# Flash-VStream: Efficient Real-Time Understanding for Long Video Streams

Haoji Zhang[1,2*]   Yiqin Wang[1,2*]   Yansong Tang[1,2✉]   Yong Liu[1,2]   Jiashi Feng[4]   Xiaojie Jin[3,4✉†]

[1]Tsinghua University   [2]Tsinghua Shenzhen International Graduate School   [3]Beijing Jiaotong University   [4]ByteDance Inc.

## Abstract

*Benefiting from the advances in large language models and cross-modal alignment, existing multimodal large language models have achieved prominent performance in image and short video understanding. However, the understanding of long videos is still challenging, as their long-context nature results in significant computational and memory overhead. Most existing work treats long videos in the same way as short videos, which is inefficient for real-world applications and hard to generalize to even longer videos. To address these issues, we propose Flash-VStream, an efficient video language model capable of processing extremely long videos and responding to user queries in real time. Particularly, we design a Flash Memory module, containing a low-capacity context memory to aggregate long-context temporal information and model the distribution of information density, and a high-capacity augmentation memory to retrieve detailed spatial information based on this distribution. Compared to existing models, Flash-VStream achieves significant reductions in inference latency. Extensive experiments on long video benchmarks and comprehensive video benchmarks, i.e., EgoSchema, MLVU, LVBench, MVBench and Video-MME, demonstrate the state-of-the-art performance and outstanding efficiency of our method. Code is available at* `https://github.com/IVGSZ/Flash-VStream`.

## 1. Introduction

Achieving robust perception and understanding of complex, dynamic environments is a crucial milestone toward Artificial General Intelligence (AGI). A core objective in this pursuit is the development of advanced large multimodal models [1, 53, 57] capable of effectively integrating diverse data types, including text, visual content, and audio.

Among multimodal tasks, long video understanding stands out as particularly significant yet challenging, primarily due to its substantial computational overhead and high GPU memory demands. Consequently, improving model efficiency is essential for making long video understanding practical. High GPU memory usage limits real-world de-
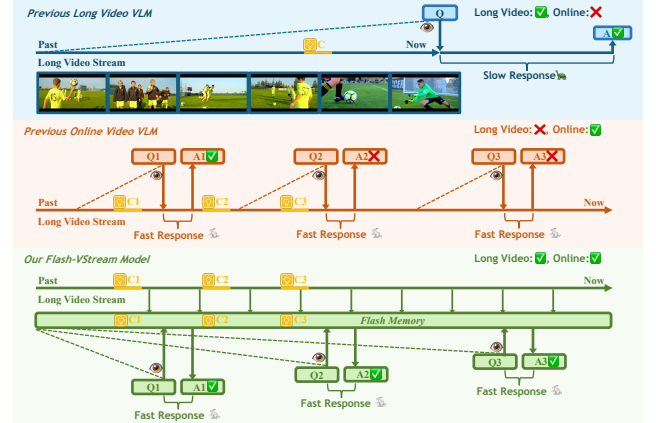


Figure 1. **Comparison with previous methods.** Flash-VStream can understand long videos accurately in an online manner. Here "C" denotes critical clues for questions.

ployment, particularly on resource-constrained edge devices. Additionally, excessive computational requirements increase inference latency, directly affecting applications requiring real-time human-computer interactions. In this context, we define a video language model (VLM) as *real-time* if it can respond to user queries within one second.

Real-time interaction capabilities are vital in numerous practical scenarios. Multimodal assistants [1, 53, 57, 64], for example, must operate in real-time to ensure fluid user experiences. Similarly, robots deployed in real-world environments benefit from VLMs capable of interpreting and reacting swiftly to dynamic situations [48, 51]. Surveillance systems also rely on real-time VLMs to continuously analyze video streams, thereby maintaining security effectively [6, 17, 44]. However, current state-of-the-art VLMs still struggle to achieve real-time responsiveness when performing question-answering tasks on long videos [8, 27, 40, 57].

To address this challenge, we start from a widely recognized observation: *temporal redundancy* is prevalent in all video types [20, 21, 60–62]. Since not all video frames are equally informative, an efficient model should allocate more computational resources preferentially to critical frames. In Sec. 3, we propose the Flash Memory module, which addresses these challenges by integrating a Context Synopsis Memory that captures information density distribution along

---

*Equal contribution.     ✉Corresponding authors.     †Project leader.

the temporal dimension, and a Detail Augmentation Memory designed to retrieve detailed content by selectively sampling key frames.

In this paper, we introduce Flash-VStream, a two-process video-language model capable of efficiently handling extremely long video streams and providing real-time responses to user queries. As illustrated in Fig. 3, the frame handler operates continuously, encoding new frames through a visual encoder and updating the Flash Memory without interruption. Concurrently, the question handler functions as a server process, triggered by incoming questions, and leverages Flash Memory to rapidly generate the first token of the answer within one second. In summary, the two-process framework ensures simultaneous video processing, memory updating and real-time response generation. The innovative memory design sets our model apart from previous works, in terms of video length and online capability ( Fig. 1).

As demonstrated in Fig. 2, Flash-VStream significantly reduces inference latency to meet real-time standards, achieving a superior balance between accuracy and efficiency, and setting new state-of-the-art performance on the full EgoSchema benchmark [43].

We further validate the generalization capability of Flash-VStream through zero-shot video question answering experiments on three long video understanding benchmarks (EgoSchema [43], MLVU [82], LVBench [58]) and two comprehensive video understanding benchmarks (MVBench [30], Video-MME [16]), as shown in Tab. 2. Additionally, detailed ablation studies in Sec. 4.4 clearly verify the effectiveness of Flash Memory. We summarize our contributions as follows:

- We introduce Flash-VStream, an efficient large video language model capable of processing extremely long video and providing real-time responses to user queries. Flash-VStream utilizes a fixed-size Flash Memory to bridge a two-process framework, ensuring efficient and timely processing of long-term video streams.
- A novel Flash Memory module is proposed to reduce temporal redundancy between consecutive frames, which includes a Context Synopsis Memory to model the distribution of information density along the temporal dimension and a Detail Augmentation Memory to retrieve detailed information from key frames.
- Extensive experiments and ablation studies on long video understanding benchmarks and comprehensive video understanding benchmarks demonstrate the outstanding performance and efficiency of Flash-VStream.

## 2. Related Work

### 2.1. Video Language Models

With recent advances in Large Language Models (LLMs) [4, 14, 45, 55, 56, 73], and Multimodal Large Language Models
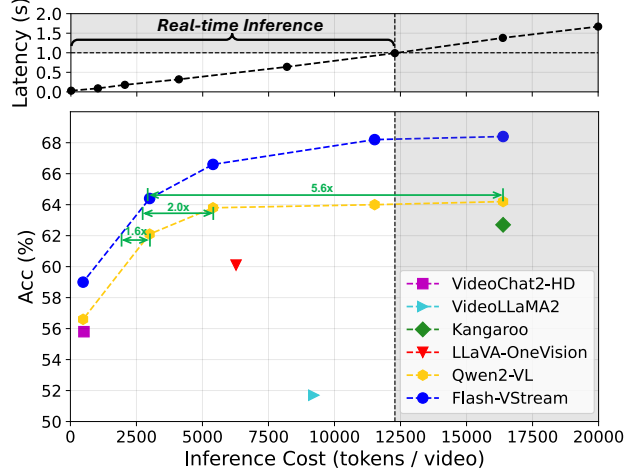


Figure 2. **Response latency / Accuracy on EgoSchema v.s. Inference cost.** Inputting more than 12000 tokens will not meet real-time requirements. Flash-VStream can respond to user queries in real time while maintaining outstanding performance.

(MLLMs) [11, 27–29, 32, 33, 54], many works begin to build Video Language Models (VLMs) based on them. LLaMA-VID [31] represents single-frame features with only 2 tokens. Chat-UniVi [25] employs dynamic tokens to model video features of different scales. Other works [40, 63, 75, 76] use different compression techniques to represent an entire video with fewer tokens. While these methods succeeded in short video understanding, they have relatively poor performance on long video understanding benchmarks [16, 43].

### 2.2. Long Video Understanding

Earlier work MIST [18] introduces an iterative method to select the most question-related video clip. SEVILA [77] performs temporal key frame localization and video question answering simultaneously. Compared to question-aware methods [18, 31, 77], Flash-VStream does not rely on specific questions and can achieve general information aggregation based on visual information itself. MovieChat [49] proposes to merge similar frame features by average pooling. Though it is able to process long video with limited GPU memory cost, its performance is suboptimal due to its training-free framework. Recent works [71, 80] explore long context extension finetuning. However, they are computationally expensive. RETAKE [59] proposes a KV cache pruning method, which successfully lowers the knowledge redundancy in long videos. Different from them, Flash-VStream keeps most informative frames in memory to reduce temporal redundancy, resulting in higher efficiency.

