

StreamMem: Query-Agnostic KV Cache Memory for Streaming Video Understanding

Anonymous CVPR submission

Paper ID ****

Abstract

001 Multimodal large language models (MLLMs) have made
 002 significant progress in visual-language reasoning, but their
 003 ability to efficiently handle long videos remains limited. De-
 004 spite recent advances in long-context MLLMs, storing and
 005 attending to the key-value (KV) cache for long visual con-
 006 texts incurs substantial memory and computational over-
 007 head. Existing visual compression methods require either
 008 encoding the entire visual context before compression or
 009 having access to the questions in advance, which is im-
 010 practical for long video understanding and multi-turn con-
 011 versational settings. In this work, we propose **Stream-**
 012 **Mem**, a query-agnostic KV cache memory mechanism for
 013 streaming video understanding. Specifically, StreamMem
 014 encodes new video frames in a streaming manner, com-
 015 pressing the KV cache using attention scores between vi-
 016 sual tokens and generic query tokens, while maintaining
 017 a fixed-size KV memory to enable efficient question an-
 018 swering (QA) in memory-constrained, long-video scenar-
 019 ios. Evaluation on three long video understanding and three
 020 streaming video question answering benchmarks shows that
 021 StreamMem achieves state-of-the-art performance in query-
 022 agnostic KV cache compression and is competitive with
 023 query-aware compression approaches.

024 1. Introduction

025 Recent advances in Multimodal Large Language Models
 026 (MLLMs) [2, 5, 17, 49, 55] enable the capability to reason
 027 across textual and visual contents. Despite fast improve-
 028 ments, their capabilities to capture fine-grained details of
 029 actions, motions, object locations, interactions between ob-
 030 jects, and spatial-temporal orders of events in long videos
 031 are still limited [58]. There are two main reasons for this.
 032 Firstly, encoding the frames in a long video often gener-
 033 ates a large number of visual tokens, exceeding the context
 034 length of the underlying Large Language Model (LLM).
 035 Secondly, storing the KV cache of these large number of

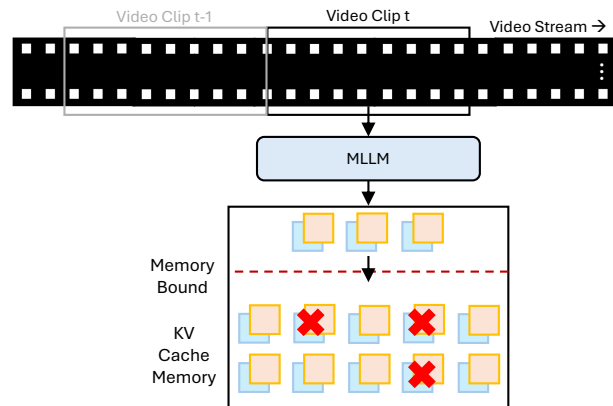


Figure 1. **Query-agnostic key-value (KV) cache compression in streaming video.** StreamMem addresses the challenge of streaming video processing under a memory budget by introducing a query-agnostic KV compression strategy.

visual tokens and attending to them during decoding poses
 significant memory and computational overhead. While the
 first issue has been alleviated by recent progress in long-
 context LLMs [38, 46, 51], the memory and compute ef-
 ficiency of dealing with long videos remains a challenge,
 especially for real-world applications on edge devices.

A number of recent works explored video token com-
 pression strategies to tackle long video understanding, in-
 cluding temporal compression [39, 40], spatial compres-
 sion [3, 53], and hybrid methods [14, 35, 41, 56]. These
 approaches can suffer significant information loss. For ex-
 ample, the action information in the video often cannot be
 captured with any single frame in the video. Many such
 methods also rely on having access to the text query for
 visual compression [16, 23, 24], which is often unknown at
 the time of video processing in real-world applications [20].

In parallel to these efforts, recent works start to explore
 streaming video processing with MLLMs, a paradigm in
 which video frames are incrementally encoded as they ar-
 rive, without prior knowledge of the video’s full length or
 the downstream query. Compared to offline video process-

057 ing, the streaming video processing setup is much more
058 flexible, as the model does not need to know the text query
059 or the length of the video when encoding visual informa-
060 tion. ReKV [8] is a leading work in this direction. It en-
061 codes new video frames in the stream with sliding window
062 attention and stores the KV cache. When the model receives
063 a question, it retrieves the most relevant KV cache in each
064 layer with in-context retrieval. While this method is shown
065 to be effective, it consumes significant memory to store all
066 the KV cache. Offloading the KV cache to memory or disk
067 and reloading them upon retrieval could also be very ineffi-
068 cient as the video becomes longer. LiveVLM [30] proposes
069 a KV compression mechanism to reduce the KV cache size
070 by 70%. While LiveVLM alleviates the issue of memory
071 consumption, it simply throws out the KV cache of earlier
072 tokens when the memory upper bound is reached, which
073 can lead to complete forgetting of earlier parts of the video.

074 To enable efficient long video processing in memory-
075 constrained environments, we introduce StreamMem, a
076 training-free and query-agnostic KV cache memory system
077 for streaming video understanding with MLLMs. Stream-
078 Mem maintains a bounded memory footprint by contin-
079 uously compressing the KV cache after each incoming
080 video clip, thus preventing out-of-memory (OOM) errors
081 and avoiding costly memory offloading regardless of video
082 length. To achieve effective and efficient memory reten-
083 tion, StreamMem leverages a novel saliency metric based
084 on cross-attention scores between visual tokens and *chat*
085 *template* tokens, allowing it to select and preserve infor-
086 mative visual content in a query-agnostic manner. In ad-
087 dition, it incorporates an input frame compression module
088 to reduce frame-level redundancy prior to MLLM encod-
089 ing, and a frame-wise KV merging mechanism that con-
090 structs prototype representations for each observed frame.
091 Together, these components produce a diverse yet com-
092 pact KV cache that supports accurate and memory-efficient
093 streaming question answering.

