# StreamingTOM: Streaming Token Compression for Efficient Video Understanding

Xueyi Chen[1,2], Keda Tao[1,3], Kele Shao[1,4,3], Huan Wang[1,*]

[1]Westlake University, [2]The Chinese University of Hong Kong, [3]Zhejiang University, [4]SII
[*]Corresponding author: wanghuan@westlake.edu.cn
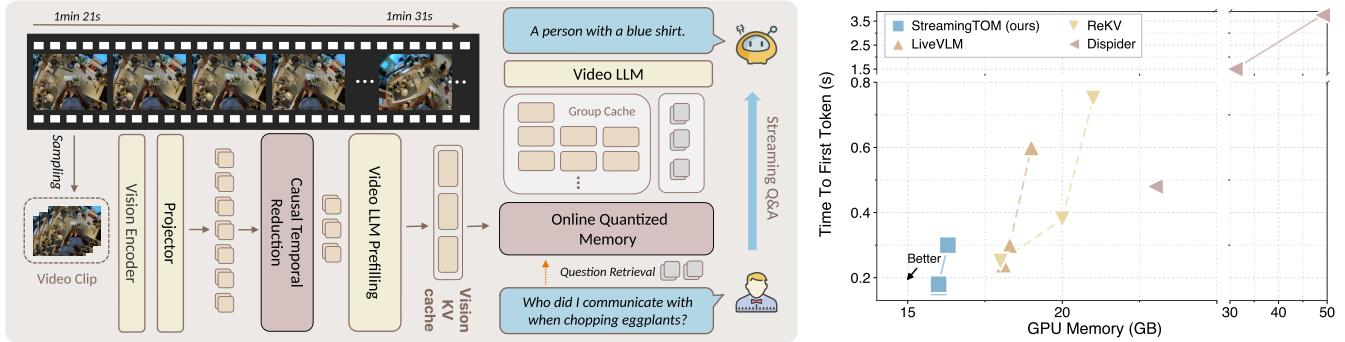https://yige24.github.io/StreamingTOM

Figure 1. **Left**: StreamingTOM (streaming token compression) is a training-free, two-stage framework for efficient streaming video understanding: *Causal Temporal Reduction* selects pre-LLM tokens with a causal, fixed-budget policy, and *Online Quantized Memory* maintains compact kv-cache via 4-bit quantization with on-demand retrieval. **Right**: on LLaVA-OV-7B, StreamingTOM achieves a 15.7× kv-cache compression ratio; compared to prior SOTA (LiveVLM), it delivers 1.2× lower peak memory and 2× faster TTFT.

## Abstract

*Unlike offline processing, streaming video vision-language models face two fundamental constraints: causality and accumulation. Causality prevents access to future frames that offline methods exploit, while accumulation causes tokens to grow unbounded, creating efficiency bottlenecks. However, existing approaches only regulate post-LLM kv-cache, leaving costly pre-LLM prefill unchanged. We introduce StreamingTOM, a training-free, plug-and-play two-stage framework that addresses both pre-LLM and post-LLM bottlenecks with predictable latency. Causal Temporal Reduction imposes a fixed per-frame budget and selects tokens based on adjacent-frame changes and token saliency, drastically reducing per-frame prefill cost by processing only a compact subset of visual tokens per frame instead of all visual tokens. Online Quantized Memory stores tokens in 4-bit format, retrieves relevant groups on demand, and dequantizes them, keeping the active kv-cache bounded regardless of stream length. Experiments demonstrate our method achieves 15.7× kv-cache compression, 1.2× lower peak memory and 2× faster TTFT compared to prior SOTA. StreamingTOM maintains state-of-the-art accuracy among training-free methods with an average of 63.8% on offline benchmarks and 55.8%/3.7 on RVS. These results highlight the practical benefits of our two-stage approach for efficient streaming video understanding with bounded growth.*

## 1. Introduction

Recent advances in multimodal and large language models have demonstrated impressive capabilities across a diverse range of vision and language tasks [1, 11, 19, 36, 37]. Building on these foundations, offline video understanding investigates models for deep semantic reasoning on prerecorded videos with known temporal boundaries [2, 3, 14, 33, 34, 45]. A dominant efficiency bottleneck in these models is the computational cost of processing the heavy spatiotemporal redundancy in visual tokens. This motivated established compression techniques such as spatiotemporal token merging and dynamic pruning [16, 25], which leverage global context, including future frames, to reduce computation with little or no loss in accuracy [6, 24, 29, 41]. These techniques inherently assume access to a global context and stable clip boundaries.

While offline techniques have matured, emerging applications like autonomous driving, embodied AI, and live video assistants demand streaming processing, where two constraints emerge: causality and accumulation [15, 35, 44]. Causality means the system cannot access future frames. Accumulation arises because the effective sequence length $L(t)$ grows as frames arrive at a fixed frame rate.

Dense attention grows quadratically with $L(t)$, while the kv-cache grows linearly with time [18]. As compute and memory grow unbounded over time, token compression shifts from an optional offline optimization to a fundamental prerequisite for streaming viability under fixed budgets.

Contemporary streaming methods either require costly model-specific training [4, 40, 44] or operate training-free but focus solely on post-LLM kv-cache management [7, 18, 42]. While the latter reduces memory, it preserves the $O(tNLd^2)$ pre-LLM computation as all visual tokens must traverse the full transformer stack. Thus, training-free and strictly causal pre-LLM token reduction, which is essential for true streaming efficiency, remains largely unexplored.

We propose **StreamingTOM**, a training-free framework that reduces visual tokens before the LLM for streaming video understanding. Our central insight is that temporal redundancy is abundant in both offline and streaming videos, but streaming combines it with strict causality and ever-growing context. To address this while respecting streaming constraints, effective compression must be pre-LLM and strictly causal. We impose a fixed per-frame budget $G$ to stabilize latency and ensure predictable compute and memory consumption across all frames. Accordingly, we implement this design with a strictly causal one-pass selector that makes decisions from adjacent frames $t-1$ and $t$. **Causal Temporal Reduction (CTR)** addresses the prefill bottleneck by selecting $G$ tokens per frame based on adjacent-frame changes and token saliency, reducing per-frame prefill cost by a factor of $N/G$ through processing only $G$ selected tokens instead of all $N$ visual tokens.[1] Meanwhile, **Online Quantized Memory (OQM)** bounds post-LLM memory growth by storing retained tokens in 4-bit format and retrieving relevant groups on demand, keeping the active kv-cache bounded regardless of stream length. Together, these two stages enable long-horizon stability under fixed budgets with predictable latency, while remaining plug-and-play across different backbones.

