

STREAMING VIDEO QUESTION-ANSWERING WITH IN-CONTEXT VIDEO KV-CACHE RETRIEVAL

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose **ReKV**, a novel, training-free approach that integrates seamlessly with existing Video Large Language Models (Video-LLMs) to enable efficient streaming video question-answering (StreamingVQA). Traditional VideoQA systems struggle with long videos, as they must process the entire video before responding to queries, and repeat this process for each new question. In contrast, our approach analyzes long videos in a streaming fashion, allowing for prompt responses as soon as user queries are received. Building on a common Video-LLM, we first incorporate a sliding-window attention mechanism, ensuring that input frames attend to a limited number of preceding frames, thereby reducing computational overhead. To prevent information loss, we store processed video key-value caches (KV-Caches) in RAM and disk, reloading them into GPU memory as needed. Additionally, we introduce a retrieval method that leverages an external retriever or the parameters within Video-LLMs to retrieve only query-relevant KV-Caches, ensuring both efficiency and accuracy in question answering. ReKV enables the separation of video analyzing and question-answering across different processes and GPUs, significantly enhancing the efficiency of StreamingVQA. Through comprehensive experimentation, we validate the efficacy and practicality of our approach, which significantly boosts efficiency and enhances applicability over existing VideoQA models.

1 INTRODUCTION

In the literature, video understanding tasks, such as action recognition (Caba Heilbron et al., 2015; Goyal et al., 2017; Kay et al., 2017), visual object tracking (Huang et al., 2019; Muller et al., 2018), and video question-answering (Xu et al., 2017; Jang et al., 2017; Xiao et al., 2021; Li et al., 2024c), have primarily focused on short clips lasting from a few seconds to minutes. However, as vision models increasingly find applications in real-world scenarios like robotics, surveillance, and live broadcasts, the research in the vision community has gradually shifted towards understanding continuous video streams, where long-term contexts and real-time interaction are crucial.

In this paper, we consider the problem of **streaming video question-answering (StreamingVQA)**. As shown in Figure 1(a), it involves continuously processing long video streams and promptly responding to queries about the visual content at any moment. It can be treated as a generalization of the standard offline VideoQA, where the model processes the entire video and all questions simultaneously. By definition, such task of StreamingVQA presents three core challenges: (i) **Efficient Video Encoding**: Unlike traditional offline VideoQA, where models have access to the entire video clip, StreamingVQA demands real-time analysis of continuous streams. Models must efficiently process incoming frames without access to future frames or frequent revisiting of distant past frames. (ii) **Video Context Preservation**: To accurately answer questions posed later in the stream, models must preserve relevant information from earlier frames, making long-term context retention a key challenge. (iii) **Real-Time Response**: The model must provide accurate answers with minimal delay, requiring efficient retrieval of video context and rapid question-answering.

Current Video-LLMs often struggle to encode long video streams due to the large volume of video tokens, forcing most models to process only a sparse subset of frames (Maaz et al., 2024; Zhang et al., 2024c; Li et al., 2024a). This results in limited video lengths or a significant loss of fine-grained visual information. While techniques like average pooling (Li et al., 2024d) and memory

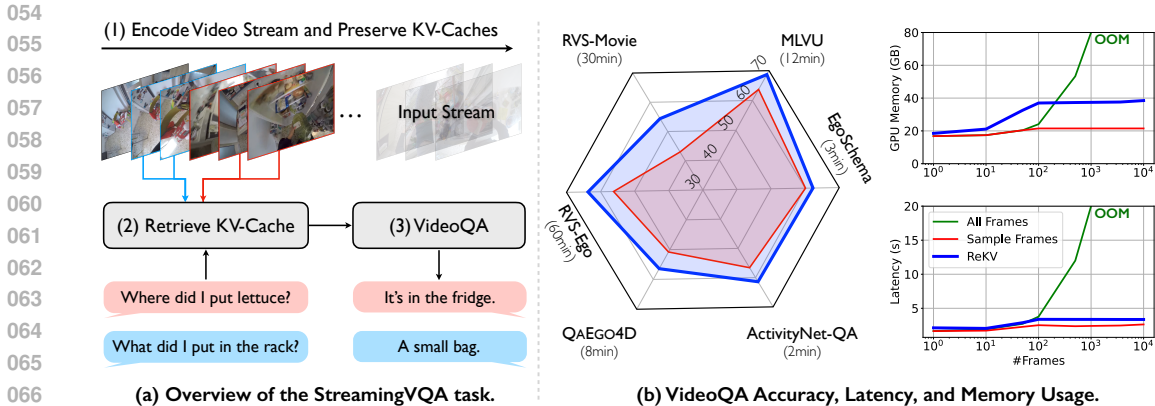


Figure 1: **Overview of the StreamingVQA task and our proposed ReKV.** (a) StreamingVQA requires a model to continuously process video streams and answer questions about previously viewed content at any moment. (b) We propose ReKV to enhance efficiency and accuracy in StreamingVQA. Tested with LLaVA-OV-7B on an H800 (80GB) GPU, ReKV maintains stable latency and GPU memory usage, preventing out-of-memory (OOM) errors as frames increase. It also improves the accuracy on six long-form VideoQA benchmarks compared to the uniform sampling baseline. Further details are provided in Section 3.

compression (Wu et al., 2022; Wang et al., 2023; He et al., 2024; Zhang et al., 2024a; Qian et al., 2024) reduce token volume, they come at the cost of losing details, particularly in temporal and lower-level visual features that are essential for complex question answering.

To address the challenges, we propose **ReKV (Retrieve In-context Video KV-Cache)**, a framework that integrates seamlessly with existing Video-LLMs (Maaz et al., 2024; Zhang et al., 2024c; Li et al., 2024a) without additional training. Our method employs two strategies for aggregating both short- and long-term temporal information. For short-term temporal context, the model adopts causal attention with a sliding-window mechanism (Han et al., 2023), where tokens attend only to a limited set of preceding tokens during encoding. For recalling long-term information, we enable dynamic access to any point within the video sequence via retrieval. Specifically, our method retains and reuses past computations (KV-Cache) to avoid redundant processing while enhancing long-term reasoning without sacrificing detail. For extremely long videos, KV-Caches can be offloaded to RAM or disk to prevent memory overflow.

To ensure real-time and accurate responses, we retrieve a fixed number of KV-Caches relevant to the current question. This design strikes a balance between efficiency and accuracy by avoiding the need to process all past frames, while still accessing the most critical information. We experimented with two retrieval methods: one using external CLIP-like models (Radford et al., 2021; Zhai et al., 2023) for semantic matching, and another leveraging internal attention weights for faster, more integrated, and potentially stronger retrieval (Xiao et al., 2024a; Li et al., 2024e).