094 We evaluate StreamMem across three offline and
095 three streaming long video understanding benchmarks
096 (EgoSchema [29], MLVU [60], VideoMME [10]; RVS-
097 Ego and RVS-Movie [50]; OVO-Bench [31]) using three
098 open-source pre-trained MLLMs (LLaVA-OneVision [21],
099 Qwen2-VL [42], and Qwen2.5-VL [2]). Results show that
100 StreamMem consistently retains high utility while keeping
101 the KV cache compact across videos of varying lengths and
102 question types. It not only surpasses state-of-the-art stream-
103 ing video models, but also achieves competitive perfor-
104 mance with methods that rely on significantly larger mem-
105 ory budgets. Comprehensive ablation studies confirm the
106 contribution of each component in the StreamMem frame-
107 work. By enabling continuous, scalable memory compres-
108 sion without fine-tuning, StreamMem provides a crucial
109 step toward building real-time MLLM agents capable of

continuous video understanding in open-world settings. 110

2. Related Work 111

Streaming video understanding with MLLMs. Stream- 112
ing video understanding refers to the setting where the 113
model continuously processes video frames in real-time. 114
The model does not know the length of the video before- 115
hand and therefore cannot sample a fixed number of frames 116
uniformly from the video. VideoLLM-online [3] presents 117
an MLLM that supports efficient streaming video process- 118
ing and real-time dialogues. However, it aggressively down- 119
samples each video frame to only include 10 visual to- 120
kens, limiting its understanding of fine-grained details in the 121
video. Flash-VStream [50] and Dispider [34] use external 122
memory modules to compress and organize visual tokens. 123
Upon receiving a question, the model retrieves relevant vi- 124
sual tokens, combines them with the text tokens, and feeds 125
them through the MLLM. Recent works start to explore 126
KV cache compression and retrieval for video understand- 127
ing. ReKV [8] encodes the video in streaming fashion and 128
stores all the KV cache by offloading to memory or disk, 129
and performs in-context retrieval of the relevant KV cache 130
for each layer when answering a question. The offloading of 131
KV cache could incur a lot of memory and is not scalable 132
to ultra-long videos. LiveVLM [30] designs a KV cache 133
compression strategy for MLLMs to significantly reduce 134
memory usage and improve question answering speed com- 135
pared to ReKV. However, it uses a fixed compression ratio 136
throughout the video and relies on first-in-first-out (FIFO) 137
strategy to maintain a constrained memory, which leads to 138
forgetting of earlier information in long videos, even though 139
they might be informative. StreamMem resolves this is- 140
sue by compressing the KV cache memory and the KVs 141
from the new frames together and ensures a fixed-size KV 142
cache memory throughout the video stream. Concurrent 143
work InfiniPot-V [20] also studies streaming video process- 144
ing with constrained memory consumption. Different from 145
StreamMem, they used a combination of two compression 146
mechanisms, temporal-axis redundancy reduction and value 147
norm-based selection. 148

Long video understanding with MLLMs. Long video 149
understanding has been a great challenge for MLLMs given 150
their constrained context length. Early models such as 151
LLaVA [25, 26] can only process a very small number of 152
frames, leading to significant information loss. A number of 153
training-based methods [27, 35, 36, 51] have been proposed 154
to reduce the number of visual tokens needed to represent 155
each video frame. In addition, recent foundation models 156
like Gemini [5] and Qwen2.5-VL [2] also have inherent 157
long visual context processing capability. These training- 158
based methods, however, are often very computationally 159
expensive, especially when fine-tuning large MLLMs, and 160

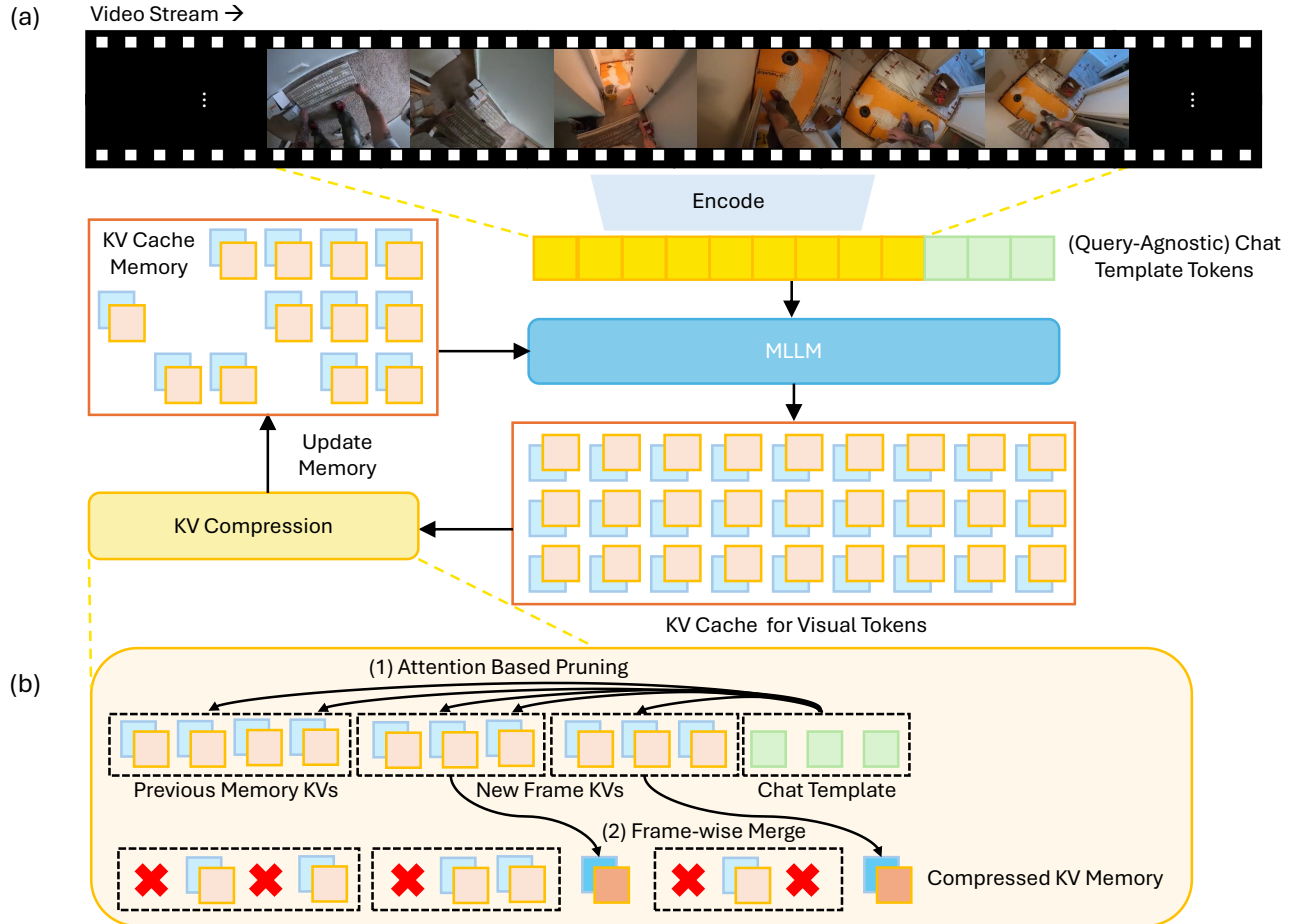


Figure 2. (a) The overall workflow of StreamMem for streaming video understanding. Incoming frames are first filtered to reduce redundancy, then passed through the vision encoder and integrated with the existing KV memory via cross-attention. The resulting KV cache is compressed to maintain a fixed memory budget, enabling continual processing of future frames or downstream question answering. (b) Detailed illustration of the KV compression module. Some KV cache in the memory and the new frames are pruned according to the attention score between the keys and the proxy queries. In addition, we aggregated the key-value pairs for each new frame into a single frame-level prototype via weighted merging (shown in darker squares). This combination of pruning and merging ensures compact yet expressive memory representations for long video sequences.