StreamingTOM achieves $15.7\times$ kv-cache compression ratio and delivers $1.2\times$ lower peak memory and $2\times$ faster time to first token (TTFT) compared to prior training-free SOTA (LiveVLM [18]) on streaming setting. It maintains state-of-the-art accuracy among training-free methods with $63.8\%$ average on VideoMME [9], MLVU [46], and EgoSchema [17], and $55.8\%/3.7$ on RVS [44]. These gains persist over long horizons as our pre-LLM token reduction and bounded kv-cache prevent unbounded growth with video length. On LLaVA-OV-7B, a one-hour stream reduces kv-cache memory from $18.8\,\mathrm{GB}$ to $1.2\,\mathrm{GB}$, confirming bounded growth over extended sessions.

In summary, our contributions are threefold.

- We establish that effective streaming compression re-

quires streaming-native design, introducing a strictly causal approach that reduces complexity from $O(tNLd^2)$ to $O(tGLd^2)$ with fixed per-frame budget $G$.
- We propose StreamingTOM, a training-free two-stage framework combining CTR for pre-LLM compression with OQM for bounded post-LLM memory management.
- We achieve state-of-the-art training-free accuracy ($63.8\%$ offline average, $55.8\%/3.7$ on RVS) with $15.7\times$ kv-cache compression, $1.2\times$ lower peak memory, and $2\times$ faster TTFT over training-free SOTA, demonstrating both effectiveness and efficiency.

## 2. Related Work

**Video Token Compression.** Training-free methods achieve substantial token reduction via hierarchical approaches [38, 41], adaptive pruning [12, 23], dynamic compression [26, 29], and holistic merging [24]. However, these approaches rely on global temporal analysis or bidirectional attention that leverages future frames. This look-ahead capability is incompatible with streaming causality, where compression must depend only on past context and cannot be revised. Training-based long-video compression methods also exist [27], but they require model retraining and are outside our training-free scope. Beyond token-level reduction, quantization techniques [21, 30, 31] offer memory compression through reduced precision, complementing spatial and temporal compression strategies.

**Streaming Video Understanding.** Training-based methods [4, 5, 8, 10, 22, 32, 39, 40, 44] achieve strong performance through response-timing control, memory decoupling, perception-decision disentanglement, multi-round interaction, visual feedback, and speech integration, but require costly model-specific retraining. Notably, TimeChat-Online [43] observes that 80% of visual tokens are redundant in streaming videos, motivating aggressive compression. In contrast, training-free approaches [7, 13, 18, 42] focus exclusively on post-LLM memory optimization through kv-cache retrieval, eviction, and compression strategies. However, all such methods preserve the full prefill computation cost because tokens have already traversed the full transformer stack, and the pre-LLM prefill remains intact.

To our knowledge, no prior training-free streaming method performs strictly causal pre-LLM token reduction. This gap is critical because pre-LLM compression cuts the prefill cost by reducing token count before the transformer, while post-LLM optimization only manages memory after the computational bottleneck. We address this with StreamingTOM, which couples Causal Temporal Reduction for pre-LLM compression with Online Quantized Memory for efficient kv-cache management under strict causality.

---

[1]**Notation:** $N$=tokens per frame, $L$=number of transformer layers, $d$=hidden width, $G$=per-frame retained tokens selected by CTR.
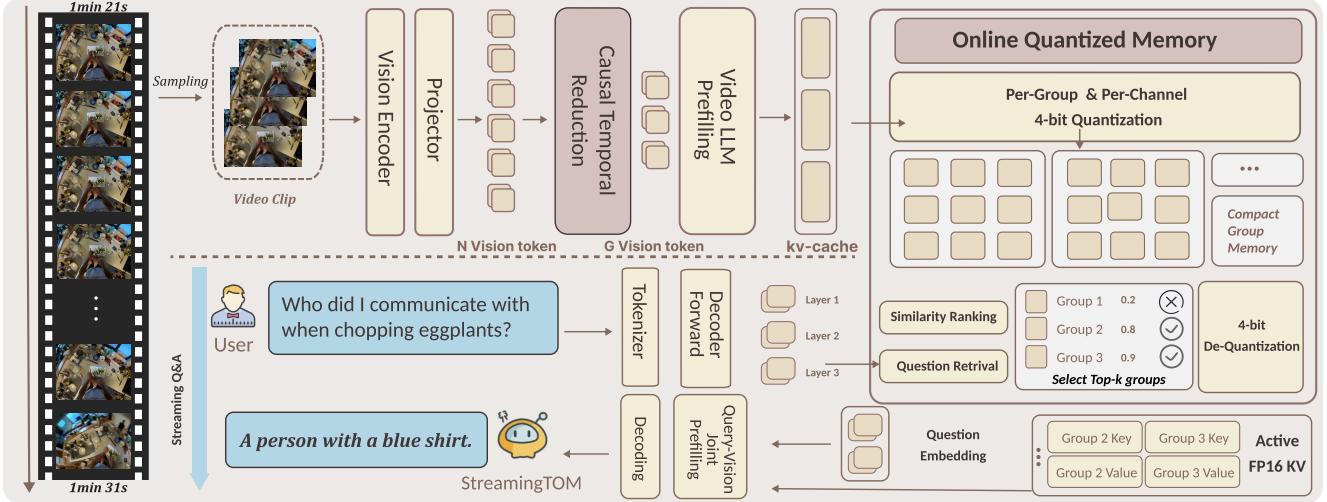
Figure 2. Overview of the StreamingTOM architecture. The framework consists of two coordinated pipelines. The vision pipeline encodes each frame and applies Causal Temporal Reduction to condense redundant tokens into compact groups, which are written to an online memory for reuse. The query pipeline processes questions and drives the decoder to interact with the memory through Online Quantized Memory, which stores groups at 4-bit precision, retrieves at most $k$ groups on demand, and dequantizes them for efficient generation.

## 3. Methodology

### 3.1. Preliminary

Streaming video understanding must operate under strict causality when processing an unbounded stream $\mathcal{V} = \{v_1, v_2, \ldots, v_t, \ldots\}$: at time $t$, the system can only access past frames $\mathcal{V}_{\leq t}$ while remaining prepared for arbitrary future queries within bounded memory $M$. This constraint rules out iterative refinement or retrospective reprocessing, since both rely on unavailable future information; instead, each frame must be processed in a single pass and finalized immediately. In this context, most training-free streaming adaptations of VLMs adopt kv-cache reuse, leveraging the causal transformer decoder to realize the efficient paradigm of "process once, reuse many".