In summary, ReKV efficiently encodes long video streams, preserves and retrieves in-context KV-Caches to address complex video question-answering. In addition, ReKV separates video encoding from question-answering into distinct processes, further enhancing efficiency. As shown in Figure 1(b), ReKV improves VideoQA accuracy while maintaining stable inference latency and memory usage as frames increase. The remainder of the paper is organized as follows: Section 4 provides an overview of the relevant literature. Section 2 formulates the StreamingVQA task and describes our proposed method in detail. In Section 3, we present ablation studies and comparisons to validate our approach. Consequently, our approach not only enhances accuracy on long VideoQA benchmarks, including MLVU (Zhou et al., 2024a), QAEGO4D_{MC} (Di & Xie, 2024), EgoSchema (Mangalam et al., 2023), and ActivityNet-QA (Yu et al., 2019), as well as StreamingVQA benchmarks (Zhang et al., 2024a), but also reduces inference latency and memory usage.

2 METHOD

This paper considers the problem of streaming video question-answering (StreamingVQA), where a model continuously processes a video stream and can respond to questions about past visual content at any moment. In Section 2.1, we formally define the task and outline the design principles for our proposed solution. Next, Section 2.2 introduces ReKV, an approach that integrates seamlessly with a Video-LLM to enable efficient StreamingVQA without requiring additional training. Overall, ReKV efficiently encodes the video stream, maintains its KV-Caches, retrieves relevant caches based on the given question, and uses them for accurate question-answering.

2.1 TASK DEFINITION AND DISCUSSION

Given a video stream $\mathcal{V}^T := [v_1, v_2, \dots, v_T]$ consisting of T frames and a set of N questions $\mathcal{Q} := \{q_1, q_2, \dots, q_N\}$, **StreamingVQA** is to answer a question q_i at any time step t ($1 \leq t \leq T$) using only the frames seen up to that point, $\mathcal{V}^t := [v_1, v_2, \dots, v_t]$.

Discussion: StreamingVQA vs. OfflineVQA. StreamingVQA involves continuously analyzing an incoming video stream and answering questions based on the observed visual content at any moment. In contrast, conventional video question-answering models Yang et al. (2022); Maaz et al. (2024); Zhang et al. (2024c); Li et al. (2024a) operate in an offline mode, referred to as OfflineVQA. The two paradigms differ in that: 1) StreamingVQA processes a continuous video stream, while OfflineVQA handles a predefined video input, and 2) StreamingVQA allows questions to be asked at any point during the stream, whereas OfflineVQA processes questions only after the entire video has been viewed. Notably, OfflineVQA can be considered a special case of StreamingVQA, where all questions are posed after the video is fully processed.

Conventional approaches typically employ a visual encoder Radford et al. (2021); Zhai et al. (2023); Fang et al. (2023) and a projection module Zhang et al. (2024c); Li et al. (2023) to process video frames \mathcal{V}^t . The output is concatenated with the tokenized question to form a sequence $[\mathcal{V}_t, q_i]$ ¹, which is then passed to an LLM Decoder to predict an answer. However, this approach is impractical due to the high computational cost associated with processing a large number of frames (T).

A common workaround is sparse frame sampling (Maaz et al., 2024; Zhang et al., 2024c; Li et al., 2024a), but this introduces new problems: (i) loss of critical visual information, leading to incomplete or inaccurate responses, and (ii) the need to reprocess frames for different questions, since questions asked at different time points require distinct frame samples. This becomes increasingly inefficient as T and N grow.

Given these challenges, current OfflineVQA methods fall short when applied to StreamingVQA scenarios. Therefore, designing a new approach optimized for StreamingVQA is crucial to handling video streams more efficiently, enabling real-time question answering and unlocking more interactive video analysis applications.

Discussion: Design Principles for Efficient StreamingVQA. To tackle the aforementioned challenges, we can exploit the causal nature of the LLM decoder to avoid redundant computations and strike a balance between accuracy and speed. During attention calculations, causal masking prevents the model from accessing future tokens, ensuring that video tokens are encoded independently of the questions. This allows us to **decouple video encoding from question-answering**.

For video encoding, we leverage the KV-Cache optimization to accelerate inference. However, as t grows large, handling the massive number of video tokens becomes increasingly inefficient and may exceed the model’s capacity Han et al. (2023); Xiao et al. (2024b). To address this, we adopt a sliding-window attention mechanism (Han et al., 2023), which limits the attention scope to only the most recent frames.

Regarding question-answering, Video KV-Caches are stored and can be reused as context to answer different questions. However, long video sequences produce a substantial amount of KV-Caches, leading to excessive GPU memory consumption, computational overhead, and unnecessary distractions if all are used. To address this, we introduce an efficient retrieval method that selects the most

¹We maintain the original notation for simplicity.

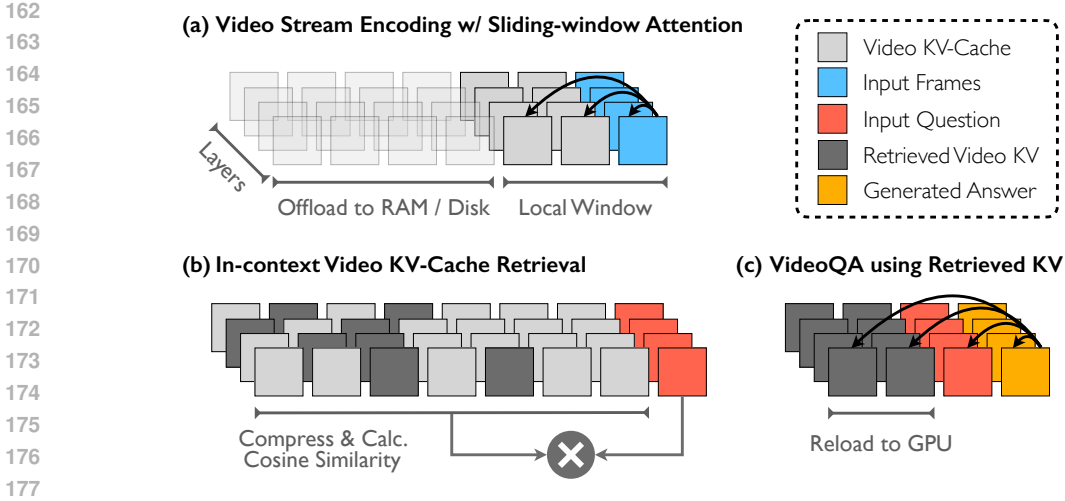


Figure 2: **Overview of ReKV.** We modify the attention mechanism in Decoder-based Video-LLMs: (a) The video stream is encoded with sliding-window attention (Equation 1), with out-of-window Video KV-Caches offloaded to RAM or disk. (b) Upon receiving a question, relevant key-value vectors are retrieved based on cosine similarity, with compressed vectors to accelerate retrieval (Equation 2). (c) The retrieved key-value vectors are reloaded onto the GPU and utilized for autoregressive answer generation (Equation 3).

relevant video key-value vectors from the video KV-Caches. These selected vectors then serve as context, enabling efficient and scalable StreamingVQA.