### 2.3. Real-Time Video Stream Understanding

Real-time video stream understanding requires models to process video streams and finish specific tasks in real-time. Most existing methods are designed to perform a specific
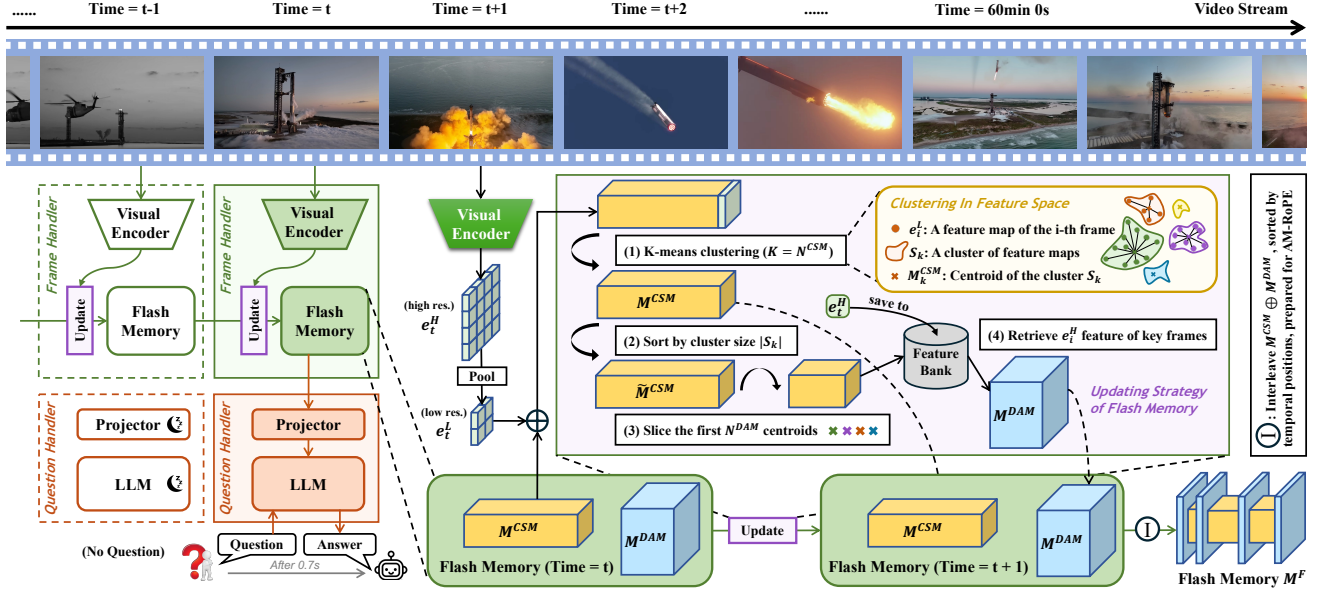
Figure 3. **Overview of Flash-VStream two-process framework.** The frame handler process continuously encodes new frames. The question handler process asynchronously responds to human inquiries in real-time. Flash Memory is composed of interleaved Context Synopsis Memory and Detail Augmentation Memory, organized in chronological order. CSM is updated by clustering low resolution feature maps on an inter-frame level. DAM is updated by retrieving high resolution feature maps of the most informative frames from a feature bank.

vision task, such as real-time object tracking [22, 68], action recognition [39, 65, 79], and segmentation [38, 66, 74]. Zhou et al. [83] design a streaming framework for dense video captioning. VideoLLM-Online [7] is designed for streaming video narration and action anticipation. Considering natural language is becoming a general interface for various modalities [19, 29, 67] and tasks [2, 36, 84], our work focuses on real-time video question answering (VQA) upon free-form user queries, which is more challenging.

## 2.4. Memory for Long Sequence Modeling

Memory mechanisms are widely used to store and retrieve information in all forms of long sequence modeling tasks, such as time series forecasting [5], recommendation system [52], and video object segmentation [9, 35]. For video understanding, MovieChat [24] utilizes a long-term memory and short-term memory framework. MC-ViT [3] proposes a memory-consolidated vision transformer for long video understanding. GLMGIR [72] proposes a multi-granularity memory to solve fine-grained video captioning. In comparison, our method uses two memory modules that focus on temporal information modeling and spatial detail enhancement, resulting in a synergic improvement of efficiency.

## 3. Flash-VStream

Flash-VStream improves model efficiency by allocating more computation to the most informative and representative key frames. As illustrated in Fig. 3, Flash Memory

is updated iteratively to retain key information from both current and historical frames. The Flash Memory comprises a Context Synopsis Memory (CSM) for long-term temporal information aggregation and representative key frame localization. It also contains a Detail Augmentation Memory (DAM) for enhancing the spatial details of the key frames.

### 3.1. Model Architecture

In order to lower the inference latency, Flash-VStream decouples vision processing and language processing into two processes. As presented in Fig. 3, the frame handler process is responsible for continuously frame encoding and memory consolidation, while the question handler process remains online, waiting for user queries. These processes collaborate by reading from and writing to shared memory, namely the Flash Memory. Following common practices, we adopt a Vision Transformer [13] as the visual encoder, a 2-layer MLP [33] as the projector, and a Qwen2-7b LLM [73] as the language decoder. Formally, the visual encoder takes the $t$-th frame $V_t \in \mathbb{R}^{H \times W \times 3}$ as input and outputs a high resolution feature map $e_t^{\mathrm{H}} \in \mathbb{R}^{h \times w \times d}$ and a low resolution feature map $e_t^{\mathrm{L}} \in \mathbb{R}^{h' \times w' \times d}$:

$$e_t^{\mathrm{H}} = f_{\mathrm{enc}}(V_t) \tag{1}$$

$$e_t^{\mathrm{L}} = f_{\mathrm{enc}}(\mathrm{pool}(V_t)) \tag{2}$$

Here $h \times w, h' \times w'$ denotes the number of patches, $d$ represents the hidden dimension size of visual encoder and "pool" stands for an average pooling layer with ratio $R_{\mathrm{pool}} = 4$.

3

| Configuration | Notation | CSM | DAM |
|---|---|---|---|
| Input Frames | - | 120 | 60 |
| Input Resolution | $H \times W$ | $224 \times 224$ | $448 \times 448$ |
| Temporal Size | $N$ | 60 | 30 |
| Spatial Size | $h \times w$ | 256 | 1024 |
| LLM tokens | $N_{\text{Vtokens}}$ | $60 \times 64$ | $30 \times 256$ |

Table 1. **Flash Memory configurations.** The shape of an input frame can be rectangular, as long as pixel amount is less than Input Resolution. Temporal Size is the number of feature maps of memory. Spatial Size is the number of ViT tokens in a feature map.

## 3.2. Context Synopsis Memory

CSM is designed for efficient long-term understanding. To represent as many frames as possible within limited resources, CSM ($M^{\text{CSM}}$) uses a set of compressed low resolution feature maps with size $S^{\text{CSM}} = N^{\text{CSM}} \times h' \times w' \times d$. CSM integrates semantically similar frames, which form a cluster of similar frames, i.e., the *context*. As defined in Eq. (3), each item in CSM is the centroid of a cluster of low resolution feature maps, i.e., the *synopsis* of the context.

$$M^{\text{CSM}} = \left\{ \frac{1}{|S_k|} \sum_{i \in S_k} e_i^{\text{L}} \right\}_{k=1}^{N^{\text{CSM}}}, 1 \leq i \leq t \qquad (3)$$

$$M^{\text{CSM}} = \text{cluster}(M^{\text{CSM}} \oplus e_{t+1}^{\text{L}}) \qquad (4)$$

Here $S_k$ is the k-th cluster set with size=$|S_k|$, and $\oplus$ denotes feature concatenation. CSM is initialized with the feature maps of the first $N^{\text{CSM}}$ frames. In Eq. (4), CSM is updated by clustering algorithm using a new feature map of the $(t+1)$-th frame, following [3, 25, 49]. Since we need to limit the number of clusters to $N^{\text{CSM}}$, K-means clustering [42] is employed as an effective and efficient clustering algorithm.

Therefore, CSM maintains the cluster centroids to reduce additional computation, and the clusters themselves serve as an implicit representation of information density.

## 3.3. Detail Augmentation Memory

While CSM effectively aggregates long-term temporal information, it compromises spatial details critical for fine-grained video understanding. To address this, we propose DAM as a complementary module to retain and augment spatial details crucial for precise video understanding. Considering the prohibitive computational costs of high resolution feature maps, DAM operates by selectively storing high resolution feature maps of key frames.