Under this paradigm, each incoming frame $v_t$ is encoded once into visual tokens $\mathbf{H}_t = \mathcal{E}_v(v_t) \in \mathbb{R}^{N \times d}$ and, across $L$ transformer layers, mapped to key and value pairs $(K_t^{(l)}, V_t^{(l)})$, and the cache accumulates by concatenation $\mathcal{K}_t^{(l)} = [\mathcal{K}_{t-1}^{(l)}; K_t^{(l)}]$ and $\mathcal{V}_t^{(l)} = [\mathcal{V}_{t-1}^{(l)}; V_t^{(l)}]$. This creates a fundamental trade-off: single-pass processing avoids repeated work, but because all tokens traverse the full stack, the per-frame compute remains $O(NLd^2)$; meanwhile, as frames accumulate, the kv-cache grows without bound. This growth can be formalized as

$$\frac{d\,\text{memory}}{dt} = 2 \cdot L \cdot N \cdot d \cdot \text{sizeof(dtype)} \cdot \text{fps}. \quad (1)$$

Integrating Eq. (1) under typical settings yields a total kv-cache footprint of about **18.8 GB** for one hour of video[2],

which substantially exceeds typical GPU memory capacity.

This unbounded growth creates two critical bottlenecks: compute and memory. Computationally, query processing must attend to all $O(tN)$ accumulated tokens; each decoding step therefore requires $O(tNLd^2)$ flops, a prohibitive cost for real-time applications. Post-LLM compression provides no relief because the tokens have already traversed the full transformer stack. On the memory side, the kv-cache grows linearly over time according to Eq. (1), quickly exceeding GPU capacity. Even with advanced eviction strategies, sustaining FP16 precision over long streaming sessions remains infeasible. Consequently, compression must be performed under a strict, query-agnostic, and causal constraint, optimizing both compute and memory without prior knowledge of the future query distribution. Let $\mathcal{L}(q, \mathcal{K}, \mathcal{V})$ denote the task loss when answering query $q$ based on the kv-cache $(\mathcal{K}, \mathcal{V})$ derived from the preceding video frames. The compression objective is then:

$$\min_{(\hat{\mathcal{K}}, \hat{\mathcal{V}}) \in \mathcal{F}_M} \mathbb{E}_{q \sim p} \mathcal{L}(q, \hat{\mathcal{K}}, \hat{\mathcal{V}}), \quad (2)$$

where $p$ denotes the query distribution, $\mathcal{F}_M = \{(\hat{\mathcal{K}}, \hat{\mathcal{V}}) : |\hat{\mathcal{K}}| + |\hat{\mathcal{V}}| \leq M\}$ defines the feasible set for memory budget $M$, and $|\cdot|$ measures memory footprint. The challenge is to compress to mitigate compute and memory growth while preserving causal validity.

### 3.2. Our Method: StreamingTOM

The dual bottlenecks identified above demand a coordinated solution, with computational constraints requiring pre-LLM

---

[2]Based on LLaVA-OV-7B: $L$=28, $N$=196, $H_{\text{kv}}$=4, $d_h$=128, FP16, 0.5 fps (1800 frames/hour): $2 \cdot 28 \cdot 1800 \cdot 196 \cdot 4 \cdot 128 \cdot 2B \approx 18.8\text{GB}$.
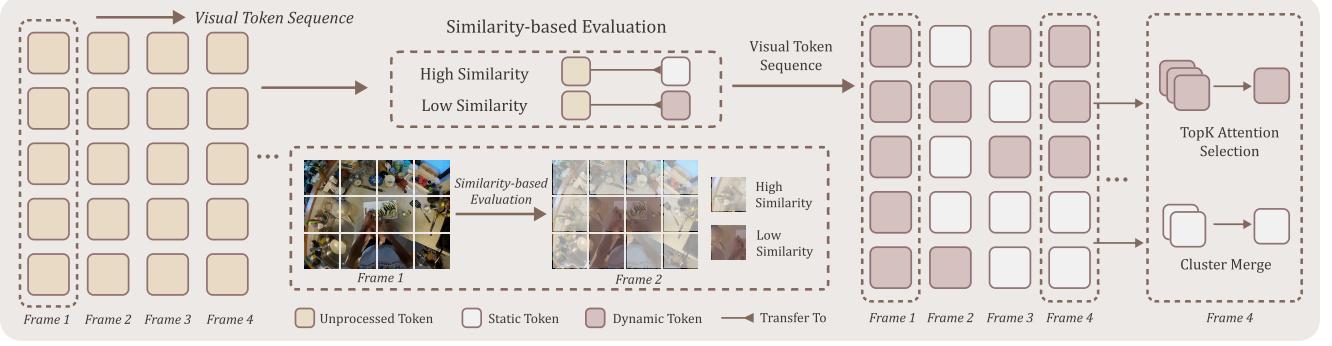
Figure 3. Overview of the CTR compression pipeline. The algorithm processes visual tokens from consecutive frames using a 2-frame window (current and previous), producing a binary classification (static, dynamic) through similarity comparison. The adaptive budget allocation dynamically distributes compression resources based on content, followed by dual-path processing: dpc clustering for static tokens and attention-based selection for dynamic tokens.

compression and memory constraints requiring post-LLM compression. Either approach in isolation is insufficient, as pre-LLM compression cannot stop kv-cache accumulation, while post-LLM compression cannot reduce the computational cost already incurred. The order is also fixed, as tokens must be selected before entering the LLM to reduce computation, and only after LLM processing can they be quantized for storage.

We propose **StreamingTOM** (<u>streaming</u> <u>to</u>ken <u>com</u>pression), a unified two-stage framework that integrates Causal Temporal Reduction (CTR) for pre-LLM computational optimization and Online Quantized Memory (OQM) for post-LLM memory management. As shown in Figure 2, the two modules are connected through a fixed per-frame budget of $G$ tokens that pins both prefill compute and per-frame kv-writes, making per-frame latency predictable. Formally, StreamingTOM is defined as the sequential composition of CTR and OQM:

$$\text{StreamingTOM} = \text{OQM}_{16 \to 4} \circ \text{CTR}_{N \to G}, \quad (3)$$

where $\text{CTR}_{N \to G}$ reduces per-frame tokens from $N$ to a fixed quota $G$ before the LLM, and $\text{OQM}_{16 \to 4}$ quantizes stored tensors from FP16 to 4-bit after the LLM while keeping the hidden width $d$ unchanged.

Central to StreamingTOM is the **group** abstraction, a fixed-size unit of $G$ tokens per frame that serves as the interface between CTR and OQM, decoupling token reduction from storage optimization: CTR selects the most informative $G$ tokens, while OQM handles their efficient storage and retrieval. The fixed budget $G$ ensures predictable per-frame latency and enables stable batching across frames with varying content complexity. Overall, prefill computation reduces from $O(TNLd^2)$ to $O(TGLd^2)$, where $T$ denotes the total number of processed frames, and storage reduces from $O(TN \cdot d \cdot 16)$ bits to $O(TG \cdot d \cdot 4)$ bits, resulting in a combined compression ratio of $\frac{4N}{G}$.

