.69 28.40 41.05 33.46 20.82 15.72 19.15 15.14 43.00 82.57 38.44 1.50 70.00 62.86 56.69 35.82		FullKV	26.04	47.76	53.33	59.23	43.37	36.05	33.66	24.05	24.79	71.50	90.21	44.43	2.00	100.00	69.39	65.57	49.46
SapaKV 19.14 32.65 45.99 54.81 38.95 26.59 17.66 20.83 16.04 49.50 87.10 38.90 3.50 99.50 64.62 58.29 42.13																			
PyramidKV 15.57 30.19 41.84 46.01 35.73 19.57 16.51 19.67 14.86 47.00 83.51 35.56 2.50 92.00 62.14 53.07 38.48 LAQ 22.74 42.15 53.55 57.89 42.84 36.74 21.33 22.25 18.34 64.50 89.55 40.93 30.0 100.00 66.74 61.70 46.52 SpecKV 23.03 37.14 53.58 56.77 42.24 31.82 21.33 22.86 19.04 60.00 88.31 41.50 3.50 100.00 66.74 61.70 46.52 LookaheadKV 26.06 44.30 53.24 58.78 42.79 35.89 52.29 22.95 21.13 66.50 88.95 41.64 3.50 99.50 65.95 62.88 47.46 \$		StreamingLLM				41.05	33.46	20.82	15.72	19.15			82.57				62.86		
LAQ 22.74 42.15 53.55 57.89 42.84 36.74 21.33 22.25 18.34 64.50 89.55 40.93 3.00 100.00 66.74 61.70 46.52 SpecKV 23.03 37.14 53.85 56.77 42.24 31.82 21.33 22.85 19.04 60.00 88.31 41.50 3.50 100.00 66.82 61.96 45.02 LookaheadKV 26.06 44.30 53.24 58.78 42.79 35.89 25.29 22.95 21.13 66.50 88.95 41.64 3.50 99.50 65.95 62.88 47.46 StreamingLLM 81.8 28.53 28.52 42.81 33.58 21.34 18.63 19.20 17.76 48.00 85.58 40.08 1.00 69.00 65.50 59.41 37.32 SpecKV 23.03 38.32 51.04 57.36 40.67 32.82 21.51 21.89 18.97 59.50 89.46 41.06 2.00 100.00 65.50 59.41 37.32 PyramidKV 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 85.20 38.98 3.50 100.00 65.51 57.32 43.05 SpecKV 22.58 41.09 53.89 59.85 42.42 34.50 24.22 23.38 20.38 70.00 89.05 42.47 3.00 100.00 68.37 64.0 47.42 LookaheadKV 25.88 45.40 52.68 58.47 44.05 36.13 27.77 23.71 22.88 69.00 89.05 43.32 2.00 100.00 66.39 64.47 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.05 43.32 2.00 100.00 68.39 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.52 48.00 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.52 48.00 StreamingLLM 21.25 32.82 53.01 58.86 42.32 33.43 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.52 48.00 StreamingLLM 21.25 32.82 53.01 58.86 42.32 33.43 27.77 23.71 22.88 69.00 89.55 43.13 2.00 100.00 68.39 64.52 48.00 StreamingLLM 21.27 33.45 59.10 43.35 37.26 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 6																			
SpecKV 23.03 37.14 53.58 56.77 42.24 31.82 21.33 22.86 19.04 60.00 88.31 41.50 3.50 100.00 66.82 61.96 45.05																			
CookaheadKV 26.06 44.30 53.24 58.78 42.79 35.89 25.29 22.95 21.13 66.50 88.95 41.64 3.50 99.50 65.95 62.88 47.46																			
StreamingLLM 18.18 28.53 28.52 42.81 33.58 21.34 18.63 19.20 17.76 48.00 85.58 40.08 1.00 69.00 65.50 59.41 37.32 SuprKV 23.03 38.32 51.04 57.36 40.67 32.82 21.51 21.89 18.97 59.50 89.46 41.06 2.00 100.00 67.62 61.88 45.45 PyramidKV 18.47 34.87 47.44 55.68 37.89 26.67 20.43 20.92 17.43 58.50 85.20 38.98 3.50 100.00 65.51 57.32 43.05 SpeckV 22.58 41.09 53.89 59.85 42.42 24.34 34.50 24.53 23.64 21.25 68.00 88.13 43.12 3.00 100.00 68.37 64.0 47.42 LookaheadKV 25.88 45.40 25.68 58.47 44.05 36.13 27.73 23.05 23.54 21.25 68.00 88.13 43.12 3.00 100.00 67.83 64.17 48.11 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 25.73 20.25 23.50 62.00 88.71 41.18 0.50 44.00 68.39 63.65 39.19 SmapKV 24.26 46.13 52.48 58.52 42.66 36.89 28.39 23.61 23.33 69.00 89.55 43.13 2.00 100.00 68.37 64.57 SpeckV 24.98 46.56 54.00 57.01 43.52 35.46 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 68.81 48.48 SpeckV 24.98 46.56 54.00 59.00 43.32 37.46 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 69.97 66.33 48.98 SpeckV 24.98 46.56 54.00 59.00 43.33 23.48 23.34 23.88 23.88 23.88																			