It is natural to borrow the results of CSM for free and use them to guide key frame selection. Given that a CSM cluster centroid $M_k^{\text{CSM}}$ is a synopsis of multiple frames, DAM retrieves fine-grained spatial features based on these centroids. DAM adopts a *Feature-Centric* retrieval policy for key frame localization. Specifically, DAM is a set of high resolution
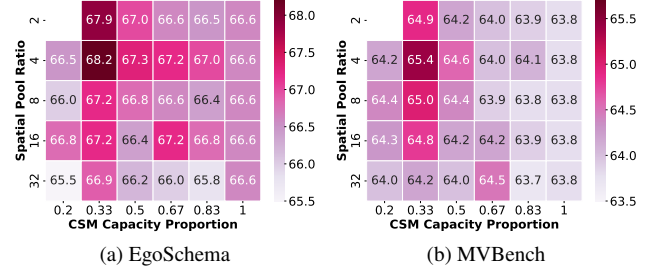


(a) EgoSchema    (b) MVBench

Figure 4. **Impact of Pool Ratio $R_{\text{pool}}$ and CSM Capacity Proportion $R_{\text{CSM}}$.** The upper left grid is blank since its setting is invalid.

feature maps $M^{\text{DAM}}$ with size $S^{\text{DAM}} = N^{\text{DAM}} \times h \times w \times d$.

$$\tilde{M}^{\text{CSM}} = \text{sort}(M^{\text{CSM}}, \text{key} = \{|S_k|\}_{k=1}^{N^{\text{CSM}}}) \qquad (5)$$

$$M^{\text{DAM}} = \{e_{f(k)}^{\text{H}}\}_{k=1}^{N^{\text{DAM}}}, f(k) = \underset{i}{\arg\min} \, D(\tilde{M}_k^{\text{CSM}}, e_i^{\text{L}}) \qquad (6)$$

As shown in Fig. 3, DAM takes the top-$N^{\text{DAM}}$ largest cluster centroids as anchors. A frame is considered a key frame if its $e_i^{\text{L}}$ is the nearest to a centroid anchor in the feature space. In Eqs. (5) and (6), $\tilde{M}_k^{\text{CSM}}$ represents the centroid of the $k$-th largest cluster, sorted by the cluster size $|S_k|$, i.e., the number of frames in the cluster. $f(k)$ denotes the index of the $k$-th important key frame. $D(\cdot, \cdot)$ denotes the Euclidean distance function. More ablation studies are provided in Sec. 4.4.

A Feature Bank $E_t^{\text{H}} = \{e_1^{\text{H}}, e_2^{\text{H}}, ..., e_t^{\text{H}}\}$ is maintained for the retrieval of high resolution feature maps, where $t$ is the number of current frames. The feature bank can be offloaded to disk to avoid memory overflow.

In short, CSM captures long-term temporal information and DAM augments it with more fine-grained spatial details of key frames. They complement each other to provide a comprehensive understanding of long videos. As shown in Fig. 3, Flash Memory $M^F$ is the interleaved form of CSM and DAM, sorted by temporal positions of feature maps. Formally, temporal positions can be calculated as:

$$P_k^{\text{CSM}} = \frac{1}{|S_k|} \sum_{i \in S_k} i, 1 \leq k \leq N^{\text{CSM}} \qquad (7)$$

$$P_k^{\text{DAM}} = f(k), 1 \leq k \leq N^{\text{DAM}} \qquad (8)$$

$$M^F = \text{sort}(M^{\text{CSM}} \oplus M^{\text{DAM}}, \text{key} = P^{\text{CSM}} \oplus P^{\text{DAM}}) \qquad (9)$$

## 3.4. Adaptive Multimodal RoPE

Rotary Position Embedding (RoPE) [50] is a widely used position embedding method in LLMs [14, 55, 56]. We improve upon the original Multimodal RoPE (M-RoPE) [50, 57] to support flexible positions, resulting in the Adaptive Multimodal RoPE (AM-RoPE). M-RoPE first splits the hidden dimension to three groups representing time, height and width axes, so that the position index becomes a triplet $(n_t, n_h, n_w)$. AM-RoPE is designed for flexible relative

| Model | Max $N_{\text{Vtokens}}$ | EgoSchema | MLVU$_{\text{dev}}$ | LVBench | MVBench | Video-MME w/o subs | Video-MME w/ subs |
|---|---|---|---|---|---|---|---|
| *Offline video language models* | | | | | | | |
| MovieChat [49] | 32 | - | - | 22.5 | - | - | - |
| TimeChat [47] | 96 | 33.0 | 30.9 | 22.3 | - | - | - |
| LLaMA-VID [31] | 2fps | 38.5 | 33.2 | 23.9 | 41.9 | - | - |
| ChatUni-Vi [25] | 896 | - | - | - | - | 40.6 | 45.9 |
| ShareGPT4-video [8] | 9216 | - | 46.4 | - | 51.2 | 39.9 | 43.6 |
| Video-Chat2-HD [30] | 512 | 55.8 | 47.9 | - | 62.3 | 45.3 | 55.7 |
| VideoLLaMA2 [10] | 9216 | 51.7 | 48.5 | - | 54.6 | 47.9 | 50.3 |
| LongVA [80] | 18432 | - | 56.3 | - | - | 52.6 | 54.3 |
| LLaVA-OneVision [27] | 6272 | 60.1 | 64.7 | - | 56.7 | 58.2 | 61.5 |
| LongVILA [71] | 50176 | - | - | - | - | 57.5 | 61.8 |
| Kangaroo [34] | 16384 | 62.7 | 61.0 | - | 61.0 | 56.0 | 57.6 |
| Oryx-MLLM [37] | 14400 | - | - | 30.4 | 63.9 | 58.3 | 62.6 |
| Qwen2-VL* [57] | 24576 | 64.8 | 66.0 | 41.4 | 65.1 | 61.1 | 65.9 |
| *Online video language models* | | | | | | | |
| VideoLLM-Online [7] | 2fps | 32.8 | 35.2 | 24.0 | 33.9 | 26.9 | 29.9 |
| VideoLLM-MOD [69] | 1fps | - | - | - | - | 49.2 | - |
| VideoStreaming [46] | 256 | 44.1 | - | - | - | - | - |
| Qwen2-VL-online† | 11520 | 64.0 | 62.9 | 39.8 | 63.3 | 59.4 | 65.1 |
| Flash-VStream (Ours) | 11520 | **68.2** (↑ 4.2) | **66.3** (↑ 3.4) | **42.0** (↑ 2.2) | **65.4** (↑ 2.1) | **61.2** (↑ 1.8) | **67.0** (↑ 1.9) |

Table 2. **Comparison with state-of-the-art video language models on video question answering benchmarks.** We conduct experiments on five mainstream multiple-choice benchmarks. $N_{\text{Vtokens}}$ denotes the number of video tokens used during evaluation. *subs* is short for subtitles. Qwen2-VL* denotes the reproduced results under the official setting. Qwen2-VL-online† denotes Qwen2-VL tested under real-time restriction ($N_{\text{Vtokens}} <= 11520$). The best two results are **bold-faced** and underlined, respectively.

position embedding. The key is to use average positions to represent the compressed cluster feature. In Flash-VStream, there are two types of video token, the CSM token and DAM token. For a DAM token at position $(x, y)$ of $M_k^{\text{DAM}}$, the AM-RoPE is calculated as: $n_t = P_k^{\text{DAM}} = t(k), n_h = y \times 2, n_w = x \times 2$ to accommodate the effect of average pooling. For a CSM token at position $(x, y)$ of $M_k^{\text{CSM}}$, $n_t = P_k^{\text{CSM}}, n_h = y, n_w = x$.

When a new question is posed, the question handler process begins to infer based on current Flash Memory $M^F$ containing $N_{\text{Vtokens}} = N^{\text{CSM}} \times h' \times w' + N^{\text{DAM}} \times h \times w$ tokens. By adjusting this budget, it is possible to ensure efficient real-time response in our asynchronous framework. We conduct a speed test in Sec. 4.2 and find that a 7b model can achieve real-time inference by limiting $N_{\text{Vtokens}} \leq 12000$, as presented in Fig. 2.

## 4. Experiments

### 4.1. Experimental setup

**Implementation Details.** Following Qwen2-VL, we employ a ViT with 3D patch embedding layer as the visual encoder, a merger projector as the projector. Therefore, every two adjacent frames are temporally pooled before being encoded, and four adjacent ViT tokens are spatially pooled to one LLM token. The visual encoder, projector and LLM are initialized from a pretrained MLLM, Qwen2-VL-7b [57]. Tab. 1 shows detailed configurations of the proposed Flash

Memory and input frames. As discussed in Sec. 4.4, we empirically assign 1/3 tokens to CSM and 2/3 tokens to DAM. The Flash Memory costs 11520 LLM tokens in total.

**Training Settings.** To enhance video understanding ability based on Flash Memory, we adopt a LoRA [23] instruction tuning stage. With the parameters of the visual encoder frozen, all linear layers of projector and LLM are LoRA finetuned. We adopted a 9k subset from LLaVA-Video dataset [81], which contains instruction-following tasks like captioning, open-ended VQA and multiple-choice VQA. More training details are provided in the supplementary.

**Evaluation Settings.** For more reliable analysis, we conduct zero-shot multiple-choice VQA experiments on three long video benchmarks and two comprehensive video benchmarks. We report the multiple-choice accuracy on each benchmark. EgoSchema [43] is a long-form VQA dataset that is specifically designed for understanding first-person behaviors. MLVU [82] is a multi-task long video benchmark covering various video genres. LVBench [58] is an extreme long video understanding benchmark designed to test long-term comprehension capabilities of models. MVBench [30] contains amounts of temporal-related tasks, aiming at testing temporal understanding ability. Video-MME [16] is a high-quality comprehensive video analysis benchmark, which contains videos ranging from 11 seconds to 1 hour.