## 3.3. Causal Temporal Reduction (CTR)

CTR addresses these streaming constraints through three design principles: strict causality with a 2-frame window, single-pass processing, and a fixed per-frame budget $G$ for predictable latency. As illustrated in Figure 3, the CTR pipeline evaluates temporal similarity between consecutive frames, classifies tokens into static/dynamic sets, and applies dual-path processing with adaptive budget allocation.

CTR relies on two complementary signals: temporal similarity $s_t^{(i)}$ captures cross-frame redundancy, while spatial saliency $\alpha_t^{(i)}$ identifies informative content within the current frame. Since the encoder maintains fixed patch indexing across frames, tokens at position $i$ in frames $t$ and $t-1$ correspond to the same spatial location. For temporal similarity, given features $\mathcal{F}_t, \mathcal{F}_{t-1} \in \mathbb{R}^{N \times d}$, CTR measures per-token consistency by cosine similarity,

$$s_t^{(i)} = \frac{\mathcal{F}_t^{(i)} \cdot \mathcal{F}_{t-1}^{(i)}}{\|\mathcal{F}_t^{(i)}\| \, \|\mathcal{F}_{t-1}^{(i)}\|}. \quad (4)$$

For spatial saliency, CTR leverages pre-computed attention-based scores $\alpha_t^{(i)} \in [0, 1]$ from the vision encoder as a zero-cost byproduct. To avoid memory peaks in streaming scenarios, these scores are computed via chunked attention without materializing the full $N \times N$ attention matrix.

To partition tokens based on temporal similarity, CTR induces a measurable decomposition of the index space $[N]$ via threshold-parameterized indicators $I_s^{(i)}, I_d^{(i)} : \mathbb{R} \to \{0, 1\}$, yielding disjoint static/dynamic subsets,

$$
\begin{aligned}
& I_s^{(i)} = \Bbbk[s_t^{(i)} > \tau_c], \ I_d^{(i)} = 1 - I_s^{(i)}, \\
& \mathcal{S}_t = \{i \in [N] \mid I_s^{(i)} = 1\}, \ \mathcal{D}_t = [N] \setminus \mathcal{S}_t
\end{aligned}
\quad (5)
$$

where $\tau_c$ is the similarity threshold. This deterministic partition enables compositional budget allocation. To compress from $N$ tokens to exactly $G$ tokens per frame, CTR

allocates the budget proportionally based on the observed static/dynamic distribution,

$$k_s = \left\lfloor G \cdot \frac{|\mathcal{S}_t|}{|\mathcal{S}_t| + |\mathcal{D}_t|} \right\rfloor, \qquad k_d = G - k_s. \quad (6)$$

This allocation adapts to content dynamics: allocating more budget to static compression when motion is minimal and prioritizing dynamic selection when changes are salient. With the budget allocated, CTR applies dual-path processing: for dynamic tokens representing new information, it selects the top-$k_d$ based on saliency scores; for static tokens containing redundant information, it consolidates them through density-based clustering to $k_s$ tokens, yielding

$$\mathcal{G}_d = \text{TopK}_{k_d}(\{\alpha_t^{(i)}\}_{i \in \mathcal{D}_t}), \quad \mathcal{G}_s = \text{Merge}_{k_s}(\{\mathcal{F}_t^{(i)}\}_{i \in \mathcal{S}_t}). \quad (7)$$

The per-frame retained group $\mathcal{G}_t$ combines both paths with $|\mathcal{G}_t| = G$. This fixed-size output is essential for streaming viability: unlike offline methods, where variable compression ratios are acceptable, streaming systems require predictable per-frame latency and stable batching to maintain real-time guarantees. CTR maintains only the previous-frame features as state, ensuring decisions are invariant to batch size and enabling consistent results across devices with different memory constraints. The state memory is $O(Nd)$ independent of stream length, preventing unbounded growth. Per-frame complexity is $O(N + G^2)$ with $G < N$, dominated by the similarity computation. This results in a pre-LLM bound that reduces prefill from $O(TNLd^2)$ to $O(TGLd^2)$.

### 3.4. Online Quantized Memory (OQM)

Although CTR alleviates computational bottlenecks through pre-LLM compression, memory in the LLM's kv-cache still grows linearly, creating a trade-off between preserving complete history for future queries and operating within bounded capacity. OQM addresses this by combining quantization with selective retrieval: CTR-produced groups are stored in 4-bit form to reduce storage overhead, while only the groups relevant to a given query are selectively dequantized and activated, enabling complete history preservation with bounded active memory. This coordination requires a unified operational unit. The group, the frame-level retained set of exactly $G$ tokens from CTR, serves as the atomic unit for both storage and retrieval, avoiding token-level indexing while preserving frame-aligned semantics.

OQM applies incremental, group-aligned quantization: each incoming group tensor $\mathbf{X}_t \in \mathbb{R}^{H \times G \times d}$ is independently compressed into a 4-bit representation along with a compact retrieval key, without revisiting historical groups,

$$\mathcal{Q}_4(\mathbf{X}_t) = \big(\text{uint4}(\hat{\mathbf{X}}_t), \boldsymbol{s}_t, \boldsymbol{m}_t, \bar{\mathbf{k}}_t\big), \quad (8)$$

where per-head, per-channel scale and offset are computed across tokens within each group, and a representative key $\bar{\mathbf{k}}_t \in \mathbb{R}^d$ is computed by averaging the keys before quantization to enable efficient retrieval,

$$\boldsymbol{s}_{t,h,d} = \frac{\max_{j \in [1,G]} \mathbf{X}_{t,h,j,d} - \min_{j \in [1,G]} \mathbf{X}_{t,h,j,d}}{15},$$
$$\boldsymbol{m}_{t,h,d} = \min_{j \in [1,G]} \mathbf{X}_{t,h,j,d} \quad (9)$$

and the 4-bit codes are packed into `uint8` for storage. This yields a dual-structure memory: quantized groups with their representative keys are appended over time, while system tokens, such as instructions and task descriptions, remain in full precision,

$$\mathcal{M}_t = \{\mathcal{KV}_{\text{sys}}^{\text{FP16}}\} \cup \{(\mathcal{Q}_4(\mathcal{G}_i), \bar{\mathbf{k}}_i)\}_{i=1}^t, \quad (10)$$

and the per-frame storage ratio from FP16 to 4-bit can be computed as

$$\frac{\text{FP16 storage}}{\text{4-bit storage}} \approx \frac{N}{G} \cdot \frac{16}{4} = \frac{4N}{G}, \quad (11)$$

where the baseline stores $N \cdot d \cdot 16$ bits and OQM stores $G \cdot d \cdot 4$ bits per frame.