StreamingLLM 18.18 28.53 28.52 42.81 33.58 21.34 18.63 19.20 17.76 48.00 85.58 40.08 1.00 69.00 65.50 59.41 37.32 37.58		LookaneadKV	26.06	44.30	55.24	58.78	42.79	35.89				66.50	88.95	41.64	3.50	99.50	65.95	62.88	47.46
SpecKV 22.58 41.09 53.89 59.85 42.42 34.50 34.50 24.22 23.56 20.35 70.00 89.05 42.47 3.00 100.00 68.39 64.04 47.42	m																		
SpecKV 22.58 41.09 53.89 59.85 42.42 34.50 34.50 24.22 23.56 20.35 70.00 89.05 42.47 3.00 100.00 68.39 64.04 47.42	99																		
SpecKV 22.58 41.09 53.89 59.85 42.42 34.50 34.50 24.22 23.56 20.35 70.00 89.05 42.47 3.00 100.00 68.39 64.40 47.42	'en																		
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	á																		
LookaheadKV 25.88 45.40 52.68 58.47 44.05 36.13 27.77 23.71 22.88 69.00 89.05 43.32 2.00 100.00 67.83 64.71 48.31 StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 25.73 20.25 23.50 62.00 88.71 41.18 0.50 44.00 68.39 63.65 39.19 SpramidKV 23.77 42.89 53.01 58.86 42.32 35.47 72.32 23.07 22.72 71.00 89.95 43.13 2.00 100.00 68.39 63.65 39.19 PyramidKV 23.77 42.89 53.01 58.86 42.32 35.47 27.22 27.10 89.95 43.13 2.00 100.00 68.16 45.25 48.00 LQ 26.11 47.27 53.45 57.01 43.52 37.26 29.50 23.88 23.47 71.50 89.63 44.00 2.00 <td>-</td> <td></td>	-																		
StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 25.73 20.25 23.50 62.00 88.71 41.18 0.50 44.00 68.39 63.65 39.19 SnapKV 24.26 46.13 52.48 58.52 42.66 36.89 82.89 23.61 23.33 69.00 89.55 43.13 2.00 100.00 69.05 66.27 48.45 PyramidkV 23.77 42.89 53.01 58.86 42.32 35.47 27.32 23.07 22.72 71.00 89.95 42.56 2.00 100.00 68.81 64.25 48.00 LAQ 26.11 47.27 53.45 57.01 43.52 372.66 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 69.09 66.83 48.84 SpeckV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 3.00 100.00 69.09 66.33 48.98																			
StreamingLLM 21.25 32.82 31.44 45.94 34.38 23.34 25.73 20.25 23.50 62.00 88.71 41.18 0.50 44.00 68.39 63.65 39.19 SnapKV 24.26 46.13 52.48 58.52 42.66 36.89 28.39 23.61 23.33 69.00 89.55 43.13 2.00 100.00 69.05 66.27 48.45 PyramidKV 23.77 42.89 53.01 58.86 42.52 33.47 27.32 23.07 22.72 71.00 89.95 42.56 2.00 100.00 69.05 66.27 48.45 LAQ 26.11 47.27 53.45 57.01 43.52 37.26 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 69.09 66.38 44.25 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 30.00 100.00 69.09 66.33 48.98 SpecKV 24.98 46.56 54.07 59.04 48.35 48.00 48		Lookancaukv	25.00	75.70	32.00	30.47	44.05	50.15				07.00	07.05	45.52	2.00	100.00	07.05	04.71	40.51
SapkV 24.26 46.13 52.48 58.52 42.66 36.89 28.39 23.61 23.33 69.00 89.55 43.13 2.00 100.00 69.05 66.27 48.45		Strooming I I M	21.25	22.92	21 44	45.04	24.20	22.24				62.00	00 71	41 10	0.50	44.00	69 20	62.65	20.10
PyramidKV 23.77 42.89 53.01 58.86 42.32 35.47 27.32 23.07 22.72 71.00 89.95 42.56 2.00 100.00 68.81 64.25 48.00 LAQ 26.11 47.27 53.45 57.01 43.52 37.26 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 67.94 64.83 48.84 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 3.00 100.00 69.09 66.53 48.98																			
LÁQ 26.11 47.27 53.45 57.01 43.52 37.26 29.50 23.88 23.47 71.50 89.63 44.00 2.00 100.00 67.94 64.83 48.84 SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 3.00 100.00 69.09 66.53 48.98																			
SpecKV 24.98 46.56 54.07 59.04 43.37 34.12 29.32 24.18 23.68 71.00 90.11 44.56 3.00 100.00 69.09 66.53 48.98																			