2.2 REKV: RETRIEVE IN-CONTEXT VIDEO KV-CACHE

We aim to enable Video-LLMs, trained on limited frames, to perform StreamingVQA without additional training. As illustrated in Figure 2, the proposed **ReKV** consists of three components: video stream encoding, video KV-Cache retrieval, and question-answering using the retrieved key-value vectors.

Video Stream Encoding with Sliding-window Attention. We encode the video stream \mathcal{V}^T incrementally, processing it chunk by chunk. At each step, the inputs include past key-value vectors $\mathbf{P} = \{(\mathbf{k}_j, \mathbf{v}_j)\}_{j=1}^{l_P}$ and the current tokens $\mathbf{X} = \{\mathbf{t}_{i+l_P}\}_{i=1}^{l_X}$, where l_P denotes the lengths of past key-values, and l_X refers to the chunk size. The local key-value vectors within a window l_L can thus be derived as $\mathbf{L} = \mathbf{P}_{[l_P-l_L+1:l_P]}$. The attention calculation is then formulated as:

$$\mathbf{O} = \text{Attn}(\mathbf{W}_Q \mathbf{X}, [\mathbf{L}_k, \mathbf{W}_K \mathbf{X}], [\mathbf{L}_v, \mathbf{W}_V \mathbf{X}]), \quad (1)$$

where \mathbf{W}_Q , \mathbf{W}_K , and \mathbf{W}_V are the attention layer parameters, \mathbf{L}_k and \mathbf{L}_v correspond to the key and value vectors in \mathbf{L} . All video KV-Caches are stored for future retrieval. Note that, for extremely long videos, we offloaded KV-Caches to RAM or disk to manage memory constraints, as in Xiao et al. (2024a).

External Video KV-Cache Retrieval. Here, we utilize an external CLIP-like model (Radford et al., 2021; Zhai et al., 2023) to retrieve question-relevant video KV-Cache, primarily as a baseline to assess whether retrieval can enhance VideoQA performance, as demonstrated in Section 3. Specifically, a CLIP-like model transforms each video frame into a vector $\mathbf{v}_t = f_v(v_t) \in \mathbb{R}^D$, where f_v represents the visual encoder, D denotes the vector dimension. Similarly, the question is encoded as $\mathbf{q}_i = f_t(q_i) \in \mathbb{R}^D$, where f_t is the text encoder. We then compute the cosine similarity between the embeddings of frame and question:

$$\text{Sim}(\mathbf{v}_t, \mathbf{q}_i) = \frac{\mathbf{v}_t \cdot \mathbf{q}_i}{\tau \|\mathbf{v}_t\| \|\mathbf{q}_i\|} \quad (2)$$

where τ is a learnable temperature parameter. This similarity is calculated at the frame level, rather than at the token level. Alternatively, we can group b consecutive frames into blocks by averaging their frame vectors and then compute block-level similarity scores. Finally, the r most relevant video

frames or $\lceil r/b \rceil$ video blocks are retrieved. The corresponding video KV-Cache, denoted as \mathbf{R} , is subsequently loaded onto the GPU for question-answering.

Internal Video KV-Cache Retrieval. Building on recent advancements in handling long sequences with LLMs (Xiao et al., 2024a; Li et al., 2024b; Fountas et al., 2024), we further explore using self-attention layers within Video-LLMs for retrieval. Similar to external retrieval, internal retrieval is still performed at the level of video frames or blocks.

During video modeling, the average of the key vectors of a frame is computed as its representative frame vector: $\mathbf{v}_t = \frac{1}{N_f} \sum_{j=1}^{N_f} \mathbf{k}_j \in \mathbb{R}^{D'}$, where N_f is the number of tokens per frame, and \mathbf{k}_j is the j -th key vector. To reduce computational overhead, we do not differentiate between attention heads and instead concatenate them into a single vector, with D' as the resultant dimension. Similarly, the question vector is computed as $\mathbf{q}_i = \frac{1}{N_q} \sum_{k=1}^{N_q} \mathbf{q}_{i,k} \in \mathbb{R}^{D'}$, where N_q is the number of tokens in the question, and $\mathbf{q}_{i,k}$ is its k -th query vector. The similarity computation and video KV-Cache retrieval are identical to that of external retrieval, except that τ is set to 1.

Note that, internal retrieval offers several advantages over external retrieval. First, it operates independently within each self-attention layer, allowing different layers to retrieve different video blocks.² This allows for a broader capture of video context. Additionally, internal retrieval reuses already computed hidden representations and does not introduce extra parameters, which reduces the computational overhead compared to external retrieval.

Question-answering using Retrieved KV. The retrieved Video KV-Caches serve as the context for video question-answering. Formally, the attention calculation is formulated as:

$$\mathbf{O} = \text{Attn}(\mathbf{W}_Q \mathbf{X}, [\mathbf{R}_k, \mathbf{W}_K \mathbf{X}], [\mathbf{R}_v, \mathbf{W}_V \mathbf{X}]), \quad (3)$$

where \mathbf{X} represents either the question tokens or the current token being decoded, and \mathbf{R}_k and \mathbf{R}_v are the key and value vectors from the context, which includes the retrieved video, question, and previously generated tokens.

3 EXPERIMENTS

3.1 BENCHMARK AND METRICS

MLVU_{dev-mc} (Zhou et al., 2024a) is the multiple-choice subset of the MLVU-dev benchmark. It focuses on evaluating the long-form video understanding of MLLMs. The question-answer pairs are manually labeled and can be divided into 3 groups: single-detail, multi-detail, and holistic. The evaluation metric is Accuracy.