Only selected groups are restored to full precision via the dequantization operator,

$$\mathcal{Q}_4^{-1}\big(\text{uint4}(\hat{\mathbf{X}}); \boldsymbol{s}, \boldsymbol{m}\big) = \text{depack}(\hat{\mathbf{X}}) \odot \boldsymbol{s} + \boldsymbol{m}, \quad (12)$$

where $\text{depack}(\cdot)$ unpacks 4-bit codes to integers in $\{0, \ldots, 15\}$ and $\odot$ denotes elementwise multiplication. To keep active memory bounded, OQM performs group-level retrieval using the pre-computed representative keys without dequantizing all groups. For a decoder state $\mathbf{q}$, OQM retrieves the stored representative keys $\bar{\mathbf{k}}_i$ for all groups and computes their similarity to $\mathbf{q}$ in full precision. It then selects the top-$k$ relevant groups based on cosine similarity,

$$\mathcal{R} = \text{TopK}\big\{\text{sim}(\mathbf{q}, \bar{\mathbf{k}}_i)\big\}_{i=1}^t, \quad (13)$$

then dequantizes only the selected groups to form the active kv for decoding,

$$\mathcal{KV}_{\text{active}} = \bigcup_{i \in \mathcal{R}} \mathcal{Q}_4^{-1}(\mathcal{G}_i). \quad (14)$$

This design enables two critical properties: total storage can scale as $O(T \cdot G \cdot d)/4$ to preserve complete history, while the active kv during decoding remains bounded by $O(k \cdot G \cdot d)$ with $k \ll T$, ensuring real-time decoding regardless of stream length. The query-time cost consists of retrieval over $T$ groups plus selective reconstruction $O(k \cdot G)$. This retrieval-then-dequantize paradigm complements CTR's fixed-budget design: CTR ensures predictable prefill by outputting exactly $G$ tokens per frame, while OQM ensures predictable decoding by retrieving at most $k$ groups, together maintaining bounded latency under streaming causality.

Table 1. Offline evaluation across three long-video benchmarks. Methods are grouped by design paradigm and training requirement, with "+" indicating LLaVA-OV-7B implementations. The **best** and second best results within the streaming-designed training-free category are in bold and underlined, respectively.

| Method | Frames | VideoMME | | | | MLVU | EgoSchema | Avg. |
|---|---|---|---|---|---|---|---|---|
| | | Long | Medium | Short | Overall | | | |
| *Offline-Designed (Training-free)* | | | | | | | | |
| LLaVA-OV-7B [14] | 32 | 48.8 | 56.4 | 70.1 | 58.4 | 64.7 | 60.1 | 61.0 |
| +*DyCoke* [29] | 32 | – | – | – | 54.3 | – | 59.5 | – |
| +*VisionZip* [41] | 32 | – | – | – | 58.2 | – | 60.3 | – |
| +*HoliTom* [24] | 32 | – | – | – | 58.9 | – | 61.2 | – |
| *Streaming-Designed (Training-based)* | | | | | | | | |
| MovieChat-7B [28] | 2048 | 33.4 | – | – | 38.2 | 25.8 | 53.5 | 39.2 |
| Dispider-7B [22] | 1 fps | 49.7 | 53.7 | – | 56.5 | 61.7 | 55.6 | 57.9 |
| *Streaming-Designed (Training-free)* | | | | | | | | |
| LLaVA-OV-7B [14] | 32 | 48.8 | 56.4 | 70.1 | 58.4 | 64.7 | 60.1 | 61.0 |
| +*ReKV* [7] | 0.5 fps | – | – | – | – | **68.5** | 60.7 | - |
| +*LiveVLM* [18] | 0.5/0.2 fps | 48.8 | 56.4 | 66.7 | 57.3 | 66.3 | 59.0 | 60.9 |
| +*StreamMem* [42] | 0.5/0.2 fps | 50.1 | 56.6 | **71.5** | 59.4 | 66.9 | 63.0 | 63.1 |
| +*StreamingTOM (ours)* | 0.5/0.2 fps | **50.6** | **57.8** | 71.3 | **59.9** | 67.9 | **63.7** | **63.8** |

## 4. Experiments

### 4.1. Experimental Setups

**Benchmarks.** We evaluate StreamingTOM on two tracks: (i) offline long-video understanding and (ii) online streaming qa. For offline evaluation, we adopt VideoMME [9], MLVU [46], and EgoSchema [17], consistent with streaming methods [18, 42]. Videos are sampled at 0.5 frames per second (fps) for clips under 30 minutes and 0.2 fps for longer sequences. We report VideoMME without subtitles and evaluate MLVU/EgoSchema on official dev splits. For online streaming, we utilize RVS-Ego and RVS-Movie [44], where questions arrive after their end timestamps and require causal reasoning based on preceding content. These benchmarks are specifically designed to evaluate real-time streaming capabilities under strict causal constraints.

**Baselines.** We compare against four categories of methods. As the offline reference, we use LLaVA-OV-7B [14], which leverages global temporal access but lacks streaming-specific memory control. For training-free streaming approaches, we evaluate four methods managing post-LLM kv-cache. ReKV [7] retrieves historical segments from external storage, StreamMem [42] maintains fixed-size memory through attention-based compression, while LiveVLM [18] and InfiniPot-V [13] restructure the cache via selective retention and compression. Additionally, we compare with training-based streaming systems that learn specialized architectures, including Flash-VStream [44], Dispider-7B [22], and MovieChat-7B [28]. As offline compression references, we also compare with training-free methods including DyCoke [29], VisionZip [41], and HoliTom [24], evaluated on fixed 32-frame clips with 25% before-LLM token retention [24].

Table 2. Evaluation results on RVS streaming benchmarks. Acc denotes accuracy (%), and Score represents response quality rated on a 1–5 scale. The **best** result for each metric is in bold.

| Method | RVS-Ego | | RVS-Movie | | Avg | |
|---|---|---|---|---|---|---|
| | Acc | Score | Acc | Score | Acc | Score |
| ReKV [7] | 63.7 | 4.0 | 54.4 | 3.6 | 59.0 | 3.8 |
| ReKV w/o off. [7] | 55.8 | 3.3 | 50.8 | 3.4 | 53.3 | 3.4 |
| Flash-VStream [44] | 57.0 | **4.0** | 53.1 | 3.3 | 55.0 | 3.6 |
| InfiniPot-V [13] | 57.9 | 3.5 | 51.4 | 3.5 | 54.6 | 3.5 |
| StreamMem [42] | 57.6 | 3.8 | 52.7 | 3.4 | 55.2 | 3.6 |
| **StreamingTOM (ours)** | **58.3** | 3.9 | **53.2** | **3.5** | **55.8** | **3.7** |

**Implementation Details.** We follow the ReKV streaming evaluation protocol [7] and only vary our module hyperparameters; unless noted, prompts and decoding settings follow the corresponding baselines. All experiments run on a single NVIDIA A6000 (48 GB) GPU using FP16 mixed precision. We use greedy decoding for deterministic outputs. For evaluation, long-video benchmarks utilize full conversation templates, while short-video benchmarks employ simplified templates to reduce verbosity.