161 needs re-training for new foundation models. Training-free
 162 long video understanding methods [28, 43, 44, 56]
 163 compress the input visual tokens or the KV cache without the
 164 need to fine-tune the model, providing more flexibility for
 165 plug-and-play usage in new and more powerful MLLMs.
 166 StreamMem draws inspirations from the training-free meth-
 167 ods for KV cache compression of video tokens, but focuses
 168 on the streaming setting where neither the length of the
 169 video nor the query is known during memory-constrained
 170 video encoding.

171 **KV cache compression in LLMs.** KV cache compres-
 172 sion methods aim to greatly improve both memory and
 173 time efficiency of LLMs when operated in long input con-
 174 texts. A number of methods explored leveraging the cross-

175 attention weights between the query and the context to
 176 identify the most important entries in the KV cache for
 177 LLMs [11, 22, 47, 59]. This strategy is also adopted in
 178 MLLMs for efficient visual understanding [4, 54]. How-
 179 ever, the query might not be available when the model
 180 processes the long context in many real-world scenar-
 181 ios, limiting the applicability of the query-dependent ap-
 182 proach. To eliminate this dependency, some recent works
 183 explored query-agnostic KV cache compression mecha-
 184 nisms [7, 12, 15, 18, 32]. Similar to this work, Zhang
 185 et al. [52] and Arif et al. [1] explored using the attention
 186 weights of the [CLS] token for KV cache compression in
 187 MLLMs. In between query-dependent and query-agnostic
 188 methods, there are also methods which use task instructions
 189 or task-specific proxy prompts [6, 19]. StreamMem belongs

190 to the most flexible category of query-agnostic methods and
 191 does not need full-context encoding, making it suitable for
 192 streaming encoding of long videos.

193 3. Preliminaries

194 **Offline video understanding with MLLMs.** The stan-
 195 dard approach to offline video understanding with MLLMs
 196 proceeds as follows. Given a long video, a fixed num-
 197 ber of frames f_1, \dots, f_T are uniformly sampled from the
 198 video, where T is determined based on the model’s con-
 199 text length or computational and memory constraints. The
 200 frames are then passed through the model’s vision encoder
 201 (typically comprising a Vision Transformer (ViT) back-
 202 bone [9, 48, 57] and a projection layer) to get N visual
 203 tokens. The visual tokens are concatenated with the text
 204 tokens, including system prompts (preceding the visual to-
 205 kens) and user queries (following the visual tokens), and
 206 the entire sequence is fed into the LLM. The LLM then gen-
 207 erates a response via autoregressive decoding. To acceler-
 208 ate decoding, key-value (KV) caches are constructed during
 209 this process.

210 KV cache compression for MLLMs in streaming video.

211 In streaming video processing with MLLMs, the video
 212 length is typically unknown in advance, precluding uniform
 213 frame sampling strategies used in the offline settings. At
 214 each time step t , the model receives a new video clip v_t
 215 (a fixed-length frame segment), encodes it into a sequence
 216 of visual tokens, and forwards them through the LLM. The
 217 model then generates the corresponding key and value mat-
 218 rices K_t^i and V_t^i at each transformer layer i , by attending
 219 to all accumulated visual tokens from prior clips.

220 However, naively storing all keys and values over time
 221 leads to linear growth in memory, which is infeasible for
 222 long videos. This motivates the need for KV cache com-
 223 pression mechanisms that maintain a fixed memory foot-
 224 print. We denote the compressed key and value matrices at
 225 time step t and layer i as $K_t^{i'}$ and $V_t^{i'}$, respectively. The
 226 objective is to compute compressed representations:

$$227 \quad K_t^{i'}, V_t^{i'} = \text{Compress}(K_{t-1}^{i'}, K_t^i, V_{t-1}^{i'}, V_t^i),$$

228 subject to the global memory constraint:

$$229 \quad \sum_{i=1}^L \|K_t^{i'}\|_0 \leq M,$$

230 where L is the number of transformer layers in the MLLM
 231 and M is the total memory budget for all layers combined.

232 The design of effective compression strategies that retain
 233 essential temporal information while bounding memory use-
 234 age is a key challenge in streaming long video processing
 235 with MLLMs. Existing approaches such as ReKV [8] and

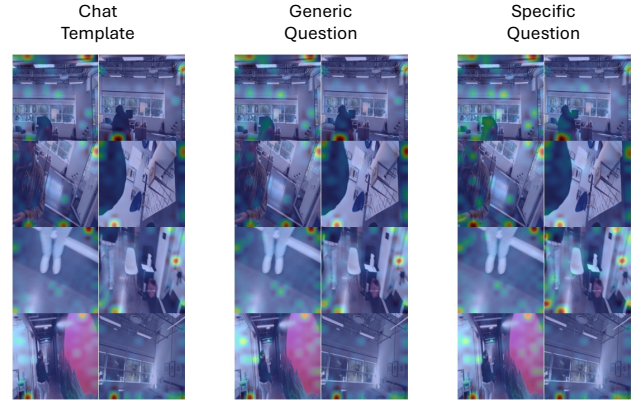


Figure 3. **Visualization of visual tokens attended to by different text queries.** Red indicates higher attention scores. Despite minor variations, different text queries attend to largely overlapping regions of the input images. The “Generic Question” is “What is happening in the video?”, while the “Specific Question” is “What occurs just before reading the magazines?” Attention scores are averaged across all layers and heads, and then interpolated from 14×14 to 384×384 to match the image resolution. The MLLM used is LLaVA-OneVision, and the video clip is sourced from the RVS-Ego benchmark (which uses videos from the Ego4D dataset [13]).

LiveVLM [30] do not address this constraint effectively, as
 their KV cache grows linearly over time, resulting in un-
 bounded memory consumption for long videos.

239 4. Method

240 We now describe the key components of StreamMem, our
 241 proposed framework for efficient streaming video under-
 242 standing with MLLMs. At each time step t , a new seg-
 243 ment of frames is received from the video stream. These
 244 frames first undergo an input filtering step to remove tem-
 245 poral redundancy. The filtered frames are then encoded by
 246 the vision encoder and processed by the MLLM to produce
 247 key-value (KV) representations $\{K_t^i, V_t^i\}_{i=1}^L$ at each trans-
 248 former layer i .