Flash-VStream is evaluated on these benchmarks in an online setting. First, frames are extracted at 1 fps. Each

5

| Model | Max $N_{\text{Vtokens}}$ | w/o subtitles | | | | w/ subtitles | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Short | Middle | Long | Average | Short | Middle | Long | Average |
| ChatUni-Vi [25] | 896 | 45.7 | 40.3 | 35.8 | 40.6 | 51.2 | 44.6 | 41.8 | 45.9 |
| ShareGPT4-video [8] | 9216 | 48.3 | 36.3 | 35.0 | 39.9 | 53.6 | 39.3 | 37.9 | 43.6 |
| VideoLLaMA2 [10] | 9216 | 56.0 | 45.4 | 42.0 | 47.9 | 59.4 | 47.6 | 43.8 | 50.3 |
| LongVA [80] | 18432 | 61.1 | 50.4 | 46.2 | 52.6 | 61.1 | 53.6 | 47.6 | 54.3 |
| LongVILA [71] | 147456 | 69.3 | 56.1 | 47.0 | 57.5 | 70.8 | 60.6 | 54.0 | 61.8 |
| Kangaroo [34] | 16384 | 66.1 | 55.3 | 46.7 | 56.0 | 68.0 | 55.4 | 49.3 | 57.6 |
| Qwen2-VL* [57] | 24576 | <u>71.3</u> | **61.3** | **50.7** | <u>61.1</u> | 71.1 | **69.0** | 57.6 | <u>65.9</u> |
| Qwen2-VL-online† | 11520 | 70.3 | 59.8 | 48.2 | 59.4 | <u>71.4</u> | 66.0 | <u>57.9</u> | 65.1 |
| Flash-VStream (Ours) | 11520 | **72.0** | <u>61.1</u> | <u>50.3</u> | **61.2** | **72.4** | <u>67.0</u> | **61.4** | **67.0** |

Table 3. **Comparison with state-of-the-art video language models on Video-MME benchmark.** $N_{\text{Vtokens}}$ denotes the number of video tokens used during evaluation. Qwen2-VL* denotes the reproduced results under the official setting. Qwen2-VL-online† denotes Qwen2-VL tested under real-time restriction ($N_{\text{Vtokens}}$ <= 11520). The best two results are **bold-faced** and <u>underlined</u>, respectively.

frame is then fed to the Flash-VStream model sequentially, with the question posed at the end of the video frame stream.

## 4.2. Computational Efficiency

We first measure the response latency of the Flash-VStream model by counting the response wall time of the question handler process under different cost limitations, as presented in Fig. 2. Cost limitation confines the maximum number of video tokens during inference, namely $N_{\text{Vtokens}}$. All experiments are tested on a single A100 GPU with BFloat16 precision and FlashAttention-2 [12]. The experimental results show that our model can achieve real-time response with $N_{\text{Vtokens}}$ <= 12000. Models like Kangaroo [34], Qwen2-VL [57], LongVA [80] and LongVILA [71] require far more than 12k tokens for inference, resulting in latency higher than 1 second. In contrast, Flash-VStream leverages a multi-process framework with an efficient memory design and outperforms all competitive methods at the same token cost.

## 4.3. Main Results

In comparison to previous competitive models [7, 8, 10, 25, 30, 31, 34, 37, 46, 47, 49, 57, 71, 80] , Flash-VStream shows superior understanding capabilities on various challenging video understanding benchmarks, as detailed in Tab. 2. For Qwen2-VL [57], we re-evaluate it under official setting and find it difficult to reproduce the reported results. Although Flash-VStream performs slightly behind Qwen2-VL on MLVU dataset, our method is more efficient and saves 53% of video tokens during inference. Flash-VStream significantly promotes the development of online VLMs, which are able to respond to human instructions in real time.

For a fair comparison, we slightly modify the testing setting of Qwen2-VL [57] by limiting the maximum number of video tokens to 11520, resulting in a "Qwen2-VL-online" baseline. As reported in Tab. 2, Flash-VStream significantly surpasses Qwen2-VL-online on both short video and long video benchmarks under equal-cost condition, which demonstrates the effectiveness of the proposed Flash Memory.

Tab. 3 provides a comprehensive comparison of model performance on the Video-MME benchmark. Leveraging the CSM and DAM memory, Flash-VStream achieves an effective trade-off between accuracy and efficiency across various video types.

## 4.4. Ablation Study

**Flash Memory.** We conduct an ablation study to evaluate the effectiveness of CSM and DAM in the Flash Memory. From the results in Tab. 4, one can observe that both CSM and DAM bring significant and highly generalizable improvements across three benchmarks. There are three groups of experiments in Tab. 4. The first group evaluates the impact of removing components from the baseline setting row ①. CSM and DAM help improve the average accuracy by 2.0% and 0.7% compared to uniform sampling, respectively. The second group further evaluates the influence of CSM with the existence of uniformly sampled DAM in row ⑤ (since full DAM relies on full CSM). The third group compares the implementation details of Qwen2-VL-online setting in Tab. 2, which show different preferences for each benchmark.

**Memory Capacity Allocation.** We investigate the capacity allocation strategy for CSM and DAM with a fixed total memory tokens budget. For **real-time inference**, the amount of visual tokens should not exceed 12000 (Fig. 2). Under this fixed budget condition, we adjust the proportion of CSM capacity in total memory $R_{\text{CSM}} = S^{\text{CSM}}/(S^{\text{CSM}} + S^{\text{DAM}})$ by controlling $N^{\text{CSM}}$ and $N^{\text{DAM}}$; and adjust the pool ratio $R_{\text{pool}} = (h \times w)/(h' \times w')$ by controlling $h \times w$ and $N^{\text{DAM}}$. The grid search result in Fig. 4 shows that allocating around one-third of the total memory capacity to CSM when pool ratio = 4 yields the best performance for EgoSchema and MVBench. This suggests that a balanced allocation strategy, where CSM is given sufficient capacity to capture long-term temporal information while DAM retains enough capacity to preserve detailed spatial information, is crucial for optimal performance. More ablation studies on memory structure configuration can be found in the supplementary.

| ID | Memory Component Settings | | | | | Evaluation Results | | | |
|---|---|---|---|---|---|---|---|---|---|
| | CSM | DAM | $N_{Vtokens}$ | CSM Size | DAM Size | MVBench | EgoSchema | Video-MME(w/o) | Average |
| ① | ✓ | ✓ | 11520 | $60 \times 64$ | $30 \times 256$ | 65.4 | 68.2 | 61.2 | 64.9 |
| ② | ✓ | Uni. Smp. | 11520 | $60 \times 64$ | $30 \times 256$ | 64.3 | 67.8 | 60.6 | 64.2 |
| ③ | ✓ | ✗ | 3840 | $60 \times 64$ | 0 | 64.0 | 66.8 | 60.1 | 63.6 |
| ④ | Uni. Smp. | ✗ | 3840 | $60 \times 64$ | 0 | 62.4 | 63.4 | 59.0 | 61.6 |
| ⑤ | ✓ | Uni. Smp. | 11520 | $60 \times 64$ | $30 \times 256$ | 64.3 | 67.8 | 60.6 | 64.2 |
| ⑥ | Uni. Smp. | Uni. Smp. | 11520 | $60 \times 64$ | $30 \times 256$ | 63.8 | 66.0 | 59.6 | 63.1 |
| ⑦ | ✗ | Uni. Smp. | 7680 | 0 | $30 \times 256$ | 63.1 | 65.7 | 59.3 | 62.7 |
| ⑧ | Uni. Smp. | ✗ | 11520 | $180 \times 64$ | 0 | 63.3 | 64.0 | 59.4 | 62.3 |
| ⑨ | ✗ | Uni. Smp. | 11520 | 0 | $45 \times 256$ | 63.2 | 65.1 | 59.0 | 62.4 |

Table 4. **Analysis of the design of Flash Memory.** We investigate the effects of eliminating the two components of Flash Memory: CSM and DAM. "Uni. Smp." stands for "Uniform Sample". The first row is the default baseline setting of our model. The penultimate row is the same as the Qwen2-VL-online setting in Tab. 2.

| ID | Clustering Policy | MVB | EGO | MME |
|---|---|---|---|---|
| ③ | K-means [42] | 64.0 | 66.8 | 60.1 |
| | DBScan [15] | 63.8 | 66.6 | 59.7 |
| | GMM | 59.6 | 65.8 | 59.5 |
| | Neighbor Merge [49] | 63.7 | 65.0 | 59.4 |
| | Neighbor Drop | 63.8 | 62.4 | 59.4 |
| ④ | Uniform Sample | 62.4 | 63.4 | 59.0 |

Table 5. **Ablation study on CSM clustering policy.** This experiment compares K-means with other clustering methods. DAM is removed to isolate the effect of CSM.

| ID | Retrieval Policy | Selection Policy | MVB | EGO | MME |
|---|---|---|---|---|---|
| ① | Feature-Centric | Top-k largest | 65.4 | 68.2 | 61.2 |
| | Temporal-Centric | Top-k largest | 64.4 | 67.8 | 60.8 |
| | Cosine Similarity | Top-k largest | 64.1 | 66.5 | 60.2 |
| ② | Uniform Sample | k frames | 64.3 | 67.8 | 60.6 |
| ① | Feature-Centric | Top-k largest | 65.4 | 68.2 | 61.2 |
| | Feature-Centric | Top-k smallest | 63.9 | 67.6 | 59.6 |
| | Uniform Sample | All frames | 65.5 | 68.6 | 61.2 |

Table 6. **Ablation study on DAM sampling policy.** This experiment compares different retrieval methods and selection methods. Here $k = N^{DAM}$ and the last row uses 19200 visual tokens.