Pre-LLM reduction with CTR retains 50 tokens per frame based on adjacent-frame changes and token saliency, using a strictly causal two-frame decision without look-ahead; we set the similarity threshold to 0.9 to determine frame-level redundancy. Post-LLM memory with OQM stores retained tokens in 4-bit quantized groups and retrieves relevant groups on demand with a total budget of 12k tokens. For fair comparison with ReKV, we align this budget to match ReKV baseline's 64 frames at 196 tokens per frame. Cached tokens are organized into groups of size 50 to stabilize memory. The streaming encoder batch size defaults to 32; our compression is batch-agnostic and supports flexible deployment.

### 4.2. Main Results

**Offline Video Evaluation.** In Table 1, StreamingTOM achieves the best overall performance among all streaming systems, regardless of whether they require training. Our training-free method reaches an average of **63.8**, surpassing the strongest training-free baseline StreamMem [42] at 63.1 and outperforming training-based Dispider-7B. StreamingTOM demonstrates consistently strong performance across individual benchmarks, achieving **59.9** on VideoMME-Overall, **67.9** on MLVU, and **63.7** on EgoSchema, consistently outperforming the training-free baseline StreamMem. Within VideoMME splits, StreamingTOM excels particularly on Long and Medium subsets while showing comparable performance on Short videos, demonstrating robust capability across different temporal scales. These results clearly show that under the same backbone and protocol, StreamingTOM sets a new SOTA for streaming video un-

Figure 4. Two qualitative examples from RVS Movie requiring long-horizon reasoning. StreamingTOM provides faithful answers consistent with ground truth, accurately capturing fine-grained semantics such as "subway station" and multi-party interactions. These results illustrate the model's ability to maintain causal reasoning and long-horizon consistency in real streaming scenarios.

derstanding despite being training-free.

Beyond streaming comparisons, StreamingTOM surpasses offline compression methods that use 32-frame clips on the same backbone, outperforming HoliTom [24], VisionZip [41], and DyCoke [29] on VideoMME-Overall and EgoSchema. This demonstrates that streaming-designed methods, by processing frames causally at 0.5/0.2 fps, adaptively cover longer temporal spans than offline methods constrained to fixed 32-frame clips.

**Online Video Evaluation.** We evaluate StreamingTOM on RVS-Ego and RVS-Movie, using GPT-3.5-turbo-0125 [20] for grading and limiting GPU memory to 28 GB without CPU offloading for fair comparison. Under this memory-constrained setting, StreamingTOM achieves the strongest overall performance with **55.8 Acc** and **3.7 Score**, surpassing the strongest training-free baseline Stream-Mem at 55.2/3.6 and outperforming other methods including InfiniPot-V, ReKV, and training-based Flash-VStream. While ReKV with CPU offloading reports higher accuracy (59.0), this strategy relies on CPU-GPU memory transfers that are impractical for real-time streaming deployment. We therefore compare under memory-constrained settings without CPU offloading, where StreamingTOM achieves the strongest performance. StreamingTOM achieves **58.3/3.9** on RVS-Ego and **53.2/3.5** on RVS-Movie, demonstrating consistent performance across both. We further provide qualitative examples on RVS Movie (Figure 4),
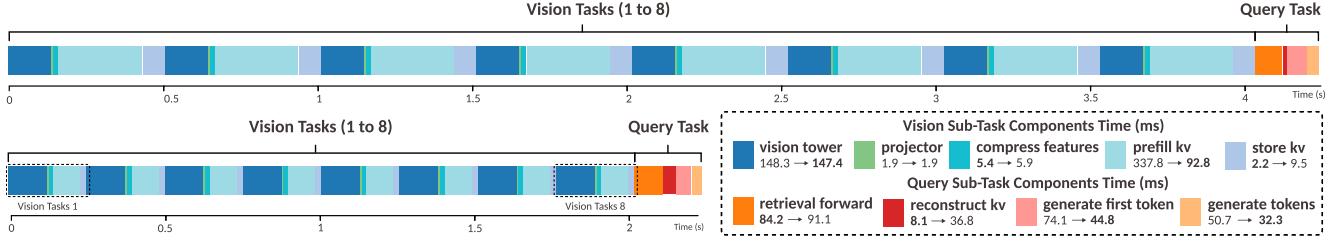
Table 3. StreamingTOM efficiency metrics across frame budgets and sampling rates under two batch sizes. Bounded memory growth and consistent throughput confirm real-time capability.

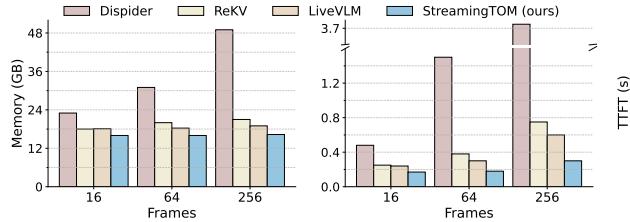| Metric | Frames | | | | Streaming | |
|---|---|---|---|---|---|---|
| | **16** | **64** | **256** | **512** | **0.2 fps** | **0.5 fps** |
| *Batch Size: 8* | | | | | | |
| GPU Mem. / GB ↓ | 16.0 | 16.0 | 16.3 | 16.7 | 16.9 | 18.4 |
| TTFT / s ↓ | 0.17 | 0.20 | 0.30 | 0.30 | 0.31 | 0.36 |
| Throughput / tokens·s$^{-1}$ ↑ | 36.7 | 32.6 | 20.8 | 20.9 | 20.8 | 20.9 |
| *Batch Size: 32* | | | | | | |
| GPU Mem. / GB ↓ | 16.7 | 18.6 | 18.9 | 19.3 | 19.5 | 21.1 |
| TTFT / s ↓ | 0.18 | 0.18 | 0.28 | 0.29 | 0.29 | 0.30 |
| Throughput / tokens·s$^{-1}$ ↑ | 36.5 | 32.4 | 20.9 | 20.9 | 21.0 | 20.9 |

where StreamingTOM produces faithful answers aligned with ground truth, demonstrating its ability to capture fine-grained semantics under causal streaming constraints.