E.2 RESULTS ON RULER

We report the RULER results across all six models tested, with cache budget settings at 64 (Figure 9) and 128 (Figure 10).

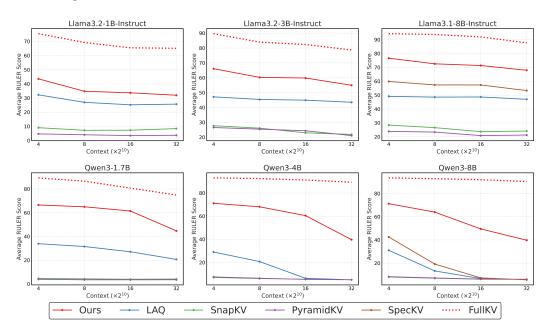


Figure 9: Full RULER results across context lengths (budget = 64)

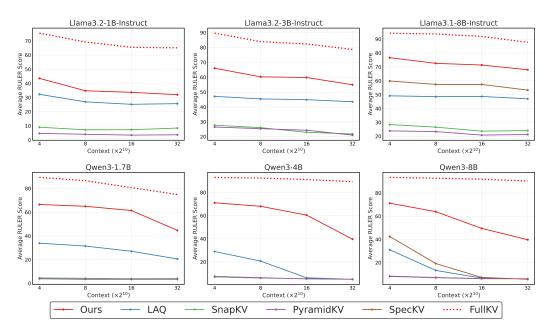


Figure 10: Full RULER results across context lengths (budget = 128)

E.3 Additional Efficiency Analysis

We show the full results of the latency analysis that were omitted in the main paper due to space limitation in this section. Note that the empirical TTFT overheads for some methods, SnapKV in particular, can be larger than theoretical estimations. These are probably due to inefficient implementation of these methods in KVCache-Factory or their official implementation. Better implementations may reduce these overheads significantly, more in line with the theoretical cost.

Table 5: Theoretical and Practical Analysis across various context lengths and methods.