249 To prevent unbounded memory growth over time, the
 250 newly computed KVs are merged with the compressed KV
 251 memory from the previous time step, $\{K_{t-1}^{i'}, V_{t-1}^{i'}\}_{i=1}^L$, and
 252 passed through a compression module. This module ap-
 253 plies two complementary strategies: (1) a novel attention-
 254 based pruning method that leverages cross-attention scores
 255 between proxy query tokens and visual tokens, and (2) a
 256 frame-wise KV merging mechanism that condenses spatial
 257 information into compact prototype representations. The
 258 output of the compression module forms the updated mem-
 259 ory $\{K_t^{i'}, V_t^{i'}\}_{i=1}^L$, which is used by the MLLM at the next
 260 time step. An overview of the full pipeline is illustrated in
 261 Figure 2, and the KV compression procedure is detailed in

262 Algorithm 1.

263 4.1. Input Frame Filtering

264 Before processing by the MLLM, each incoming video
265 clip (a chunk of consecutive frames) is passed through a
266 lightweight filtering step to reduce temporal redundancy.
267 Given a sequence of frames, we compute their visual em-
268 beddings using the vision encoder. For each consecutive
269 pair of frames, we measure the cosine similarity between
270 their embeddings. If the similarity exceeds a predefined
271 threshold δ , the two frames are deemed redundant and their
272 representations are merged by simple averaging.

273 This lightweight filtering step is similar to the temporal
274 compression used in LongVU [35]. In contrast to previ-
275 ous streaming approaches that rely on sliding window at-
276 tention [8, 30], our method explicitly reduces redundancy
277 in the input space. This ensures that highly similar frames
278 (common in static scenes or high-frame-rate videos) do not
279 overwhelm the KV cache with repetitive information, ulti-
280 mately preserving the diversity and informativeness of the
281 stored memory.

282 4.2. KV Cache Memory

283 After frame-wise token compression, the retained visual to-
284 kens of the current video segment v_t are concatenated with
285 a set of auxiliary query tokens and passed into the MLLM.
286 The model computes key-value pairs $\{K_t^i, V_t^i\}_{i=1}^L$ in each
287 transformer layer i , attending over both the current tokens,
288 the tokens in the previous time step, and the compressed KV
289 cache from previous time step, $\{K_{t-1}^{i'}, V_{t-1}^{i'}\}_{i=1}^L$.

290 To guide the KV cache compression process, we rely on
291 the cross-attention scores between the auxiliary query to-
292 kens and the visual tokens. This attention-based saliency
293 measure has proven effective in prior works [4, 44] for real
294 user queries. However, unlike those settings, our method
295 operates under a query-agnostic streaming setup, where the
296 user query is unavailable at the time of visual token selec-
297 tion.

298 To approximate a generic query, we leverage the sys-
299 tem’s chat template tokens as a proxy. Specifically, we use
300 the tokens: `<|im_end|><|im_start|>assistant`
301 `\n`, which we append after the visual tokens. Due to the
302 prevalence of video captioning data during the MLLM pre-
303 training, this implicitly prompts the MLLM to generate a
304 generic video description even in the absence of an explicit
305 question. As a result, we expect the model to implicitly at-
306 tend to informative visual content in this setup.

307 Formally, let $Q \in \mathbb{R}^{q \times d}$ be the query representation of
308 the chat template tokens at a given layer, and K_t be the key
309 matrices of the visual tokens from the KV memory and the
310 current clip. The cross-attention scores are computed as:

$$311 A_t^i = \text{Softmax} \left(\frac{Q(K_t)^{\top}}{\sqrt{d}} \right), \quad (1)$$

where A_t^i denotes the attention weights from chat template
tokens to visual tokens. We aggregate these scores (e.g.,
by averaging over q) to obtain an importance score for each
visual token, which we use to select the top- k most salient
visual tokens to retain in the compressed cache from each
layer. The memory budget is even distributed across all lay-
ers.

In addition to pruning, StreamMem further compresses
memory via KV merging. Inspired by frame-level merging
in MLLMs [14] and visual token merging [56], we compute
a *prototype* key and value representation for each frame.
This is done by computing a weighted average of the keys
and values based on the normalized attention scores:

$$\bar{K}_t^i = \sum_{j=1}^n \alpha_j^i \cdot K_{t,j}^i, \quad \bar{V}_t^i = \sum_{j=1}^n \alpha_j^i \cdot V_{t,j}^i, \quad (2)$$

where α_j^i denotes the normalized importance score of the
 j -th visual token at layer i .

These prototype representations \bar{K}_t^i, \bar{V}_t^i are inserted at
the end of the selected token sequence from v_t , preserving
frame-wise temporal alignment via position IDs. Therefore,
the final compressed cache $\{K_t^{i'}, V_t^{i'}\}$ consists of a mix of
salient visual tokens and frame prototypes, enabling both
fine-grained and global memory retention.

4.3. Positional Embedding

MLLMs are typically not extensively trained on long video
sequences due to the scarcity of high-quality, long-form
video-text data. As a result, despite the long context lengths
supported by the underlying language models, MLLMs of-
ten struggle to generalize effectively in long video under-
standing scenarios. To address this limitation, we adopt the
YaRN context window extension technique [33, 45], origi-
nally proposed for language models, to extend the visual
context capacity of MLLMs for streaming video process-
ing.

Prior works on streaming processing of long videos with
MLLMs [8, 20, 30] reassign positional IDs to the visual to-
kens that are retained after KV cache compression. How-
ever, this reassignment discards the original spatial and
temporal information associated with these tokens, poten-
tially degrading performance. We demonstrate that apply-
ing YaRN with a properly chosen scaling factor (based on
the MLLM’s visual context window length) allows us to
preserve positional consistency across streaming segments
and improves performance compared to naively reassigning
position embeddings.

5. Experiments

5.1. Experiment Setup

Benchmarks. We evaluate StreamMem on a number
of widely used offline long video understanding bench-

Method	Frames/FPS	KV Size	MLVU	EgoSchema	VideoMME		
					Medium	Long	All
GPT-4o	-	-	64.6	72.2	70.3	65.3	71.9
MovieChat+	2048	-	25.8	53.5	-	33.4	38.2
Dispider	1 fps	-	61.7	55.6	53.7	49.7	57.2
LongVU	1 fps/400	-	65.4	67.6	58.2	59.5	60.6
LLaVA-OneVision-7B	32	6K	64.7	60.1	54.7	46.2	56.9
+ ReKV [†]	0.5 fps	353K/h	68.5	60.7	-	-	-
+ LiveVLM	0.5/0.2 fps	-	66.3	63.0	56.4	48.8	57.3
+ StreamMem (Ours)	0.5/0.2 fps	6K	66.9	63.0	56.6	50.1	59.4
Qwen2-VL-7B	768	50K	65.8	65.2	-	-	63.9
+ Uniform Sample	-	6K	57.0	64.4	53.3	48.7	58.1
+ SnapKV	768	6K	60.7	62.6	55.3	51.3	59.6
+ InfiniPot-V	768	6K	65.8	65.6	60.8	53.4	62.8
+ StreamMem (Ours)	4.0/0.5 fps	6K	65.9	67.2	62.4	52.3	62.1
Qwen2.5-VL-3B	768	50K	63.3	64.4	-	-	60.3
+ Uniform Sample	-	6K	60.6	62.0	56.9	47.8	58.3
+ InfiniPot-V	768	6K	62.1	61.8	-	-	59.3
+ StreamMem (Ours)	4.0/0.5 fps	6K	62.3	62.2	60.1	49.1	59.5

Table 1. Evaluation results of different MLLMs on offline long video understanding benchmarks. †: ReKV stores the KV cache of all seen frames so it is considered an “upper bound.”