## 4.3. Efficiency Analysis

We evaluate efficiency using three metrics (peak memory, TTFT, and generation throughput) across varying frame budgets, sampling rates, and batch sizes. Our compression is batch-agnostic, maintaining consistent accuracy across all batch sizes. Figure 6 compares StreamingTOM in detail with representative training-free baselines. Unlike LiveVLM, which compresses kv-cache post-LLM, our dual-stage design achieves superior efficiency: **CTR** reduces visual tokens *before* the LLM, cutting the $O(tNLd^2)$

Figure 5. **Complete processing timeline comparison.** The top row shows the baseline, while the bottom row shows StreamingTOM. For a 64-frame stream at batch size 8 with 50 tokens per frame, StreamingTOM incurs minimal overhead: OQM introduces 7.3ms for kv storage, 6.9ms for retrieval, and 28.7ms for 4-bit reconstruction. In contrast, CTR delivers 3.6× prefill acceleration from 337.8ms to 92.8ms, yielding an efficient 0.20s query TTFT and demonstrating substantial performance gains at negligible cost.



Figure 6. Memory and TTFT comparison across frame budgets. StreamingTOM outperforms all baselines in both metrics.

Table 4. Ablation on token retention and quantization. Ratio is relative to a 196-token, 16-bit baseline. *Results for 50 tokens with 2-bit quantization are interpolated from 48 and 52 tokens.

| Token | Quant. | Ratio (%) | VideoMME↑ | | | |
|-------|--------|-----------|-----------|--------|------|---------|
| | | | Short | Medium | Long | Overall |
| 40 | 4-bit | 5.1 | 70.1 | 56.8 | 49.9 | 58.9 |
| | 2-bit | 2.6 | 69.1 | 57.1 | 48.0 | 58.1 |
| 50 | 4-bit | 6.4 | **71.3** | **57.8** | **50.6** | **59.9** |
| | 2-bit* | 3.2 | 69.3 | 57.3 | 48.8 | 58.5 |
| 60 | 4-bit | 7.7 | 69.9 | 57.6 | 50.4 | 59.3 |
| | 2-bit | 3.8 | 68.9 | 56.8 | 50.3 | 58.7 |

prefill computation and kv-writes, while **OQM** further compresses the resulting kv-cache through 4-bit quantization and bounded retrieval. StreamingTOM exhibits significantly flatter memory growth due to a smaller and more stable active kv-cache, with advantages widening as the number of frames increases: at 256 frames, StreamingTOM achieves $1.2\times$ lower peak memory and $2\times$ faster TTFT compared to LiveVLM.

Table 3 provides a comprehensive performance profile of StreamingTOM across different configurations to guide practical deployment. Memory exhibits sublinear growth consistently due to CTR's fixed token budget, remaining nearly constant from 16.0 GB to 16.7 GB across 16-512 frames, demonstrating bounded memory under pre-LLM compression. TTFT remains consistently low and throughput stabilizes at real-time levels across varying video lengths and batch sizes, confirming scalability for practical streaming deployment. Figure 5 breaks down the complete processing pipeline, comparing baseline and StreamingTOM across vision and query stages. These results confirm that StreamingTOM maintains efficient and predictable performance across diverse streaming scenarios.

### 4.4. Ablation Study

Table 4 analyzes the trade-off between per-frame detail and temporal coverage. Our default configuration (50 tokens, 4-bit) achieves the best balance: fewer tokens (40) sacrifice critical details, while more tokens (60) reduce temporal cov-

erage under fixed memory budgets, both degrading accuracy. This configuration reduces memory footprint to 6.4% relative to the 196-token 16-bit baseline, achieving a $15.7\times$ compression ratio while maintaining strong accuracy.

## 5. Conclusion

This paper presents StreamingTOM, the first training-free framework that performs pre-LLM token reduction designed for streaming video understanding. Unlike existing methods that only manage post-LLM kv-cache, StreamingTOM addresses pre-LLM computation and post-LLM memory bottlenecks through a unified two-stage design. Causal Temporal Reduction (CTR) enforces causality with a fixed per-frame budget $G$, drastically reducing prefill complexity by processing only $G$ selected tokens instead of all $N$ visual tokens per frame, while Online Quantized Memory (OQM) bounds post-LLM memory through 4-bit quantization and retrieval. Experiments demonstrate that StreamingTOM achieves state-of-the-art accuracy among training-free methods, reaching 63.8% on offline benchmarks and 55.8%/3.7 on RVS with $15.7\times$ kv-cache compression. Critically, it sustains bounded memory growth over long horizons by reducing one-hour stream kv-cache from 18.8GB to 1.2GB, enabling practical deployment of video LLMs for real-time, long-duration applications.

# References

[1] Anthropic. Claude Sonnet 4.5. https://www.anthropic.com/news/claude-sonnet-4-5, 2024. 1

[2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 1

[3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1

[4] 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 *CVPR*, 2024. 2

[5] Joya Chen, Ziyun Zeng, Yiqi Lin, Wei Li, Zejun Ma, and Mike Zheng Shou. Livecc: Learning video llm with streaming speech transcription at scale. In *CVPR*, 2025. 2

[6] 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 *ECCV*, 2024. 1

[7] 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. In *ICLR*, 2025. 2, 6

[8] Xin Ding, Hao Wu, Yifan Yang, Shiqi Jiang, Donglin Bai, Zhibo Chen, and Ting Cao. Streammind: Unlocking full frame rate streaming video dialogue through event-gated cognition. In *ICCV*, 2025. 2

[9] 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 *CVPR*, 2025. 2, 6

[10] Shenghao Fu, Qize Yang, Yuan-Ming Li, Yi-Xing Peng, Kun-Yu Lin, Xihan Wei, Jian-Fang Hu, Xiaohua Xie, and Wei-Shi Zheng. Vispeak: Visual instruction feedback in streaming videos. *arXiv preprint arXiv:2503.12769*, 2025. 2

[11] Google DeepMind. Gemini 2.5 Pro. https://deepmind.google/models/gemini/pro/, 2024. 1

[12] Xiaohu Huang, Hao Zhou, and Kai Han. Prunevid: Visual token pruning for efficient video large language models. In *ACL*, 2025. 2

[13] 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. 2, 6

[14] 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. 1, 6

[15] Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen Luo, Peng Li, Yang Liu, and Maosong Sun. Streamingbench: Assessing the gap for mllms to achieve streaming video understanding. *arXiv preprint arXiv:2411.03628*, 2024. 1

[16] Jinming Liu, Junyan Lin, Yuntao Wei, Kele Shao, Keda Tao, Jianguo Huang, Xudong Yang, Zhibo Chen, Huan Wang, and Xin Jin. Revisiting mllm token technology through the lens of classical visual coding. *arXiv preprint arXiv:2508.13460*, 2025. 1

[17] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. In *NeurIPS*, 2023. 2, 6

[18] Zhenyu Ning, Guangda Liu, Qihao Jin, Wenchao Ding, Minyi Guo, and Jieru Zhao. Livevlm: Efficient online video understanding via streaming-oriented kv cache and retrieval. *arXiv preprint arXiv:2505.15269*, 2025. 2, 6

[19] OpenAI. GPT-5. https://openai.com/gpt-5/, 2024. 1

[20] OpenAI. GPT-3.5 Turbo. https://platform.openai.com/docs/models/gpt-3.5-turbo, 2024. 7