QAEGO4D_{test-mc} (Di & Xie, 2024) is the multiple-choice subset of the QAEGO4D-test benchmark, focusing on question-answering in long egocentric videos. The evaluation metric is Accuracy.

EgoSchema (Mangalam et al., 2023) is a diagnostic benchmark for long VideoQA, featuring over 5000 multiple-choice questions and long temporal certificate length. It challenges AI models with long-term understanding, as current state-of-the-art models achieve significantly lower accuracy compared to human performance.

ActivityNet-QA (Yu et al., 2019) encompasses human-annotated QA pairs on 5,800 videos derived from the ActivityNet (Caba Heilbron et al., 2015) dataset. This benchmark is designed to assess the

Table 1: **Summary of the evaluation benchmarks.** MC stands for multiple-choice VideoQA, while OE refers to open-ended VideoQA.

Benchmark	Duration	#Videos	#QA	Type
MLVU _{dev-mc}	12 min	1,242	2,175	MC
QAEGO4D _{test-mc}	8.3 min	148	500	MC
EgoSchema	3 min	5,031	5,031	MC
ActivityNet-QA	2 min	800	8,000	OE
RVS-Ego	60 min	10	1,500	OE
RVS-Movie	30 min	22	2,000	OE

²For simplicity, we omit the layer index in the above explanation.

capabilities of VideoQA models in long-term spatiotemporal reasoning. Our evaluation methodology aligns with that of Video-ChatGPT (Maaz et al., 2024), employing GPT-3.5-turbo to judge the accuracy of the open-ended VideoQA responses.

RVS-Ego and **RVS-Movie** (Zhang et al., 2024a) are Streaming VideoQA benchmarks, constructed using 10 long videos from the Ego4D dataset (Grauman et al., 2022) and 22 long videos from the MovieNet dataset (Huang et al., 2020), respectively. These benchmarks feature open-ended questions paired with timestamps, which are initially generated by GPT-4V (OpenAI, 2023b) and GPT-4 (OpenAI, 2023a), and subsequently refined through manual filtering.

3.2 IMPLEMENTATION DETAILS

We evaluate our approach by integrating it into LLaVA-OV-0.5B and LLaVA-OV-7B (Li et al., 2024a), chosen for their simplicity and strong performance.

All experiments are conducted on NVIDIA H800 GPUs with BF16 precision. For video modeling, we process the video stream at 0.5 FPS, in line with GPT-4o’s testing on MLVU (Zhou et al., 2024a). The local window size is set to 15K. For external video KV-Cache retrieval, we use SigLIP-SO400M (Zhai et al., 2023) as the retriever. For internal KV-Cache retrieval, we set the block size (b) to 1 and the number of retrieved frames (r) to 64 by default, with further hyperparameter variations explored in Section 3.3.

Unless otherwise specified, **ReKV** refers to the use of internal video KV-Cache retrieval.

3.3 ABLATIONS

In this section, we conduct ablation studies on the effectiveness of in-context retrieval, number of retrieved frames, and the block size.

Effectiveness of In-context Retrieval. The experiments on QAEGO4D_{test-mc}, as presented in Table 2, demonstrate the effects of various retrieval methods on VideoQA accuracy and recall. A strong positive correlation is evident: higher recall generally leads to better VideoQA performance. Uniform Sampling, which sparsely selects frames from the video, achieves the lowest recall and, correspondingly, the poorest VideoQA accuracy. In contrast, Oracle Retrieval, with perfect recall, delivers the highest VideoQA accuracy, far surpassing Uniform Sampling. While External and Internal Retrieval do not match Oracle-level precision, both outperform Uniform Sampling, with Internal Retrieval excelling due to its higher recall. Notably, the larger model (LLaVA-OV-7B) achieves superior performance with Internal Retrieval, benefiting from its stronger reasoning capabilities.

Table 2: Ablation study on QAEGO4D_{test-mc}

Retrieval Method	VideoQA Acc.	Recall
LLaVA-OV-7B		
Uniform Sampling	50.8	6.1
External Retrieval	54.4	58.1
Internal Retrieval	56.4	70.5
Oracle Retrieval	60.2	100
LLaVA-OV-0.5B		
Uniform Sampling	44.0	6.1
External Retrieval	45.8	58.1
Internal Retrieval	46.0	63.4
Oracle Retrieval	47.4	100

The MLVU benchmark (Zhou et al., 2024a) encompasses three types of VideoQA tasks: *Single Detail* requires identifying a single critical plot within a long video, *Multi Detail* necessitates the integration of multiple plots, and *Holistic* demands a comprehensive understanding of the entire video. This makes MLVU an ideal platform for evaluating our in-context retrieval method. As shown in Table 3, both External and Internal Retrieval enhance the overall VideoQA accuracy over the Uniform Sampling baseline. The enhancements are most pronounced in Single Detail tasks, demonstrating that ReKV effectively retrieves question-relevant video context. Furthermore, Internal Retrieval significantly outperforms External Retrieval in Holistic tasks, likely due to its ability to capture broader context and leverage the Video-LLM’s video modeling capabilities, as discussed in Section 2.2.

Number of Retrieved Frames. We fix the block size ($b = 1$) and evaluate the impact of varying the numbers of retrieved frames ($r \in \{8, 16, 32, 48, 64, 80\}$) on the QAEGO4D and MLVU bench-

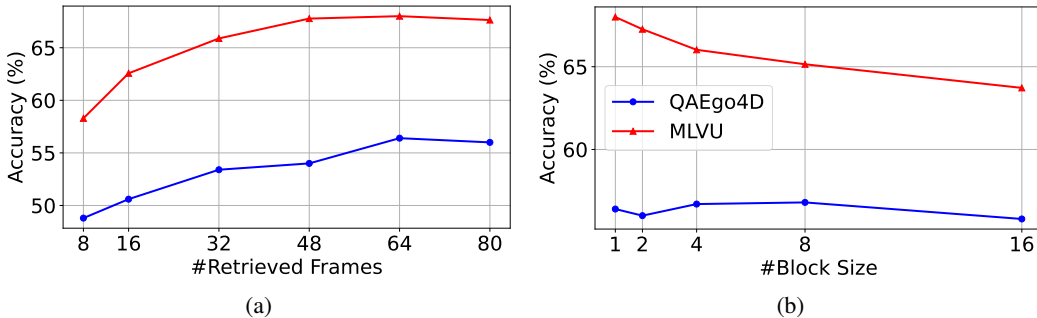


Figure 3: **Ablation study of retrieval hyperparameters:** (a) number of retrieved frames and (b) number of frames per retrieval block. Experiments are conducted with LLaVA-OV-7B.