Model	Frames	KV Size	Real-Time Visual Perception							Backward Tracing			Forward Active Responding				Overall	
			OCR	ACR	ATR	STU	FPD	OJR	Avg.	EPM	ASI	HLD	Avg.	REC	SSR	CRR	Avg.	Avg.
GPT-4o	64	-	69.8	64.2	71.6	51.1	70.3	59.8	64.5	57.9	75.7	48.7	60.8	27.6	73.2	59.4	53.4	59.5
Flash-VStream	1fps	-	24.2	29.4	28.5	33.7	25.7	28.8	28.4	39.1	37.2	5.9	27.4	8.02	67.25	60.0	45.1	33.6
Dispider	1fps	-	57.7	49.5	62.1	44.9	61.4	51.6	54.6	48.5	55.4	4.3	36.1	18.1	37.4	48.8	34.7	41.8
LongVU	1fps	-	53.7	53.2	62.9	47.8	68.3	59.8	57.6	40.7	59.5	4.8	35.0	12.2	69.5	60.8	47.5	46.7
LLaVA-OneVision	64	12K	66.4	57.8	73.3	53.4	71.3	62.0	64.0	54.2	55.4	21.5	43.7	25.6	67.1	58.8	50.5	52.7
+StreamMem	1fps	6K	70.5	56.0	71.6	49.4	70.3	63.6	63.6	53.9	59.5	17.7	43.7	26.1	66.1	61.3	51.2	52.8
Qwen2-VL	64	6K	60.4	50.5	56.0	47.2	66.3	55.4	56.0	47.8	56.1	35.5	46.5	31.7	65.8	48.8	48.7	50.4
+StreamMem	1fps	6K	71.8	56.0	64.7	47.8	65.4	60.9	61.1	48.8	58.8	32.3	46.6	31.7	65.0	55.8	50.8	52.8

Table 2. Evaluation results on OVO-Bench. All evaluated open-source models are 7B in size.

Method	RVS-Ego		RVS-Movie	
	Acc	Score	Acc	Score
ReKV	63.7	4.0	54.4	3.6
ReKV w/o offloading.	55.8	3.3	50.8	3.4
Flash-VStream	57.0	4.0	53.1	3.3
InfiniPot-V	57.9	3.5	51.4	3.5
StreamMem (Ours)	57.6	3.8	52.7	3.4

Table 3. Results of different streaming video question answering methods on RVS-Ego and RVS-Movie benchmarks.

Method	KV Size	Holistic	S.D.	M.D.	All
Full KV	50K	76.3	73.9	43.3	65.9
InfiniPot-V	6K	77.2	72.3	44.8	65.8
StreamMem (Ours)	6K	77.5	72.7	44.4	65.9
InfiniPot-V	12K	76.9	73.4	44.0	66.0
StreamMem (Ours)	12K	77.7	73.1	43.9	66.0
InfiniPot-V	24K	76.9	74.0	42.2	65.7
StreamMem (Ours)	24K	77.6	73.4	44.5	66.3

Table 4. Comparison of InfiniPot-V and StreamMem on MLVU with different KV sizes. We use Qwen2VL-7B as the base MLLM.

360 marks, including MLVU [60], EgoSchema [29], and
 361 VideoMME [10]. By default, we process the video stream at
 362 0.5 frames per second (FPS), in accordance with ReKV [8].
 363 For MLVU and EgoSchema, we report the results on the of-

ficial “dev” set. For VideoMME, we report results without
 subtitles. Each video clip is set to 8 frames. In addition
 to offline video understanding benchmarks, we also evalu-
 ate on RVS-Ego and RVS-Movie, two streaming video un-
 364
 365
 366
 367

Algorithm 1 Streaming Video Encoding and KV Cache Compression**Require:** Total KV Cache size M , Template tokens Q .Initialize cache K, V and score matrix s (one row for each transformer layer).**while** not end of video **do**Fetch a new batch of frames v_i from stream. $v'_i = \text{Filter}(v_i)$ based on frame similarity $K_i, V_i, s_i = \text{Encode}(v'_i, Q)$ **if** $|K| > M$ **then** $\mathcal{I} = \text{Topk}(s, k = M); K, V = K[\mathcal{I}], V[\mathcal{I}]$ **end if**Append K_i, V_i, s_i to K, V, s . // Equation 1Insert Merge(K_i), Merge(V_i) to K, V . // Equation 2**end while**

368 derstanding benchmarks, and OVO-Bench[31], a real-world
 369 online video understanding benchmark. All experiments
 370 can be run with one A100 GPU.

371 **Models.** We apply our method on three popular open-
 372 source MLLMs: LLaVA-OneVision-7B [21], Qwen2-VL-
 373 7B [42], and Qwen2.5-VL-3B [2].

374 **Baselines.** We evaluate StreamMem against strong base-
 375 lines, including:

- 376 • Query-agnostic streaming video-language understanding
- 377 with MLLMs, namely LiveVLM [30] and InfiniPot-
- 378 V [20], two recent streaming methods that perform KV
- 379 cache compression independently of the query.
- 380 • Online video MLLMs such as MovieChat+ [37] and
- 381 Dispider [34].
- 382 • LongVU [35], an MLLM that utilizes visual token com-
- 383 pression for long video understanding.

384 We also report the performance of simple uniform frame
 385 sampling for Qwen2-7B and Qwen2.5-3B models, and
 386 ReKV [8] for LLaVA-OneVision, which stores the KV
 387 cache of all previously seen frames without compression.
 388 While ReKV is not feasible in memory-constrained settings
 389 for long videos, it serves as an oracle-style upper bound on
 390 performance under unbounded memory.

5.2. Main Results

392 **Offline video understanding.** We report the main re-
 393 sults for offline video understanding benchmarks in Table 1.
 394 For experiments with LLaVA-OneVision, we sample videos
 395 shorter than 30 minutes at 0.5 fps and videos longer than 30
 396 minutes at 0.2 fps and constrain the GPU memory alloca-
 397 tion below 24 GB, following the setup of LiveVLM [30].
 398 For Qwen2-VL and Qwen2.5-VL experiments, we sample
 399 video less than 3 minutes at 4.0 fps (to match the uniform

sampling of 768 frames in InfiniPot-V [20] and other videos
 at 0.5 fps. We keep the KV cache size at 6K per transformer
 layer in the MLLM.