**CSM Clustering Policy.** In Tab. 5, we compare the K-means clustering algorithm with other clustering methods. We remove DAM in Flash Memory to isolate the effect of CSM and take row ③ in Tab. 4 as a baseline. For example, DBScan [15] is robust to noise and outliers. GMM (Gaussian Mixture Model) is a clustering method based on probability distributions, which is appropriate for statistical analysis. MovieChat [49] utilizes a "Neighbor Merge" method for memory updating. It merges the tokens of the two most similar adjacent frames at every step. To examine the effect of token merge, we implement a similar method named "Neighbor Drop", which randomly chooses a frame to keep and another to drop, instead of merging them. All the methods above are compared to uniform sampling.

As presented in Tab. 5, K-means outperforms other clustering methods in the EgoSchema and MVBench datasets, surpassing the uniform sampling baseline by 3.4% and 1.6%, respectively. DBScan, while robust to noise and outliers, slightly underperforms K-means. We attribute this to the fact that this is a cluster number known task, where DBScan may not be the best choice. Note that although Neighbor Merge achieves equally good accuracy on MVBench, it falls behind in long video understanding tasks in EgoSchema.

**DAM Sampling Policy.** DAM is designed for storing details of the most informative key frames. As proposed in Sec. 3.3, the default Feature-Centric retrieval method relates the importance of a frame to the distance between its feature map and the cluster centroid feature map. It selects $N^{DAM}$ frames nearest to the centroids of the top-$N^{DAM}$ largest clusters as DAM. We compare it with other methods from Tab. 4 row ①

baseline. Temporal-Centric retrieval selects frames nearest to the temporal position $P_i^{CSM}$ of top-$N^{DAM}$ largest clusters. Cosine Similarity method replaces the Euclidean distance in Eq. (6) with cosine similarity. As presented in Tab. 6, Feature-Centric retrieval demonstrates robust improvement on short video and long video understanding benchmarks. Cosine Similarity results in relatively lower performance, indicating the magnitude of features matters.

As shown in the lower part of Tab. 6, top-k largest selection method outperforms top-k smallest method and uniform-k method, being comparable to sampling all frames.

### 4.5. Memory Visualization and Case Study

We investigate the memory consolidation procedure in the deep feature space by dimension reduction with Principal Component Analysis (PCA). As illustrated on the left side of Fig. 5, each red hexagon stands for a memory token of CSM or DAM, and each blue point stands for a frame feature map ($e_t^L$ in comparison with CSM and $e_t^H$ in comparison with DAM). The PCA plots show that the memory features are distinctly separated from the frame features, suggesting that the CSM and DAM mechanisms are successful in encoding and consolidating important information. We note that CSM and DAM formulate different shapes of clusters for each video. For example, in the upper part of Fig. 5, DAM has four densely distributed areas while CSM has only two. We suppose that this is because high resolution features have more discriminability. The visualization proves that the Flash Memory features effectively capture the underlying
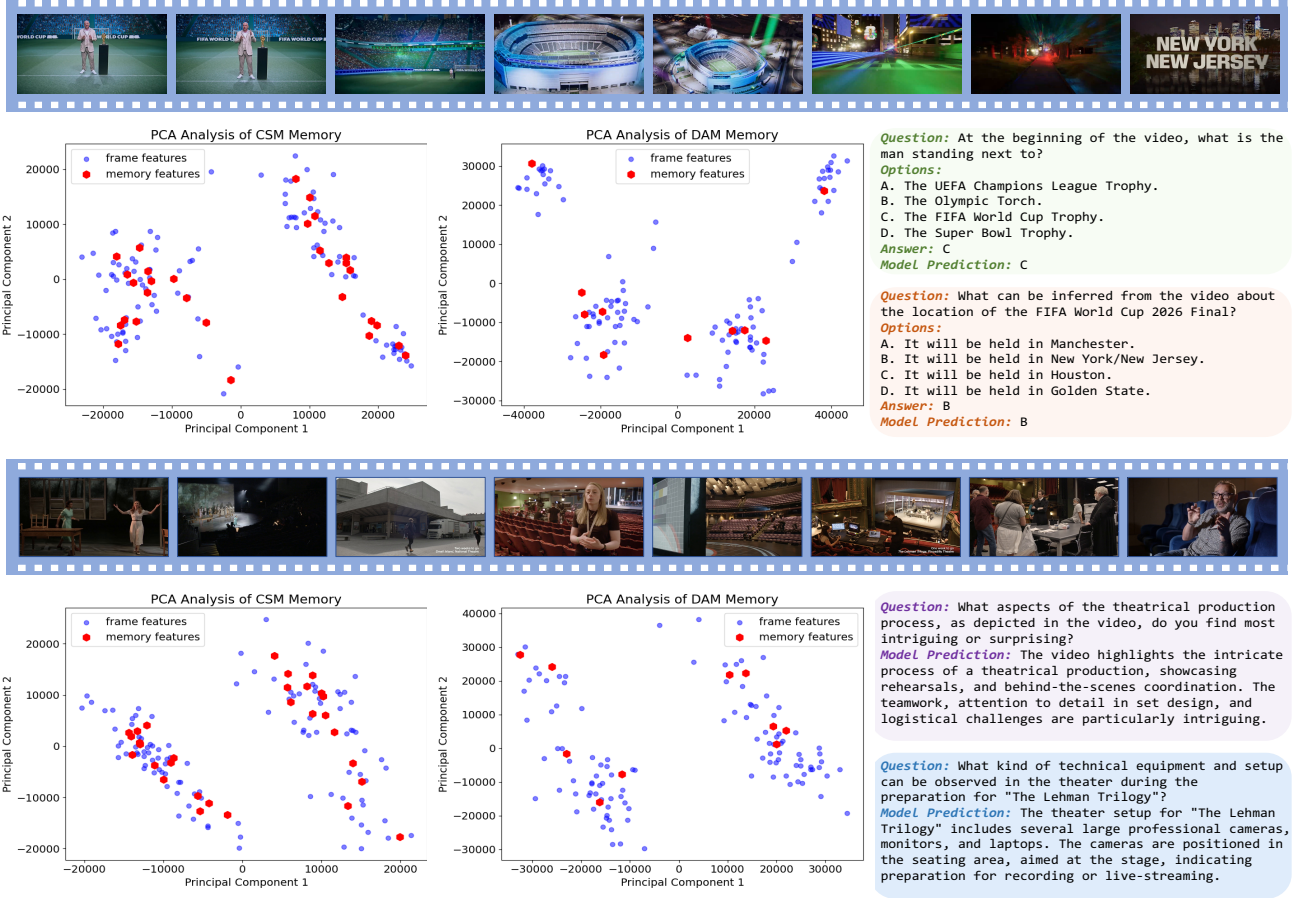
Figure 5. **Memory Distribution Visualization and Case Study.** The left side presents a PCA visualization of the Flash Memory distribution in the feature space. Each point in it stands for a feature map of a single frame or a slice of memory. The CSM and DAM appropriately represent the distributional characteristics of the feature clusters. On the right side, different types of question answering cases show exceptional proficiency of the Flash-VStream model. (upper: multiple-choice questions; lower: open-ended questions)

structure of the feature space.

On the right side of Fig. 5, we present several question-answering cases that illustrate the efficacy of the Flash-VStream model. The case study encompasses both multiple-choice questions and open-ended questions, including those related to spatial cognition, optical character recognition (OCR), and complex reasoning. These results demonstrate that Flash-VStream not only excels in understanding and processing long video content but also in providing precise and contextually relevant answers, thereby validating its practical applicability in real-world scenarios. For more case studies, please refer to the supplementary materials.

## 5. Conclusion

In conclusion, we introduce Flash-VStream, an efficient video-language model for understanding long video streams and providing real-time responses to user queries. Flash-VStream uses a two-process asynchronous framework to separate vision and language processing, ensuring real-time responses. The core innovation lies in the Flash Memory module, which comprises a Context Synopsis Memory for long-term temporal information aggregation and a Detail Augmentation Memory for retrieving detailed spatial information. This design is based on the observation that temporal redundancy is prevalent in videos. Extensive experiments on multiple comprehensive video benchmarks demonstrate the superior performance and efficiency of Flash-VStream compared to existing state-of-the-art models. We hope our work will inspire further research and advancements in the field of efficient long video understanding.