Table 3: **Ablation study on MLVU_{dev-mc}.** The experiments are based on LLaVA-OV-7B.

Task	Single Detail			Multi Detail		Holistic		Avg.
	Needle	Ego	PlotQA	Order	Count	Topic	Anomaly	
Uniform Sampling	74.1	59.7	69.8	45.9	32.0	87.9	72.0	64.7
External Retrieval	78.6	69.6	71.6	40.2	37.9	84.5	63.0	66.3
Internal Retrieval	76.3	66.0	74.2	42.5	35.0	90.2	75.0	67.7

marks. As illustrated in Figure 3(a), increasing the number of retrieved frames generally improves VideoQA accuracy, as it implies capturing more relevant visual context. However, on MLVU, this improvement plateaus as more frames are retrieved since the additional irrelevant information hinders the subsequent question-answering process. Additionally, retrieving more frames increases the computational overhead of the question-answering stage, further slowing down inference.

Retrieval Block Size. When processing video streams, we group b consecutive frames into blocks for block-level retrieval. For this experiment, we fix the number of retrieved frames at $r = 64$ and evaluate different block sizes ($b \in 1, 2, 4, 8, 16$). With a fixed r , larger block sizes result in fewer, more concentrated retrieved blocks. Figure 3(b) shows that increasing block size negatively affects accuracy on MLVU, while performance on QAEgo4D remains relatively stable. This suggests that MLVU tasks benefit from retrieving more dispersed visual cues, aligning with its design of multi-detail and holistic tasks (Zhou et al., 2024a). In contrast, QAEgo4D primarily relies on a single relevant clip per question (Di & Xie, 2024).

3.4 OFFLINE VIDEO QUESTION-ANSWERING

Streaming video understanding is a relatively under-explored area, with limited StreamingVQA benchmarks available Zhang et al. (2024a). As discussed in Section 2.1, OfflineVQA can be considered as a special case of StreamingVQA. Thus, we first evaluate our method in the offline setting using four widely adopted long-form VideoQA benchmarks, comparing our results against state-of-the-art VideoQA methods. A summary of these benchmarks can be found in Table 1.

As shown in Table 4, our proposed ReKV always enhances the performance of LLaVA-OV-0.5B and LLaVA-OV-7B without additional training. Notably, LLaVA-OV-7B+ReKV outperforms two memory-based StreamingVQA models (VideoStreaming (Qian et al., 2024) and FlashVStream (Zhang et al., 2024a)) by a large margin. While the base model already demonstrates strong performance, and we **do not** claim credit for this achievement, our method can integrate seamlessly with Video-LLMs, benefiting from their ongoing advancements.

Table 4: **Offline video question-answering on four long-form benchmarks.** “Acc.” denotes accuracy, and “Score” is the open-ended answer rating by `gpt-3.5-turbo` on a scale from 1 to 5.

Method	MLVU	QAEGO4D	EgoSchema	ActivityNet-QA	
	dev Acc.	test Acc.	Acc.	Acc.	Score
GPT-4V (OpenAI, 2023b)	49.2	-	-	57.0	-
GPT-4o (OpenAI, 2024)	64.6	-	-	-	-
Gemini-1.5-Flash (Team et al., 2023)	-	-	65.7	55.3	-
Gemini-1.5-Pro (Team et al., 2023)	-	-	72.2	57.5	-
Video-ChatGPT-7B (Maaz et al., 2024)	31.3	-	-	-	-
LLaMA-VID-7B (Li et al., 2024d)	33.2	-	-	47.4	3.30
MiniGPT4-Video-7B (Ataallah et al., 2024)	44.5	-	-	44.3	3.35
Video-LLaVA-7B (Maaz et al., 2024)	47.3	-	-	-	-
LongVA-7B (Zhang et al., 2024b)	56.3	-	-	50.0	-
VideoStreaming (Qian et al., 2024)	-	-	44.1	-	-
Flash-VStream (Zhang et al., 2024a)	50.2	38.2	38.1	51.9	3.40
LLaVA-OV-0.5B (Li et al., 2024a)	50.3	44.0	26.8	50.5	3.02
+ReKV	53.2	46.0	29.8	52.1	3.15
LLaVA-OV-7B (Li et al., 2024a)	64.7	50.8	60.1	56.6	3.29
+ReKV	67.7	56.4	61.3	60.4	3.52

Table 5: **StreamingVQA benchmark results.** All methods are tested under identical conditions. “Video Enc.” is frames encoded per second. “Latency” is measured from question input to response completion. “GPU” indicates peak GPU memory usage, and “KV-Cache” refers to the video KV-Cache size offloaded per hour.

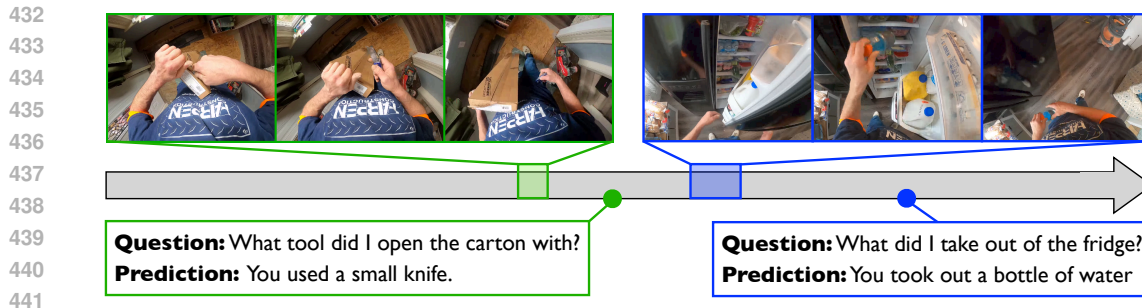
Retrieval Method	RVS-Ego		RVS-Movie		Running Speed		Memory Usage	
	Acc.	Score	Acc.	Score	Video Enc.	Latency	GPU	KV-Cache
Flash-VStream-7B Zhang et al. (2024a)	57.3	4.0	53.1	3.3	14 FPS	2.4s	20 GB	-
LLaVA-OV-7B								
Uniform Sampling	56.2	3.7	43.0	3.3	-	2.9s	21 GB	-
External Retrieval	62.4	3.9	53.6	3.5	11 FPS	5.8s	55 GB	18.8 GB/h
Internal Retrieval	63.7	4.0	54.4	3.6	11 FPS	3.3s	38 GB	18.8 GB/h
LLaVA-OV-0.5B								
Uniform Sampling	51.8	3.7	37.2	3.2	-	2.5s	7 GB	-
External Retrieval	54.1	3.8	44.7	3.4	17 FPS	4.1s	37 GB	4.0 GB/h
Internal Retrieval	54.7	3.9	44.6	3.4	17 FPS	1.6s	19 GB	4.0 GB/h

3.5 STREAMING VIDEO QUESTION-ANSWERING

We then evaluate our method on the streaming setting using the RVS-Ego and RVS-Movie benchmarks. During video stream modeling, questions are input immediately after their annotated end timestamps and answered based on the preceding video content.