From the results we observe that StreamMem outper-
 forms the baselines on all benchmarks except the “long”
 subset of VideoMME for Qwen2-VL-7B. On LLaVA-
 OneVision-7B, StreamMem significantly outperforms the
 uniform sampling baseline with a comparable KV cache
 size. This highlights the benefits of streaming processing
 compared to uniform sampling, where significant informa-
 tion loss can incur in the sampling process. On Qwen2.5-
 VL-3B, StreamMem significantly narrows the gap between
 full KV and compressed KV on the challenging MLVU
 benchmark, showing that StreamMem also works well with
 smaller MLLMs, which are especially suitable for memory-
 constrained settings.

Streaming video understanding. We evaluate Stream-
 Mem on the RVS-Ego and RVS-Movie benchmarks for
 streaming video understanding using LLaVA-OneVision-
 7B [21]. Unlike the offline video understanding bench-
 marks considered earlier, these two datasets pose open-
 ended question answering tasks that require models to re-
 ason over long visual contexts. Following the evaluation pro-
 tocol used by prior work, we assess the generated answers
 using GPT-3.5-turbo-0125, which judges both the accu-
 racy and an alignment score from 1 to 5. The results are
 provided in Table 3. Consistent with InfiniPot-V, we con-
 strain GPU memory usage to stay below 28 GB. For ReKV
 without CPU offloading, it simply discards older KVs and
 retains only recent context as “short-term memory.” The
 performance drop between ReKV with and without CPU
 offloading underscores the importance of maintaining long-
 range memory for high-quality answers.

StreamMem outperforms ReKV without offloading and
 is competitive with InfiniPot-V and Flash-VStream, demon-
 strating its effectiveness in open-ended question answering
 under constrained memory settings. These results highlight
 the method’s ability to retain and utilize salient long-term
 information throughout streaming video.

Online Video Understanding on OVO-Bench. We eval-
 uate StreamMem on OVO-Bench [31], a comprehensive
 benchmark for real-world online video understanding, and
 report the results in Table 2. Results show that enabling
 streaming video understanding with StreamMem consis-
 tently improves both LLaVA-OneVision and Qwen2-VL
 under comparable or more constrained memory budgets.
 For Qwen2-VL, integrating StreamMem at the same KV
 cache size (6K) yields a clear overall gain, improving
 the average performance from 50.4 to 52.8. For LLaVA-
 OneVision, StreamMem achieves competitive overall per-
 formance while reducing the KV cache size from 12K

Query Type	Holistic	S.D.	M.D.	All
True Query	80.8	71.6	46.4	68.1
Generic Text Query	78.1	71.3	43.0	66.7
Chat Template Query	78.8	71.5	43.0	66.9

Table 5. Ablation study on different proxy queries for attention-based KV compression.

KV Merging Strategy	Holistic	S.D.	M.D.	All
No Merging	77.3	69.7	42.8	65.6
Avg. Merging	80.5	70.4	41.3	66.3
Weighted Merging	78.8	71.5	43.0	66.9

Table 6. Ablation study on the effect of different KV merging strategies.

Memory Allocation	Holistic	S.D.	M.D.	All
Uniform	78.8	71.5	43.0	66.9
Inverse Entropy	80.3	70.4	41.3	66.3
Exponential Decay (0.97)	78.4	71.0	42.2	66.4
Exponential Decay (0.98)	80.3	70.1	42.0	66.2

Table 7. Comparison of different memory allocation strategies on MLVU with LLaVA-OneVision-7B.

451 to 6K, indicating substantially better memory efficiency.
 452 Overall, these results highlight that StreamMem boosts accuracy
 453 of frontier MLLMs over a wide range of real-world
 454 online video understanding tasks while keeping a tight KV
 455 cache constraint.

456 5.3. Ablation Studies

457 We conduct ablation studies on different components in our
 458 method. For these experiments, we use LLaVA-OneVision-
 459 7B [21] on the MLVU benchmark [60]. We report the average
 460 performance on the three subsets of the MLVU, namely
 461 holistic tasks (including Topic Reasoning and Anomaly
 462 Recognition), single detail (Needle QA, Ego Reasoning,
 463 and Plot QA), and multi-detail (including Action Order, and
 464 Action Count).

465 **Type of proxy query.** We compare the results for using
 466 different queries, including the ground truth query, the chat
 467 template query, and a generic text query (“What is happen-
 468 ing in the video?”) for the attention-based KV compression
 469 module in Table 5. We observe that the generic text query
 470 obtains similar performance to the chat template query, sug-
 471 gesting that the chat template query, while not including any
 472 real text, is implicitly acting as a generic query of video con-
 473 tent. Using the ground truth user query for KV compression
 474 still significantly outperforms the query-agnostic methods,

especially in “multi-detail” tasks, showing the challenge for
 query-agnostic methods to retain all the details required to
 answer the question without knowing the question during
 video processing.

Merging strategy. We compare the results for different
 KV merging strategies in Table 6. We observe that all
 frame-wise KV cache merging methods perform better than
 no KV merging, confirming the results from LiveVLM [30].
 StreamMem improved over LiveVLM which uses simple
 average merging and inserting to the end of each frame by
 applying weighted merging based on the attention scores
 between the chat template tokens and the visual tokens.

Memory Allocation Strategies. Table 7 compares sev-
 eral simple memory-budget allocation strategies, including
 entropy-based and decay-based weighting schemes. Results
 show that, despite prior observations suggesting that bias-
 ing toward earlier layers can be beneficial [4, 44], a uni-
 form allocation strategy consistently outperforms these al-
 ternatives. This suggests that uniform memory allocation
 remains a strong and robust baseline for managing limited
 memory budgets in streaming video settings.

496 6. Conclusion

Enabling continuous video stream processing under a
 bounded memory constraint is essential for deploying
 multimodal large language models (MLLMs) in real-
 world, embodied scenarios. Yet, most prior work in long
 video-language understanding has focused on static or
 offline settings, assuming known queries, finite video
 lengths, and full access to the visual context in advance.
 These assumptions limit their applicability in streaming
 or open-world environments. In this work, we present
 StreamMem, a training-free and query-agnostic KV cache
 compression framework tailored for streaming video under-
 standing. By using attention scores between visual tokens
 and chat template tokens as a proxy for query relevance,
 StreamMem effectively retains salient visual information
 without requiring access to future queries. When applied to
 open-source MLLMs, StreamMem achieves state-of-the-art
 performance across a diverse set of offline and streaming
 long video benchmarks. Beyond demonstrating competi-
 tive empirical results, we conduct an in-depth analysis
 of various components in our framework, including KV
 merging strategies and positional embedding techniques,
 shedding light on the design considerations for construct-
 ing a memory-bounded visual processing pipeline. These
 insights lay a foundation for future research in scaling
 MLLMs to continuously process real-world visual streams.