[21] Xiaohuan Pei, Tao Huang, and Chang Xu. Cross-self kv cache pruning for efficient vision-language inference. *arXiv preprint arXiv:2412.04652*, 2024. 2

[22] Rui Qian, Shuangrui Ding, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. Dispider: Enabling video llms with active real-time interaction via disentangled perception, decision, and reaction. In *CVPR*, 2025. 2, 6

[23] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. In *ICCV*, 2025. 2

[24] Kele Shao, Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan Wang. Holitom: Holistic token merging for fast video large language models. In *NeurIPS*, 2025. 1, 2, 6, 7

[25] Kele Shao, Keda Tao, Kejia Zhang, Sicheng Feng, Mu Cai, Yuzhang Shang, Haoxuan You, Can Qin, Yang Sui, and Huan Wang. When tokens talk too much: A survey of multimodal long-context token compression across images, videos, and audios. *arXiv preprint arXiv:2507.20198*, 2025. 1

[26] Leqi Shen, Guoqiang Gong, Tao He, Yifeng Zhang, Pengzhang Liu, Sicheng Zhao, and Guiguang Ding. Fastvid: Dynamic density pruning for fast video large language models. *arXiv preprint arXiv:2503.11187*, 2025. 2

[27] Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, et al. Longvu: Spatiotemporal adaptive compression for long video-language understanding. In *ICML*, 2025. 2

[28] 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. 6

[29] Keda Tao, Can Qin, Haoxuan You, Yang Sui, and Huan Wang. Dycoke: Dynamic compression of tokens for fast video large language models. In *CVPR*, 2025. 1, 2, 6, 7

[30] Keda Tao, Haoxuan You, Yang Sui, Can Qin, and Huan Wang. Plug-and-play 1.x-bit kv cache quantization for video large language models. *arXiv preprint arXiv:2503.16257*, 2025. 2

[31] Zhongwei Wan, Ziang Wu, Che Liu, Jinfa Huang, Zhihong Zhu, Peng Jin, Longyue Wang, and Li Yuan. Look-m: Look-once optimization in kv cache for efficient multimodal long-context inference. *arXiv preprint arXiv:2406.18139*, 2024. 2

[32] Haibo Wang, Bo Feng, Zhengfeng Lai, Mingze Xu, Shiyu Li, Weifeng Ge, Afshin Dehghan, Meng Cao, and Ping Huang. Streambridge: Turning your offline video large language model into a proactive streaming assistant. *arXiv preprint arXiv:2505.05467*, 2025. 2

[33] 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. 1

[34] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe Chen, Hongjie Zhang, Ganlin Yang, Haomin Wang, Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, Guanzhou Chen, Zichen Ding, Changyao Tian, Zhenyu Wu, Jingjing Xie, Zehao Li, Bowen Yang, Yuchen Duan, Xuehui Wang, Zhi Hou, Haoran Hao, Tianyi Zhang, Songze Li, Xiangyu Zhao, Haodong Duan, Nianchen Deng, Bin Fu, Yinan He, Yi Wang, Conghui He, Botian Shi, Junjun He, Yingtong Xiong, Han Lv, Lijun Wu, Wenqi Shao, Kaipeng Zhang, Huipeng Deng, Biqing Qi, Jiaye Ge, Qipeng Guo, Wenwei Zhang, Songyang Zhang, Maosong Cao, Junyao Lin, Kexian Tang, Jianfei Gao, Haian Huang, Yuzhe Gu, Chengqi Lyu, Huanze Tang, Rui Wang, Haijun Lv, Wanli Ouyang, Limin Wang, Min Dou, Xizhou Zhu, Tong Lu, Dahua Lin, Jifeng Dai, Weijie Su, Bowen Zhou, Kai Chen, Yu Qiao, Wenhai Wang, and Gen Luo. Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 1

[35] Meng Wei, Chenyang Wan, Xiqian Yu, Tai Wang, Yuqiang Yang, Xiaohan Mao, Chenming Zhu, Wenzhe Cai, Hanqing Wang, Yilun Chen, et al. Streamvln: Streaming vision-and-language navigation via slowfast context modeling. *arXiv preprint arXiv:2507.05240*, 2025. 1

[36] Haoyuan Wu, Xueyi Chen, Rui Ming, Jilong Gao, Shoubo Hu, Zhuolun He, and Bei Yu. Totrl: Unlock llm tree-of-thoughts reasoning potential through puzzles solving. *arXiv preprint arXiv:2505.12717*, 2025. 1

[37] Haoyuan Wu, Rui Ming, Jilong Gao, Hangyu Zhao, Xueyi Chen, Yikai Yang, Haisheng Zheng, Zhuolun He, and Bei Yu. On-policy optimization with group equivalent preference for multi-programming language understanding. *arXiv preprint arXiv:2505.12723*, 2025. 1

[38] Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui He, Jiaqi Wang, Feng Wu, and Dahua Lin. Pyramiddrop: Accelerating your large vision-language models via pyramid visual redundancy reduction. In *CVPR*, 2025. 2

[39] Haomiao Xiong, Zongxin Yang, Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Jiawen Zhu, and Huchuan Lu. Streaming video understanding and multi-round interaction with memory-enhanced knowledge. In *ICLR*, 2025. 2

[40] Ruyi Xu, Guangxuan Xiao, Yukang Chen, Liuning He, Kelly Peng, Yao Lu, and Song Han. Streamingvlm: Real-time understanding for infinite video streams. *arXiv preprint arXiv:2510.09608*, 2025. 2

[41] Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao Wang, Jingyao Li, Bei Yu, and Jiaya Jia. Visionzip: Longer is better but not necessary in vision language models. In *CVPR*, 2025. 1, 2, 6, 7

[42] Yanlai Yang, Zhuokai Zhao, Satya Narayan Shukla, Aashu Singh, Shlok Kumar Mishra, Lizhu Zhang, and Mengye Ren. Streammem: Query-agnostic kv cache memory for streaming video understanding. *arXiv preprint arXiv:2508.15717*, 2025. 2, 6

[43] Linli Yao, Yicheng Li, Yuancheng Wei, Lei Li, Shuhuai Ren, Yuanxin Liu, Kun Ouyang, Lean Wang, Shicheng Li, Sida Li, et al. Timechat-online: 80% visual tokens are naturally redundant in streaming videos. In *ACM MM*, 2025. 2

[44] Haoji Zhang, Yiqin Wang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. Flash-vstream: Memory-based real-time understanding for long video streams. In *ICCV*, 2025. 1, 2, 6

[45] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Llava-video: Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 1

[46] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang Liang, Shitao Xiao, Minghao Qin, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: Benchmarking multi-task long video understanding. In *CVPR*, 2025. 2, 6