## Acknowledgements

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1

[2] Sule Bai, Yong Liu, Yifei Han, Haoji Zhang, and Yansong Tang. Self-calibrated clip for training-free open-vocabulary segmentation. *arXiv preprint arXiv:2411.15869*, 2024. 3

[3] Ivana Balazevic, Yuge Shi, Pinelopi Papalampidi, Rahma Chaabouni, Skanda Koppula, and Olivier J Hénaff. Memory consolidation enables long-context video understanding. In *ICML*, pages 2527–2542, 2024. 3, 4

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, pages 33, 1877–1901, 2020. 2

[5] Yen-Yu Chang, Fan-Yun Sun, Yueh-Hua Wu, and Shou-De Lin. A memory-network based solution for multivariate time-series forecasting. *arXiv preprint arXiv:1809.02105*, 2018. 3

[6] Jianguo Chen, Kenli Li, Qingying Deng, Keqin Li, and S Yu Philip. Distributed deep learning model for intelligent video surveillance systems with edge computing. *IEEE Transactions on Industrial Informatics*, 2019. 1

[7] 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*, pages 18407–18418, 2024. 3, 5, 6

[8] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with better captions. *NeurIPS*, 37:19472–19495, 2025. 1, 5, 6

[9] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *ECCV*, pages 640–658. Springer, 2022. 3

[10] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 5, 6

[11] Wenliang Dai, Junnan Li, DONGXU LI, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *NeurIPS*, 36, 2024. 2

[12] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *ICLR*, 2024. 6, 13

[13] 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 *ICLR*, 2021. 3

[14] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 2, 4

[15] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, pages 226–231, 1996. 7

[16] 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*, pages 24108–24118, 2025. 2, 5

[17] Tian Gan, Xiao Wang, Yan Sun, Jianlong Wu, Qingpei Guo, and Liqiang Nie. Temporal sentence grounding in streaming videos. In *ACM MM*, pages 4637–4646, 2023. 1

[18] Difei Gao, Luowei Zhou, Lei Ji, Linchao Zhu, Yi Yang, and Mike Zheng Shou. Mist: Multi-modal iterative spatial-temporal transformer for long-form video question answering. In *CVPR*, pages 14773–14783, 2023. 2

[19] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *IJCV*, 132(2):581–595, 2024. 3

[20] Amir Ghodrati, Babak Ehteshami Bejnordi, and Amirhossein Habibian. Frameexit: Conditional early exiting for efficient video recognition. In *CVPR*, pages 15608–15618, 2021. 1

[21] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *TPAMI*, 44(11):7436–7456, 2021. 1

[22] Anfeng He, Chong Luo, Xinmei Tian, and Wenjun Zeng. A twofold siamese network for real-time object tracking. In *CVPR*, pages 4834–4843, 2018. 3

[23] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022. 5

[24] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and Dahua Lin. Movienet: A holistic dataset for movie understanding. In *ECCV*, pages 709–727, 2020. 3

[25] Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual representation empowers large language models with image and video understanding. In *CVPR*, pages 13700–13710, 2024. 2, 4, 5, 6, 13, 14

[26] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *CVPR*, pages 13299–13308, 2024. 14

[27] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1, 2, 5, 13, 14

[28] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900, 2022.

9

[29] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742, 2023. 2, 3

[30] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *CVPR*, pages 22195–22206, 2024. 2, 5, 6

[31] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In *ECCV*, pages 323–340. Springer, 2025. 2, 5, 6, 13, 14

[32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36:34892–34916, 2023. 2

[33] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *CVPR*, pages 26296–26306, 2024. 2, 3

[34] Jiajun Liu, Yibing Wang, Hanghang Ma, Xiaoping Wu, Xiaoqi Ma, Xiaoming Wei, Jianbin Jiao, Enhua Wu, and Jie Hu. Kangaroo: A powerful video-language model supporting long-context video input. *arXiv preprint arXiv:2408.15542*, 2024. 5, 6

[35] Yong Liu, Ran Yu, Fei Yin, Xinyuan Zhao, Wei Zhao, Weihao Xia, and Yujiu Yang. Learning quality-aware dynamic memory for video object segmentation. In *ECCV*, pages 468–486, 2022. 3

[36] Yong Liu, Cairong Zhang, Yitong Wang, Jiahao Wang, Yujiu Yang, and Yansong Tang. Universal segmentation at arbitrary granularity with language instruction. In *CVPR*, pages 3459–3469, 2024. 3

[37] Zuyan Liu, Yuhao Dong, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. Oryx mllm: On-demand spatial-temporal understanding at arbitrary resolution. *arXiv preprint arXiv:2409.12961*, 2024. 5, 6

[38] Zhuoyan Luo, Yicheng Xiao, Yong Liu, Shuyan Li, Yitong Wang, Yansong Tang, Xiu Li, and Yujiu Yang. Soc: Semantic-assisted object cluster for referring video object segmentation. *NeurIPS*, pages 26425–26437, 2023. 3

[39] Diogo C Luvizon, David Picard, and Hedi Tabia. Multi-task deep learning for real-time 3d human pose estimation and action recognition. *IEEE TPAMI*, 43(8):2752–2764, 2020. 3

[40] Fan Ma, Xiaojie Jin, Heng Wang, Yuchen Xian, Jiashi Feng, and Yi Yang. Vista-llama: Reducing hallucination in video language models via equal distance to visual tokens. In *CVPR*, pages 13151–13160, 2024. 1, 2

[41] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *ACL*, pages 12585–12602, 2024. 14

[42] J MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability/University of California Press*, 1967. 4, 7

[43] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *NeurIPS*, 36:46212–46244, 2023. 2, 5

[44] Khan Muhammad, Tanveer Hussain, Javier Del Ser, Vasile Palade, and Victor Hugo C De Albuquerque. Deepres: A deep learning-based video summarization strategy for resource-constrained industrial surveillance scenarios. *IEEE Transactions on Industrial Informatics*, 16(9):5938–5947, 2019. 1

[45] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *NeurIPS*, pages 27730–27744, 2022. 2

[46] Rui Qian, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Shuangrui Ding, Dahua Lin, and Jiaqi Wang. Streaming long video understanding with large language models. *NeurIPS*, 37: 119336–119360, 2025. 5, 6

[47] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multimodal large language model for long video understanding. In *CVPR*, pages 14313–14323, 2024. 5, 6

[48] Pierre Sermanet, Tianli Ding, Jeffrey Zhao, Fei Xia, Debidatta Dwibedi, Keerthana Gopalakrishnan, Christine Chan, Gabriel Dulac-Arnold, Sharath Maddineni, Nikhil J Joshi, et al. Robovqa: Multimodal long-horizon reasoning for robotics. In *ICRA*, pages 645–652. IEEE, 2024. 1

[49] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *CVPR*, pages 18221–18232, 2024. 2, 4, 5, 6, 7, 13, 14

[50] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 4

[51] James Supancic III and Deva Ramanan. Tracking as online decision-making: Learning a policy from streaming videos with reinforcement learning. In *ICCV*, pages 322–331, 2017. 1

[52] Qiaoyu Tan, Jianwei Zhang, Ninghao Liu, Xiao Huang, Hongxia Yang, Jingren Zhou, and Xia Hu. Dynamic memory based attention network for sequential recommendation. In *AAAI*, pages 4384–4392, 2021. 3

[53] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1

[54] Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*, 2025. 2

[55] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2, 4

[56] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya

Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2, 4

[57] 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, 4, 5, 6, 13, 14

[58] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu Huang, Bin Xu, Yuxiao Dong, et al. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024. 2, 5

[59] 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. 2

[60] Yulin Wang, Zhaoxi Chen, Haojun Jiang, Shiji Song, Yizeng Han, and Gao Huang. Adaptive focus for efficient video recognition. In *ICCV*, pages 16249–16258, 2021. 1

[61] Yulin Wang, Yang Yue, Yuanze Lin, Haojun Jiang, Zihang Lai, Victor Kulikov, Nikita Orlov, Humphrey Shi, and Gao Huang. Adafocus v2: End-to-end training of spatial dynamic networks for video recognition. In *CVPR*, pages 20030–20040. IEEE, 2022.

[62] Yulin Wang, Yang Yue, Xinhong Xu, Ali Hassani, Victor Kulikov, Nikita Orlov, Shiji Song, Humphrey Shi, and Gao Huang. Adafocusv3: On unified spatial-temporal dynamic video recognition. In *ECCV*, pages 226–243. Springer, 2022. 1

[63] Yiqin Wang, Haoji Zhang, Yansong Tang, Yong Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin. Hierarchical memory for long video qa. *arXiv preprint arXiv:2407.00603*, 2024. 2

[64] Yiqin Wang, Haoji Zhang, Jingqi Tian, and Yansong Tang. Ponder & press: Advancing visual gui agent towards general computer control. *arXiv preprint arXiv:2412.01268*, 2024. 1

[65] Yulin Wang, Haoji Zhang, Yang Yue, Shiji Song, Chao Deng, Junlan Feng, and Gao Huang. Uni-adafocus: Spatial-temporal dynamic computation for video recognition. *TPAMI*, 2024. 3

[66] Yuji Wang, Jingchen Ni, Yong Liu, Chun Yuan, and Yansong Tang. Iterprime: Zero-shot referring image segmentation with iterative grad-cam refinement and primary word emphasis. In *AAAI*, pages 8159–8168, 2025. 3

[67] Yuji Wang, Haoran Xu, Yong Liu, Jiaze Li, and Yansong Tang. Sam2-love: Segment anything model 2 in language-aided audio-visual scenes. In *CVPR*, pages 28932–28941, 2025. 3

[68] Zhongdao Wang, Liang Zheng, Yixuan Liu, Yali Li, and Shengjin Wang. Towards real-time multi-object tracking. In *ECCV*, pages 107–122, 2020. 3

[69] Shiwei Wu, Joya Chen, Kevin Qinghong Lin, Qimeng Wang, Yan Gao, Qianli Xu, Tong Xu, Yao Hu, Enhong Chen, and Mike Zheng Shou. Videollm-mod: Efficient video-language streaming with mixture-of-depths vision computation. *NeurIPS*, 37:109922–109947, 2024. 5