			Theoretica	En	pirical Cost			
Context Length	Method	Compute (TFLOPs)	Memory Traffic (GB)	TTFT (ms)	TTFT Overhead (ms)	TTFT (ms)	TTFT Overhead (ms)	
	Forward Pass Only	60	13	113	N/A	130	N/A	
	LookaheadKV (ours)	60	13	114	0.92	141	11.38	
4096	SnapKV	60	13	113	0.01	143	13.14	
	SpecKV	70	77	165	52.10	223	92.42	
	LAQ	61	444	347	233.81	637	506.58	
	Forward Pass Only	136	13	257	N/A	291	N/A	
8192	LookaheadKV (ours)	137	13	258	1.03	302	10.88	
	SnapKV	136	13	257	0.01	311	20.17	
	SpecKV	159	81	337	79.53	411	120.51	
	LAQ	137	445	492	234.59	800	509.38	
	Forward Pass Only	336	13	635	N/A	658	N/A	
	LookaheadKV (ours)	337	13	636	1.27	677	18.50	
16384	SnapKV	336	13	635	0.01	695	37.12	
	SpecKV	398	89	792	157.05	866	207.31	
	LAQ	337	447	871	236.15	1182	523.54	
	Forward Pass Only	928	13	1754	N/A	1760	N/A	
	LookaheadKV (ours)	929	13	1755	1.74	1798	38.04	
32768	SnapKV	928	13	1754	0.01	1838	77.67	
	SpecKV	1115	106	2156	402.80	2263	502.87	
	LAQ	930	406	1993	239.26	2314	553.68	

F HYPER-PARAMETERS

Training hyper-parameters.

Learning rate was searched for Llama and Qwen model family among $[5 \times 10^{-5}, 1 \times 10^{-4}, 2 \times 10^{-4}, 1 \times 10^{-3}]$. The final hyper-parameters for all experiments are shown in Table 6.

Table 6: Training hyperparameters.

Parameters	Values
Optimizer	Adam
β_1, β_2	0.9, 0.95
Effective Batch Size	32
Drop-out (p)	0.0
Max Sequence Length	16384 (prompt length) + 512 (response length)
Train Iters	7600
Learning rate	1×10^{-3} (for Llama), 2×10^{-4} (for Qwen)
Schedule	Cosine
Warmup steps	2%
Min LR	0.0
Gradient clipping	1.0

Eviction hyper-parameters. We use the implementations in KVCache-Factory or their official implementations (SpecKV) for all our methods, except for LAQ which we re-implement ourselves. Following prior works (Li et al., 2024; Cai et al., 2024; Galim et al., 2025), we use standard configuration settings for all baseline methods, including an observation window size of 32, maxpooling kernel size of 7, and mean reduction for GQA compatibility (Feng et al., 2024). For LookaheadKV we use the same settings, except we do not use window size, as our method does not train with the suffix window for prediction. Further, since our lookahead size $n_{\rm lookahead}$ is 32, we set the maximum generation limit of LAQ and SpecKV to 32 tokens so that the methods can be compared using the same number of draft tokens.

G DATASETS, BENCHMARKS, AND SOFTWARE

Software Our source code is available in the supplementary, and our implementation is built on KVCache-Factory.

Training Dataset Our training dataset mixture consist of random samples from publicly available datasets: 50K long_sft subset of ChatQA2-Long-SFT-data, 20K subset of tulu-3-sft-olmo-2-mixture, 7K samples from The Stack, and 3K samples from MetaMathFewshot, HellaSwag_DPO_Fewshot, and ARC_DPO_Fewshot, respectively.

Evaluation Benchmarks We used LongBench dataset as fetched and processed by KVCache-Factory, see HF Dataset for the official source. For RULER, we used RULER Github. For LongProc, we used LongProc Github.

H LLM USAGE

LLM assistants were used to refine the wording of selected sentences, while the majority of the text was written by human. All LLM-generated text was carefully inspected to ensure that it contained no harmful or controversial content. Additionally, we used LLMs to help in finding some of the related literature discussed in the paper.