Question-answering Performance. Table 5 presents the StreamingVQA performance. Both external and internal retrieval methods significantly outperform the uniform sampling baseline. Additionally, our approach enables LLaVA-OV-7B to surpass Flash-VStream Zhang et al. (2024a), demonstrating ReKV’s effectiveness for the StreamingVQA. Besides, internal retrieval consistently outperforms external retrieval, as discussed in Section 2.2.

Running Speed and Memory Usage. We also examine the running speed and memory usage under controlled conditions. Specifically, a 1-hour, 1080P video from RVS-Ego with 100 scattered questions is used. Each question is padded to 64 tokens, and the generated answers are fixed at 128



442 Figure 4: **Streaming VQA qualitative examples.** The example is drawn from the QAEGO4D benchmark. The
443 video stream is processed frame by frame. ● and ● mark the timestamps at which questions are posed. □ and
444 □ indicate the relevant video contexts that support answering these questions.

445
446 tokens in length. The video frames are pre-extracted at 0.5 FPS (1,800 frames in total) and streamed
447 to the Video-LLM frame by frame.

448
449 As illustrated in Table 5, both retrieval methods maintain high video encoding speeds, with
450 LLaVA-OV-7B achieving 11 FPS and LLaVA-OV-0.5B achieving 17 FPS. Moreover, KV-
451 Cache offloading remains manageable, with LLaVA-OV-7B at 18.8GB/h and LLaVA-OV-0.5B
452 at 4.0GB/h (see appendix for more details). External retrieval, however, introduces higher latency
453 and GPU memory usage due to additional computations in the external retriever, whereas internal
454 retrieval significantly reduces both. Figure 1 has also demonstrated that latency and GPU mem-
455 ory usage remain stable as more frames are processed. Flash-VStream also shows good efficiency.
456 However, it only maintains a relatively small memory footprint (681 tokens) (Zhang et al., 2024a),
457 leading to potential information loss when dealing with extremely long videos.

458 **Qualitative Examples.** Figure 4 presents an example of streaming video question-answering. Our
459 approach continuously processes video streams while responding to questions posed at different
460 timestamps. To improve efficiency, it stores and retrieves relevant video KV-Caches as contextual
461 information for answering these questions.

462 4 RELATED WORK

463
464
465 **LLMs for Video Understanding.** In recent years, there has been a surge of interest in leveraging
466 Large Language Models (LLMs) for video understanding, leading to the development of several
467 innovative approaches (Maaz et al., 2024; Zhang et al., 2024c; Li et al., 2024a). These models
468 typically use a Vision Encoder to extract video features, followed by a mapping step with Linear
469 Projection, MLP, or Q-Former (Li et al., 2023). The mapped features are combined with textual
470 data and fed into large language models (LLMs) to generate a text output. These models have
471 relatively simple architectures, requiring less training data and computational resources, yet they
472 achieve strong performance on short video understanding benchmarks (Xu et al., 2017; Xiao et al.,
473 2021; Li et al., 2024c). However, they employ sparse sampling or token compression techniques
474 to reduce the number of tokens, which can result in significant information loss when dealing with
475 longer or more content-rich videos. As a result, they are not well-suited for long video understanding
476 or streaming video understanding.

477 **Long Video Understanding.** A central challenge in long video understanding is effectively com-
478 pressing the information from lengthy videos. Many approaches use language as a bridge, condens-
479 ing videos into dense captions (Zhang et al., 2023; Islam et al., 2024; Zhou et al., 2024b). While
480 this achieves good results in some cases, compressing video content into text often leads to the loss
481 of crucial visual details.

482 Another line of research focuses on compressing long videos into a memory bank (Wu et al.,
483 2019; 2022; Wang et al., 2023). This approach has been integrated into Video-LLMs for **stream-**
484 **ing video understanding**, as shown in works like VideoStreaming (Qian et al., 2024) and Flash-
485 VStream (Zhang et al., 2024a). These methods dynamically update the memory during video pro-
cessing and utilize it for downstream tasks. Despite their innovation, a major limitation of these

486 methods is their failure to account for video length and information density, especially when using
487 a fixed memory size. For example, Flash-VStream compresses both 10-second clips and hour-long
488 videos into the same 681 tokens. Furthermore, these methods lack interpretability, making it diffi-
489 cult to determine how much information is being compressed into the memory or whether relevant
490 video information is being accurately retrieved during downstream tasks.

491 In pursuit of greater interpretability in long video understanding, methods such as GroundVQA (Di
492 & Xie, 2024) and GeLM (Chen et al., 2024) advocate for localizing relevant video clips while
493 responding to user queries. Drawing inspiration from these, this work refrains from excessively
494 condensing video information. By harnessing the causal capabilities of Video-LLMs, it preserves
495 the entire Video KV-Cache, allowing for the retrieval of relevant information when required. This
496 strategy effectively mitigates the substantial loss of video content while improving interpretability.

497 **Long Context Handling for LLMs.** Handling long text sequences in LLMs has been a major
498 challenge due to high computational and memory costs, leading to training constraints on shorter
499 sequences. Techniques like StreamingLLM (Xiao et al., 2024b) and LM-Infinite (Han et al., 2023)
500 use sliding window attention to process long sequences incrementally, but discard distant tokens,
501 limiting the model’s ability to capture long-range dependencies. Recent approaches (Xiao et al.,
502 2024a; Li et al., 2024b; Fountas et al., 2024) address this by storing and retrieving previously com-
503 puted KV-Caches, enabling better recall of distant contexts.