[70] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *CVPR*, pages 9777–9786, 2021. 14

[71] Fuzhao Xue, Yukang Chen, Dacheng Li, Qinghao Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Haotian Tang, Shang Yang, Zhijian Liu, et al. Longvila: Scaling long-context visual language models for long videos. *arXiv preprint arXiv:2408.10188*, 2024. 2, 5, 6

[72] Yichao Yan, Ning Zhuang, Bingbing Ni, Jian Zhang, Minghao Xu, Qiang Zhang, Zheng Zhang, Shuo Cheng, Qi Tian, Yi Xu, et al. Fine-grained video captioning via graph-based multi-granularity interaction learning. *TPAMI*, 44(2):666–683, 2019. 3

[73] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 2, 3

[74] Zhao Yang, Jiaqi Wang, Xubing Ye, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Language-aware vision transformer for referring segmentation. *TPAMI*, 2024. 3

[75] Xubing Ye, Yukang Gan, Yixiao Ge, Xiao-Ping Zhang, and Yansong Tang. Atp-llava: Adaptive token pruning for large vision language models. In *CVPR*, pages 24972–24982, 2025. 2

[76] Xubing Ye, Yukang Gan, Xiaoke Huang, Yixiao Ge, and Yansong Tang. Voco-llama: Towards vision compression with large language models. In *CVPR*, pages 29836–29846, 2025. 2

[77] Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36:76749–76771, 2023. 2

[78] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *AAAI*, pages 9127–9134, 2019. 14

[79] Bowen Zhang, Limin Wang, Zhe Wang, Yu Qiao, and Hanli Wang. Real-time action recognition with enhanced motion vector cnns. In *CVPR*, pages 2718–2726, 2016. 3

[80] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*, 2024. 2, 5, 6

[81] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 5, 12

[82] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*, 2024. 2, 5

[83] Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Nagrani, and Cordelia Schmid. Streaming dense video captioning. In *CVPR*, pages 18243–18252, 2024. 3

[84] Yixuan Zhu, Haolin Wang, Ao Li, Wenliang Zhao, Yansong Tang, Jingxuan Niu, Lei Chen, Jie Zhou, and Jiwen Lu. Instarevive: One-step image enhancement via dynamic score matching. *arXiv preprint arXiv:2504.15513*, 2025. 3

# Flash-VStream: Efficient Real-Time Understanding for Long Video Streams

## Supplementary Material

In the supplementary material, we first provide implementation details of the Flash Memory mechanism and training settings. Subsequently, we conduct an analysis experiment on model inference efficiency and more ablation studies on memory structure configurations. We then present more visual cases to provide a comprehensive understanding of the performance of models.

## A. Implementation Details

This section describes the details of the proposed Flash Memory mechanism in Sec. 3. The Flash Memory consists of Context Synopsis Memory (CSM) and Detail Augmentation Memory (DAM). CSM uses a clustering-based updating policy, while DAM uses a retrieval-based updating policy.

$$M_k^{\text{CSM}} = \frac{1}{|S_k|} \sum_{i \in S_k} e_i^{\text{L}}, 1 \le k \le N^{\text{CSM}} \qquad (10)$$

$$M^{\text{CSM}} = \text{cluster}(M^{\text{CSM}} \oplus e_{t+1}^{\text{L}}) \qquad (11)$$

### A.1. Context Synopsis Memory

As mentioned in Sec. 3.2, CSM is designed for aggregating long-context temporal information and modeling the distribution of information density. $M_k^{\text{CSM}}$ represents the centroid of the k-th cluster. $M^{\text{CSM}}$ is initialized with the first $N^{\text{CSM}}$ feature maps of the first $N^{\text{CSM}}$ frames. When the next frame arrives, a clustering algorithm is employed to consolidate its feature map into existing clusters. Here we illustrate the "cluster" operation of Eq. (11) in detail.

As shown in Alg. 1, CSM performs a temporal-wise *K-means Clustering* algorithm to condense $(N^{\text{CSM}} + 1) \times h' \times w'$ tokens to $N^{\text{CSM}} \times h' \times w'$ tokens. Each frame feature in temporal memory $M_k^{\text{CSM}} = c_k \in \mathbb{R}^{h' \times w' \times d}$ represents the centroid of the i-th feature map cluster.

### A.2. Detail Augmentation Memory

As described in Sec. 3.3, DAM aims at storing spatial details of the most informative key frames, based on the feature clusters of CSM. For DAM, we use a Feature-Centric Sampling method to calculate $M^{\text{DAM}} \in \mathbb{R}^{N^{\text{DAM}} \times h \times w \times d}$.

Alg. 2 shows the pseudo code of *Feature-Centric Retrieval*. Here $w_j$ is equal to the size of $j$-th cluster, i.e., the number of feature maps in this cluster. We choose the centroids of the top-k largest clusters as anchors. Then we select key features from the feature bank $E_t^H$. $E_t^H$ keeps high-resolution feature maps of all frames on disk, where $t$ is the current number of frames. The features nearest to these anchors in the feature map space are considered as key features, which are added to the DAM.

---

**Algorithm 1** K-means Clustering Algorithm

**Require:** Current cluster centroids $M = M^{\text{CSM}}$
**Require:** Newest frame feature $e = e_t^{\text{L}}$
**Require:** Set of all points $X = \{M_1, M_2, \ldots, M_N, e\}$
**Require:** Weights vector of points $W = \{w_1, w_2, \ldots, w_N, 1\}$
**Require:** Maximum memory length $N = N^{\text{CSM}}$
**Require:** Maximum number of iterations $T$
1: **procedure** K-MEANS($X, W, N, T$)
2:     Initialize $t \leftarrow 0$
3:     Initialize centroids $C = \{c_1, c_2, \ldots, c_N\}$ from $X$
4:     Initialize cluster assignment $S_j \leftarrow \{\}, 1 \le j \le N$
5:     **while** $t < T$ **do**
6:         **for** $x_i \in X$ **do**
7:             $j \leftarrow \underset{j}{\operatorname{argmin}} \|x_i - c_j\|^2$
8:             $S_j \leftarrow S_j \cup \{x_i\}$
9:         **end for**
10:         **for** $j = 1, 2, \ldots, N$ **do**
11:             $c_j^{\text{new}} \leftarrow \dfrac{\sum_{x_i \in S_j} w_i \cdot x_i}{\sum_{x_i \in S_j} w_i}$
12:         **end for**
13:         Clear $S$
14:         $C \leftarrow C^{\text{new}}$
15:         $t \leftarrow t + 1$
16:     **end while**
17:     **for** $j = 1, 2, \ldots, N$ **do**
18:         $w_j^{\text{CSM}} \leftarrow \sum_{x_i \in S_j} w_i$
19:     **end for**
20:     $M^{\text{CSM}} = C$
21:     $W^{\text{CSM}} = \{w_1^{\text{CSM}}, w_2^{\text{CSM}}, \ldots, w_N^{\text{CSM}}\}$
22:     **return** $M^{\text{CSM}}, W^{\text{CSM}}$
23: **end procedure**

---

**Algorithm 2** Feature-Centric Sampling

**Require:** Current feature bank $E_t^{\text{H}} = \{e_1^{\text{H}}, e_2^{\text{H}}, \ldots, e_t^{\text{H}}\}$
**Require:** Current cluster centroids $M^{\text{CSM}}$
**Require:** Weights vector of points $W = \{w_1, w_2, \ldots, w_N\}$
**Require:** Maximum memory length $N = N^{\text{DAM}}$
1: **procedure** KEY FEATURE RETRIEVAL($E^{\text{H}}, M^{\text{CSM}}, W, N$)
2:     $k \leftarrow N$
3:     $idx \leftarrow \text{argsort}(W, \text{descending=True})$
4:     $j_1, j_2, \ldots, j_k \leftarrow idx[:k]$
5:     $M^{\text{DAM}} \leftarrow \{\}$
6:     **for** $z = 1, 2, \ldots, k$ **do**
7:         $anchor \leftarrow M_{j_z}^{\text{CSM}}$
8:         $i \leftarrow \underset{i \le t}{\operatorname{argmin}} \|e_i^{\text{L}} - anchor\|^2$
9:         $M^{\text{DAM}} \leftarrow M^{\text{DAM}} \cup \{e_i^{\text{H}}\}$
10:     **end for**
11:     **return** $M^{\text{DAM}}$
12: **end procedure**