523

References

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

- [1] Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S Nikolopoulos, Hans Vandierendonck, Deepu John, and Bo Ji. Hired: Attention-guided token dropping for efficient inference of high-resolution vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1773–1781, 2025. 3
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 2, 7
- [3] Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng Gao, Dongxing Mao, and Mike Zheng Shou. Videollm-online: Online video large language model for streaming video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18407–18418, 2024. 1, 2
- [4] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *European Conference on Computer Vision*, pages 19–35. Springer, 2024. 3, 5, 8
- [5] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 2
- [6] Giulio Corallo, Orion Weller, Fabio Petroni, and Paolo Papotti. Beyond rag: Task-aware kv cache compression for comprehensive knowledge reasoning. *arXiv preprint arXiv:2503.04973*, 2025. 3
- [7] Alessio Devoto, Yu Zhao, Simone Scardapane, and Pasquale Minervini. A simple and effective L_2 norm-based strategy for kv cache compression. *arXiv preprint arXiv:2406.11430*, 2024. 3
- [8] Shangzhe Di, Zhelun Yu, Guanghao Zhang, Haoyuan Li, Tao Zhong, Hao Cheng, Bolin Li, Wanggui He, Fangxun Shu, and Hao Jiang. Streaming video question-answering with in-context video kv-cache retrieval. *arXiv preprint arXiv:2503.00540*, 2025. 2, 4, 5, 6, 7
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021. 4
- [10] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118, 2025. 2, 6
- [11] Yu Fu, Zefan Cai, Abedelkadir Asi, Wayne Xiong, Yue Dong, and Wen Xiao. Not all heads matter: A head-level kv cache compression method with integrated retrieval and reasoning. *arXiv preprint arXiv:2410.19258*, 2024. 3
- [12] Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. Model tells you what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*, 2023. 3
- [13] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18995–19012, 2022. 4
- [14] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. Ma-Imm: Memory-augmented large multimodal model for long-term video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13504–13514, 2024. 1, 5
- [15] Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Monishwaran Maheswaran, June Paik, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. Squeezed attention: Accelerating long context length llm inference. *arXiv preprint arXiv:2411.09688*, 2024. 3
- [16] Kai Hu, Feng Gao, Xiaohan Nie, Peng Zhou, Son Tran, Tal Neiman, Lingyun Wang, Mubarak Shah, Raffay Hamid, Bing Yin, et al. M-llm based video frame selection for efficient video understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13702–13712, 2025. 1
- [17] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1
- [18] Jang-Hyun Kim, Jinuk Kim, Sangwoo Kwon, Jae W Lee, Sangdoon Yun, and Hyun Oh Song. Kvzip: Query-agnostic kv cache compression with context reconstruction. *arXiv preprint arXiv:2505.23416*, 2025. 3
- [19] Minsoo Kim, Kyuhong Shim, Jungwook Choi, and Simyung Chang. Infinipot: Infinite context processing on memory-constrained llms. *arXiv preprint arXiv:2410.01518*, 2024. 3
- [20] Minsoo Kim, Kyuhong Shim, Jungwook Choi, and Simyung Chang. Infinipot-v: Memory-constrained kv cache compression for streaming video understanding. *arXiv preprint arXiv:2506.15745*, 2025. 1, 2, 5, 7
- [21] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 2, 7, 8
- [22] 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. *Advances in Neural Information Processing Systems*, 37:22947–22970, 2024. 3
- [23] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In 579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636