504 **Retrieval-Augmented Generation.** Retrieval-augmented generation (RAG) combines retrieval
505 mechanisms with generative models to enhance performance across various NLP tasks by incor-
506 porating external knowledge Guu et al. (2020); Lewis et al. (2020); Borgeaud et al. (2022) and
507 improving performance in vision-language tasks Xu et al. (2024). In-context retrieval, recently pro-
508 posed for handling long inputs Ram et al. (2023), retrieves information from the input document
509 itself rather than an external knowledge base. In-context KV-Cache retrieval further improves effi-
510 ciency by pre-encoding long documents, avoiding redundant encodings, and leveraging the LLM’s
511 own retrieval capabilities for faster, more effective performance.

512 5 CONCLUSION

513 In conclusion, this paper introduces ReKV, a training-free approach designed to enhance the ef-
514 ficiency of Video Large Language Models (Video-LLMs) for streaming video question-answering
515 (StreamingVQA). Unlike conventional video question-answering (VideoQA) systems that must pro-
516 cess entire videos before answering, ReKV enables rapid, real-time responses. By employing a
517 sliding-window attention mechanism, it ensures that the model only considers a subset of previous
518 frames while encoding the video stream, significantly cutting down on computational demands. To
519 retain key video context, we developed an in-context KV-Cache retrieval method that efficiently
520 stores and reloads key-value vectors that relevant for each query. This targeted retrieval strategy,
521 combined with the ability to perform video modeling and question-answering on separate processes
522 and GPUs, results in a highly efficient streaming VideoQA system. Extensive experiments show
523 that ReKV not only surpasses existing VideoQA models in performance but also enhances their
524 practicality for real-world streaming applications.

525 6 LIMITATIONS AND FUTURE WORK

526 While ReKV improves the accuracy and efficiency of Video-LLMs in the StreamingVQA task, it still
527 has several limitations that deserves future investigation: *First*, although the KV-Cache offloading
528 to RAM or disk is manageable, as shown in Table 5, handling extremely long video streams, such
529 as those in surveillance, may lead to an unsustainable increase in cache size. This issue can be
530 mitigated by integrating techniques such as quantization, token pruning, and compression. *Second*,
531 the use of a constant block size for grouping consecutive frames during retrieval can disrupt video
532 continuity. A more refined solution would involve segmenting videos into semantically coherent
533 blocks. *Third*, our method retrieves a fixed number of frames. Future work could explore dynamic
534 retrieval strategies that adjust the number of frames based on video context and query requirements.
535 *Finally*, StreamingVQA remains an under-explored task with few available benchmarks. Developing
536 high-quality benchmarks with precise temporal annotations is crucial for advancing future research.
537
538
539

REFERENCES

- 540
541
542 Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and
543 Mohamed Elhoseiny. Minigt4-video: Advancing multimodal llms for video understanding with
544 interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024.
- 545 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Milli-
546 can, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al.
547 Improving language models by retrieving from trillions of tokens. In *ICML*, 2022.
- 548 Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A
549 large-scale video benchmark for human activity understanding. In *CVPR*, 2015.
- 550 Qirui Chen, Shangzhe Di, and Weidi Xie. Grounded multi-hop videoqa in long-form egocentric
551 videos. *arXiv preprint arXiv:2408.14469*, 2024.
- 552
553 Shangzhe Di and Weidi Xie. Grounded question-answering in long egocentric videos. In *CVPR*,
554 2024.
- 555 Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong
556 Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale.
557 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
558 19358–19369, 2023.
- 559 Zafeirios Fountas, Martin A Benfeghoul, Adnan Omerjee, Fenia Christopoulou, Gerasimos Lam-
560 pouras, Haitham Bou-Ammar, and Jun Wang. Human-like episodic memory for infinite context
561 llms. *arXiv preprint arXiv:2407.09450*, 2024.
- 562 Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne West-
563 phal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al.
564 The” something something” video database for learning and evaluating visual common sense. In
565 *ICCV*, 2017.
- 566
567 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-
568 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in
569 3,000 hours of egocentric video. In *CVPR*, 2022.
- 570
571 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented
572 language model pre-training. In *ICML*, 2020.
- 573
574 Chi Han, Qifan Wang, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. Lm-infinite: Simple
575 on-the-fly length generalization for large language models. *arXiv preprint arXiv:2308.16137*,
576 2023.
- 577
578 Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava,
579 and Ser-Nam Lim. Ma-lmm: Memory-augmented large multimodal model for long-term video
580 understanding. In *CVPR*, 2024.
- 581
582 Lianghua Huang, Xin Zhao, and Kaiqi Huang. Got-10k: A large high-diversity benchmark for
583 generic object tracking in the wild. *TPAMI*, 2019.
- 584
585 Qingqiu Huang, Yu Xiong, Anyi Rao, Jiase Wang, and Dahua Lin. Movienet: A holistic dataset for
586 movie understanding. In *ECCV*, 2020.
- 587
588 Md Mohaiminul Islam, Ngan Ho, Xitong Yang, Tushar Nagarajan, Lorenzo Torresani, and Gedas
589 Bertasius. Video recap: Recursive captioning of hour-long videos. In *CVPR*, 2024.
- 590
591 Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-
592 temporal reasoning in visual question answering. In *CVPR*, 2017.
- 593
Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action
video dataset. *arXiv:1705.06950*, 2017.

- 594 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
595 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented gener-
596 ation for knowledge-intensive nlp tasks. *NeurIPS*, 2020.
- 597 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei
598 Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint*
599 *arXiv:2408.03326*, 2024a.
- 600
601 Jingyao Li, Han Shi, Xin Jiang, Zhenguo Li, Hong Xu, and Jiaya Jia. Quickllama: Query-aware
602 inference acceleration for large language models. *arXiv preprint arXiv:2406.07528*, 2024b.
- 603
604 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
605 pre-training with frozen image encoders and large language models. In *ICML*, 2023.
- 606 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,
607 Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In
608 *CVPR*, 2024c.
- 609 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language
610 models. In *ECCV*, 2024d.
- 611
612 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle
613 Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before
614 generation. *arXiv preprint arXiv:2404.14469*, 2024e.
- 615 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:
616 Towards detailed video understanding via large vision and language models. In *ACL*, 2024.
- 617
618 Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic bench-
619 mark for very long-form video language understanding. *NeurIPS*, 2023.
- 620 Matthias Muller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet:
621 A large-scale dataset and benchmark for object tracking in the wild. In *ECCV*, 2018.
- 622
623 OpenAI. Gpt-4, March 2023a. URL [https://cdn.openai.com/papers/
624 gpt-4-system-card.pdf](https://cdn.openai.com/papers/gpt-4-system-card.pdf).
- 625 OpenAI. Gpt-4v, September 2023b. URL [https://openai.com/index/
626 gpt-4v-system-card/](https://openai.com/index/gpt-4v-system-card/).
- 627
628 OpenAI. Gpt-4o, May 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- 629 Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi
630 Wang. Streaming long video understanding with large language models. *arXiv preprint*
631 *arXiv:2405.16009*, 2024.
- 632
633 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
634 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
635 models from natural language supervision. In *ICML*, 2021.
- 636 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and
637 Yoav Shoham. In-context retrieval-augmented language models. *ACL*, 2023.
- 638
639 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
640 hanced transformer with rotary position embedding. *Neurocomputing*, 2024.
- 641 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
642 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
643 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 644
645 Jiahao Wang, Guo Chen, Yifei Huang, Limin Wang, and Tong Lu. Memory-and-anticipation trans-
646 former for online action understanding. In *ICCV*, 2023.
- 647 Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krahenbuhl, and Ross
Girshick. Long-term feature banks for detailed video understanding. In *CVPR*, 2019.