- 637 *European Conference on Computer Vision*, pages 323–340. 695
638 Springer, 2024. 1 696
- 639 [24] Jianxin Liang, Xiaojun Meng, Yueqian Wang, Chang Liu, 697
640 Qun Liu, and Dongyan Zhao. End-to-end video question an- 698
641 swering with frame scoring mechanisms and adaptive sam- 699
642 pling. *arXiv preprint arXiv:2407.15047*, 2024. 1 700
- 643 [25] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 701
644 Visual instruction tuning. *Advances in neural information 702*
645 *processing systems*, 36:34892–34916, 2023. 2 703
- 646 [26] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 704
647 Improved baselines with visual instruction tuning. In *Pro- 705*
648 *ceedings of the IEEE/CVF conference on computer vision 706*
649 *and pattern recognition*, pages 26296–26306, 2024. 2 707
- 650 [27] Xiangrui Liu, Yan Shu, Zheng Liu, Ao Li, Yang Tian, and 708
651 Bo Zhao. Video-xl-pro: Reconstructive token compres- 709
652 sion for extremely long video understanding. *arXiv preprint 710*
653 *arXiv:2503.18478*, 2025. 2 711
- 654 [28] Xuyang Liu, Yiyu Wang, Junpeng Ma, and Linfeng Zhang. 712
655 Video compression commander: Plug-and-play inference ac- 713
656 celeration for video large language models. *arXiv preprint 714*
657 *arXiv:2505.14454*, 2025. 3 715
- 658 [29] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra 716
659 Malik. Egoschema: A diagnostic benchmark for very long- 717
660 form video language understanding. *Advances in Neural In- 718*
661 *formation Processing Systems*, 36:46212–46244, 2023. 2, 719
662 6 720
- 663 [30] Zhenyu Ning, Guangda Liu, Qihao Jin, Wenchao Ding, 721
664 Minyi Guo, and Jieru Zhao. LiveVlm: Efficient online video 722
665 understanding via streaming-oriented kv cache and retrieval. 723
666 *arXiv preprint arXiv:2505.15269*, 2025. 2, 4, 5, 7, 8 724
- 667 [31] Junbo Niu, Yifei Li, Ziyang Miao, Chunjiang Ge, Yuanhang 725
668 Zhou, Qihao He, Xiaoyi Dong, Haodong Duan, Shuangrui 726
669 Ding, Rui Qian, et al. Ovo-bench: How far is your video- 727
670 llms from real-world online video understanding? In *Pro- 728*
671 *ceedings of the Computer Vision and Pattern Recognition 729*
672 *Conference*, pages 18902–18913, 2025. 2, 7 730
- 673 [32] Junyoung Park, Dalton Jones, Matthew J Morse, Raghavv 731
674 Goel, Mingu Lee, and Chris Lott. Keydiff: Key similarity- 732
675 based kv cache eviction for long-context llm inference 733
676 in resource-constrained environments. *arXiv preprint 734*
677 *arXiv:2504.15364*, 2025. 3 735
- 678 [33] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico 736
679 Shippole. Yarn: Efficient context window extension of large 737
680 language models. *arXiv preprint arXiv:2309.00071*, 2023. 5 738
- 681 [34] Rui Qian, Shuangrui Ding, Xiaoyi Dong, Pan Zhang, Yuhang 739
682 Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. Dispider: 740
683 Enabling video llms with active real-time interaction via dis- 741
684 entangled perception, decision, and reaction. In *Proceedings 742*
685 *of the Computer Vision and Pattern Recognition Conference*, 743
686 pages 24045–24055, 2025. 2, 7 744
- 687 [35] Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng 745
688 Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyu Xiao, Bal- 746
689 akrishnan Varadarajan, Florian Bordes, et al. Longvu: Spa- 747
690 tiotemporal adaptive compression for long video-language 748
691 understanding. *arXiv preprint arXiv:2410.17434*, 2024. 1, 749
692 2, 5, 7 750
- 693 [36] Yan Shu, Zheng Liu, Peitian Zhang, Minghao Qin, Jun- 751
694 jie Zhou, Zhengyang Liang, Tiejun Huang, and Bo Zhao. 752
Video-xl: Extra-long vision language model for hour-scale
video understanding. In *Proceedings of the Computer Vision
and Pattern Recognition Conference*, pages 26160–26169,
2025. 2
- [37] Enxin Song, Wenhao Chai, Tian Ye, Jenq-Neng Hwang, Xi
Li, and Gaoang Wang. Moviechat+: Question-aware sparse
memory for long video question answering. *arXiv preprint
arXiv:2404.17176*, 2024. 7
- [38] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen
Bo, and Yunfeng Liu. Roformer: Enhanced transformer with
rotary position embedding. *Neurocomputing*, 568:127063,
2024. 1
- [39] Reuben Tan, Ximeng Sun, Ping Hu, Jui-hsien Wang, Hanieh
Deilamsalehy, Bryan A Plummer, Bryan Russell, and Kate
Saenko. Koala: Key frame-conditioned long video-llm. In
*Proceedings of the IEEE/CVF Conference on Computer Vi-
sion and Pattern Recognition*, pages 13581–13591, 2024. 1
- [40] Xi Tang, Jihao Qiu, Lingxi Xie, Yunjie Tian, Jianbin Jiao,
and Qixiang Ye. Adaptive keyframe sampling for long video
understanding. In *Proceedings of the Computer Vision and
Pattern Recognition Conference*, pages 29118–29128, 2025.
1
- [41] Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan
Wang. Dycoko: Dynamic compression of tokens for fast
video large language models. In *Proceedings of the Com-
puter Vision and Pattern Recognition Conference*, pages
18992–19001, 2025. 1
- [42] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
Ge, et al. Qwen2-vl: Enhancing vision-language model’s
perception of the world at any resolution. *arXiv preprint
arXiv:2409.12191*, 2024. 2, 7
- [43] Xiao Wang, Qingyi Si, Jianlong Wu, Shiyu Zhu, Li Cao, and
Liqiang Nie. Retake: Reducing temporal and knowledge
redundancy for long video understanding. *arXiv preprint
arXiv:2412.20504*, 2024. 3
- [44] Xiao Wang, Qingyi Si, Jianlong Wu, Shiyu Zhu, Li Cao,
and Liqiang Nie. Adaretake: Adaptive redundancy reduction
to perceive longer for video-language understanding. *arXiv
preprint arXiv:2503.12559*, 2025. 3, 5, 8
- [45] Hongchen Wei and Zhenzhong Chen. Visual context window
extension: A new perspective for long video understanding.
arXiv preprint arXiv:2409.20018, 2024. 5
- [46] Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang,
Prajwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta,
Karthik Abinav Sankararaman, Barlas Oguz, et al. Effective
long-context scaling of foundation models. *arXiv preprint
arXiv:2309.16039*, 2023. 1
- [47] Yuhui Xu, Zhanming Jie, Hanze Dong, Lei Wang, Xudong
Lu, Aojun Zhou, Amrita Saha, Caiming Xiong, and Doyen
Sahoo. Think: Thinner key cache by query-driven pruning.
arXiv preprint arXiv:2407.21018, 2024. 3
- [48] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and
Lucas Beyer. Sigmoid loss for language image pre-training.
In *Proceedings of the IEEE/CVF international conference on
computer vision*, pages 11975–11986, 2023. 4
- [49] Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu,
Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang,

- 753 Hang Zhang, Xin Li, et al. Videollama 3: Frontier multi-
754 modal foundation models for image and video understand-
755 ing. *arXiv preprint arXiv:2501.13106*, 2025. 1
- 756 [50] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi
757 Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-
758 based real-time understanding for long video streams. *arXiv*
759 *preprint arXiv:2406.08085*, 2024. 2
- 760 [51] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng,
761 Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan,
762 Chunyuan Li, and Ziwei Liu. Long context transfer from
763 language to vision. *arXiv preprint arXiv:2406.16852*, 2024.
764 1, 2
- 765 [52] Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiy-
766 ong Zhuo, Jiajun Cao, Shaobo Guo, Qi She, and Shang-
767 hang Zhang. Beyond text-visual attention: Exploiting vi-
768 sual cues for effective token pruning in vlms. *arXiv preprint*
769 *arXiv:2412.01818*, 2024. 3
- 770 [53] Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang
771 Feng. Llava-mini: Efficient image and video large mul-
772 timodal models with one vision token. *arXiv preprint*
773 *arXiv:2501.03895*, 2025. 1
- 774 [54] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng,
775 Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki
776 Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Vi-
777 sual token sparsification for efficient vision-language model
778 inference. *arXiv preprint arXiv:2410.04417*, 2024. 3
- 779 [55] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Zi-
780 wei Liu, and Chunyuan Li. Video instruction tuning with
781 synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 1
- 782 [56] Yiming Zhang, Zhuokai Zhao, Zhaorun Chen, Zenghui
783 Ding, Xianjun Yang, and Yining Sun. Beyond training:
784 Dynamic token merging for zero-shot video understanding.
785 *arXiv preprint arXiv:2411.14401*, 2024. 1, 3, 5
- 786 [57] Yiming Zhang, Zhuokai Zhao, Zhaorun Chen, Zhili Feng,
787 Zenghui Ding, and Yining Sun. Rankclip: Ranking-
788 consistent language-image pretraining. *arXiv preprint*
789 *arXiv:2404.09387*, 2024. 4
- 790 [58] Yiming Zhang, Chengzhang Yu, Zhuokai Zhao, Kun Wang,
791 Qiankun Li, Zihan Chen, Yang Liu, Zenghui Ding, and
792 Yining Sun. Circuitprobe: Dissecting spatiotemporal
793 visual semantics with circuit tracing. *arXiv preprint*
794 *arXiv:2507.19420*, 2025. 1
- 795 [59] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen,
796 Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian,
797 Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter ora-
798 cle for efficient generative inference of large language mod-
799 els. *Advances in Neural Information Processing Systems*, 36:
800 34661–34710, 2023. 3
- 801 [60] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang
802 Liang, Shitao Xiao, Minghao Qin, Xi Yang, Yongping
803 Xiong, Bo Zhang, et al. Mlvu: Benchmarking multi-task
804 long video understanding. In *Proceedings of the Computer*
805 *Vision and Pattern Recognition Conference*, pages 13691–
806 13701, 2025. 2, 6, 8