- 648 Chao-Yuan Wu, Yanghao Li, Karttikeya Mangalam, Haoqi Fan, Bo Xiong, Jitendra Malik, and
649 Christoph Feichtenhofer. Memvit: Memory-augmented multiscale vision transformer for efficient
650 long-term video recognition. In *CVPR*, 2022.
- 651 Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan
652 Liu, and Maosong Sun. Infilm: Training-free long-context extrapolation for llms with an efficient
653 context memory. In *ICML Workshop*, 2024a.
- 654 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
655 language models with attention sinks. In *ICLR*, 2024b.
- 656 Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa:next phase of question-
657 answering to explaining temporal actions. In *CVPR*, 2021.
- 658 Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang.
659 Video question answering via gradually refined attention over appearance and motion. In *ACM
660 Multimedia*, 2017.
- 661 Jilan Xu, Yifei Huang, Junlin Hou, Guo Chen, Yuejie Zhang, Rui Feng, and Weidi Xie. Retrieval-
662 augmented egocentric video captioning. In *CVPR*, 2024.
- 663 Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Zero-shot video
664 question answering via frozen bidirectional language models. *NeurIPS*, 2022.
- 665 Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-
666 qa: A dataset for understanding complex web videos via question answering. In *AAAI*, 2019.
- 667 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
668 image pre-training. In *ICCV*, 2023.
- 669 Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas
670 Bertasius. A simple llm framework for long-range video question-answering. *arXiv preprint
671 arXiv:2312.17235*, 2023.
- 672 Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin.
673 Flash-vstream: Memory-based real-time understanding for long video streams. *arXiv preprint
674 arXiv:2406.08085*, 2024a.
- 675 Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue
676 Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision.
677 *arXiv preprint arXiv:2406.16852*, 2024b.
- 678 Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and
679 Chunyuan Li. Llava-next: A strong zero-shot video understanding model, April 2024c. URL
680 <https://llava-vl.github.io/blog/2024-04-30-llava-next-video/>.
- 681 Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang,
682 Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video
683 understanding. *arXiv preprint arXiv:2406.04264*, 2024a.
- 684 Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Na-
685 grani, and Cordelia Schmid. Streaming dense video captioning. In *CVPR*, 2024b.
- 686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

In the appendix, we provide additional implementation details of our method, including multi-processing serving, the prompt templates for VideoQA, positional encoding, and KV-Cache size calculation, to support reproducibility.

A ADDITIONAL IMPLEMENTATION DETAILS

A.1 MULTI-PROCESSING SERVING

As discussed in Section 2.1, our approach enables the separation of video modeling and question-answering across different processes and GPUs, significantly enhancing efficiency in real-world applications. Specifically, we dedicate a primary process for video stream encoding, utilizing sliding-window attention to analyze the video and store the computed cache in RAM. If RAM capacity is exceeded, the data can be offloaded to disk. Additionally, a process pool is maintained, with the number of processes determined by the frequency of queries and available resources. Each process loads the same Video-LLM parameters but operates independently. The video processing continues uninterrupted, without waiting for question-answering tasks to complete. When a query is posed, we log its timestamp to ensure that video information after this point is excluded from the answer. An available process from the pool is then activated to retrieve relevant video key-value vectors using our method, loading them onto its GPU for question-answering. This approach enables efficient StreamingVQA applications, with significant potential in areas such as robotics, surveillance, augmented reality, and live broadcasting.

A.2 PROMPT TEMPLATES FOR VIDEOQA

We use the consistent prompt template on all multiple-choice VideoQA benchmarks. Text in red indicates variable inputs.

```
System:
You are a helpful assistant.
User:
<video>
Question: <question>
Options:
(A) <Option_A>
(B) <Option_B>
(C) <Option_C>
(D) <Option_D>
(E) <Option_E>
Answer with the option's letter from the given choices directly.
Assistant:
```

The prompt template for open-ended VideoQA is rather simpler:

```
System:
You are a helpful assistant.
User:
<video>
<question>
Assistant:
```

A.3 POSITIONAL ENCODING

The baseline Video-LLMs we used incorporate Rotary Position Embeddings (RoPE) [Su et al. \(2024\)](#), a widely used relative positional encoding method. Our video streaming encoding process follows LM-Infinite [Han et al. \(2023\)](#), where RoPE operates normally within the local window but is constrained by a “distance ceiling” for distant tokens. For question-answering, we do not account for the original positions of the retrieved KV-Caches, instead treating them as regular context tokens. We also experimented with the static variation from Inf-LLM [Xiao et al. \(2024a\)](#), which assigns the same position to all retrieved tokens. Results indicate that applying standard RoPE to retrieved video tokens yields better performance, likely due to the critical role of capturing temporal information in video understanding.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A.4 KV-CACHE SIZE CALCULATION

The size of the KV-Cache can be calculated using the following formula, assuming BF16 precision:

$$2 \times L \text{ layers} \times T \text{ frames} \times M \text{ tokens/frame} \times H \text{ heads} \times D \text{ dimension} \times 2 \text{ bytes.}$$

For LLaVA-OV-7B [Li et al. \(2024a\)](#), with $L = 28$, $M = 196$, $H = 4$, and $D = 128$, processing a 1-hour video at 0.5 FPS ($T = 1800$) results in a total KV-Cache size of 18.8 GB.

Similarly, for LLaVA-OV-0.5B [Li et al. \(2024a\)](#), with $L = 24$, $M = 196$, $H = 2$, and $D = 64$, processing a 1-hour video at 0.5 FPS results in a total KV-Cache size of 4.0 GB.

These theoretical calculations are consistent with the experimental results shown in [Table 5](#).