---

## B. Training Details

We train Flash-VStream on a 9k subset of LLaVA-Video [81] dataset for one epoch. During training, we freeze the parameters of visual encoder, while all linear layers of projector and
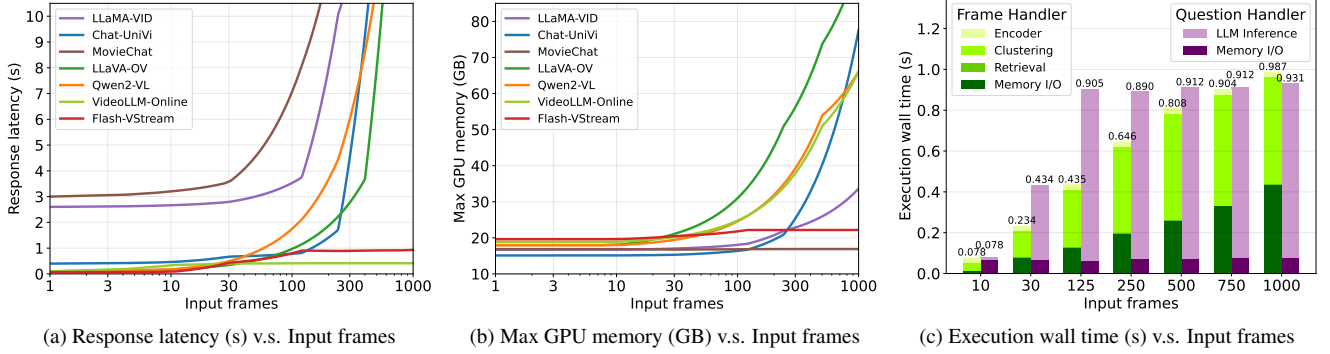
(a) Response latency (s) v.s. Input frames   (b) Max GPU memory (GB) v.s. Input frames   (c) Execution wall time (s) v.s. Input frames

Figure 6. **(a) Response latency comparison. (b) Max GPU memory comparison. (c) Execution wall time analysis.** Response latency refers to the wall time between inputting a question and outputting the first token of the answer. Max GPU memory indicates the peak GPU memory usage during inference. All experiments were conducted on A100 GPUs using BFloat16 and FlashAttention-2.

| Settings | Value |
|---|---|
| Batch Size | 64 |
| Learning Rate | 8e-4 |
| Lora Rank | 64 |
| Lora Alpha | 32 |
| Learning Schedule | Cosine decay |
| Warmup Ratio | 0.01 |
| Weight Decay | 0.1 |
| Epoch | 1 |
| Optimizer | AdamW |
| Deepspeed Stage | 2 |
| Visual Encoder | Freeze |
| Projector | Open |
| LLM | Open |

Table 7. Training settings of Flash-VStream.

LLM are LoRA finetuned. The overall training can be finished in about 10 hours on 8 A100 80G GPUs with BFloat16 automatic mixed precision and FlashAttention-2 [12]. Detailed training settings are shown in Tab. 7.

## C. Efficiency Analysis

An efficiency analysis is performed to assess the inference efficiency of Flash-VStream. Specifically, we concentrate on the response latency and GPU memory consumption of models, as discussed in Sec. 1 of the paper.

We compare Flash-VStream with other competitive video language models [25, 27, 31, 49, 57] in terms of response latency and max GPU memory. As presented in Fig. 6, Flash-VStream demonstrates superior performance in both efficiency metrics. Fig. 6a shows the response latency comparison, where Flash-VStream consistently exhibits lower latency across varying numbers of input frames. This indicates that Flash-VStream is more efficient in processing video inputs, resulting in faster response times (less than 1 second). Fig. 6b illustrates the maximum GPU memory usage. Flash-VStream maintains a relatively stable and lower GPU memory consumption compared to other models, even as the number of input frames increases. This efficiency in memory usage makes Flash-VStream more scalable and suitable for deployment in resource-constrained environments.

From a systematic perspective, we measure the execution wall time of each process in Fig. 6c. The result shows that the question handler process stays fast enough (< 1s) regardless of the number of input frames. This is because the question handler only relies on size-fixed Flash Memory. The execution time of the frame handler process grows up to more than 1 second when the number of frames exceeds 1000. Although this may result in delayed updates of visual information, it would not affect the response latency.

Overall, the results highlight the efficiency advantages of Flash-VStream in terms of both response latency and GPU memory consumption, making it a competitive choice for real-time long video understanding tasks.

## D. Ablation Study on Memory Structure

In Sec. 4.4 and Fig. 4, we initially explored the relationship between memory allocation strategy and pool ratio of CSM and DAM. Empirically, we found the best setting for these configurations under the fixed-budget constraint. In this section, we aim to answer the following questions:

**Q1:** How sensitive is the model performance to cluster numbers of CSM, i.e., $N^{\text{CSM}}$?

**Q2:** How sensitive is the model performance to key frame numbers of DAM, i.e., $N^{\text{DAM}}$?

As presented in Tab. 8, we conduct two groups of experiments to investigate the model's sensitivity to memory structure configurations, i.e., memory sizes $N^{\text{CSM}}$ and $N^{\text{DAM}}$. In each group, we compare different memory size choices to the baseline row ① and row ③ in Table 4. The results show a scaling trend of accuracy with different memory sizes.

| ID | Memory Component Settings | | | | | Evaluation Results | | | |
|---|---|---|---|---|---|---|---|---|---|
| | CSM | DAM | $N_{Vtokens}$ | CSM Size | DAM Size | EgoSchema | MVBench | Video-MME(w/o) | Average |
| | ✓ | ✓ | 19200 | $60 \times 64$ | $60 \times 256$ | 68.6 | 65.5 | 61.2 | 65.1 |
| | ✓ | ✓ | 15360 | $60 \times 64$ | $45 \times 256$ | 68.3 | 65.3 | 61.0 | 64.9 |
| ① | ✓ | ✓ | 11520 | $60 \times 64$ | $30 \times 256$ | 68.2 | 65.4 | 61.2 | 64.9 |
| | ✓ | ✓ | 7680 | $60 \times 64$ | $15 \times 256$ | 67.5 | 64.9 | 60.8 | 64.4 |
| ③ | ✓ | ✗ | 3840 | $60 \times 64$ | 0 | 66.8 | 64.0 | 60.1 | 63.6 |
| | ✓ | ✗ | 5760 | $90 \times 64$ | 0 | 66.6 | 63.9 | 61.0 | 63.8 |
| ③ | ✓ | ✗ | 3840 | $60 \times 64$ | 0 | 66.8 | 64.0 | 60.1 | 63.6 |
| | ✓ | ✗ | 1920 | $30 \times 64$ | 0 | 65.7 | 63.6 | 58.8 | 62.7 |
| | ✓ | ✗ | 960 | $15 \times 64$ | 0 | 63.0 | 63.0 | 58.3 | 61.5 |

Table 8. **Ablation study of memory structure configurations.** We investigate the model's sensitivity to cluster numbers of CSM and key frame numbers of DAM.

| Score | 0 | 1 | 2 | 3 | 4 | 5 | Total | Average Score |
|---|---|---|---|---|---|---|---|---|
| Right | 8 | 0 | 26 | 111 | 1916 | 2732 | 4793 | 4.53 |
| Wrong | 355 | 290 | 1712 | 82 | 82 | 686 | 3207 | 2.41 |
| Total | 363 | 290 | 1738 | 193 | 1998 | 3418 | 8000 | 3.68 |

Table 9. **Score distribution of a GPT-3.5-based evaluation.** We tested Qwen2-VL-7b on ActivityNet-QA benchmark, using GPT-3.5-turbo-0125 for evaluation. It is observed that many wrong predictions are assigned with a high score "5", leading to a biased result.

Therefore, the results of grid search experiment illustrated in Fig. 4 are reasonable.

## E. Case Study

In this section, we conduct a case study to provide a comprehensive understanding of the performance of models. This study presents a series of visual cases involving various types of videos, each accompanied by a specific question and multiple-choice options to evaluate the performance of three different models: Qwen2-VL [57], LLaVA-OV [27], and the proposed Flash-VStream.

Figs. 7 to 12 present different genres of videos, including documentaries, cartoons, commercials, sports programs and tutorial videos. As shown in these cases, Flash-VStream exhibits strong understanding capabilities in object recognition, action recognition, action reasoning, temporal reasoning, object counting and object reasoning.

## F. Limitations

### F.1. Fail Case Analysis

Flash-VStream may produce incorrect predictions in certain scenarios, such as text-intensive long videos (see Fig. 13) and long videos with rapid scene changes (see Fig. 14). We suggest that these heavily edited video have a different information density distribution compared to native videos, making efficient timing modeling more difficult.

### F.2. GPT-3.5-based Metric for Open-ended VQA

It is worth noting that some previous works [25, 31, 49] follow Video-ChatGPT [41] to test models on open-ended VQA benchmarks [70, 78] based on GPT-3.5-based judgment (GPT accuracy and GPT score). However, we notice that these metrics are highly unstable and prone to bias, so we try to avoid evaluating models on these *LLM-as-a-judge* benchmarks. Since GPT APIs are proprietary and upgrade over time, this evaluation approach lacks reliability, stability and reproducibility [26]. Furthermore, the evaluation can be disturbed by the hallucination of GPT, leading to a biased evaluation result [49]. As presented in Tab. 9, there is always a discrepancy between the distribution of GPT accuracy and GPT score. Therefore, it is still challenging to benchmark the open-ended VQA ability of MLLMs.

## G. Future Work

Future work could focus on enhancing the models' ability to understand edited videos with intensive text or rapid scene transitions, while maintaining the overall efficiency. Another interesting direction for future work would be to investigate reliable evaluation methods for open-ended VQA. Additionally, the techniques developed in this study could be adapted for use in other fields such as robotics and surveillance systems. We hope that our work will inspire further innovations and improvements in these fields, ultimately leading to more intelligent and versatile systems.
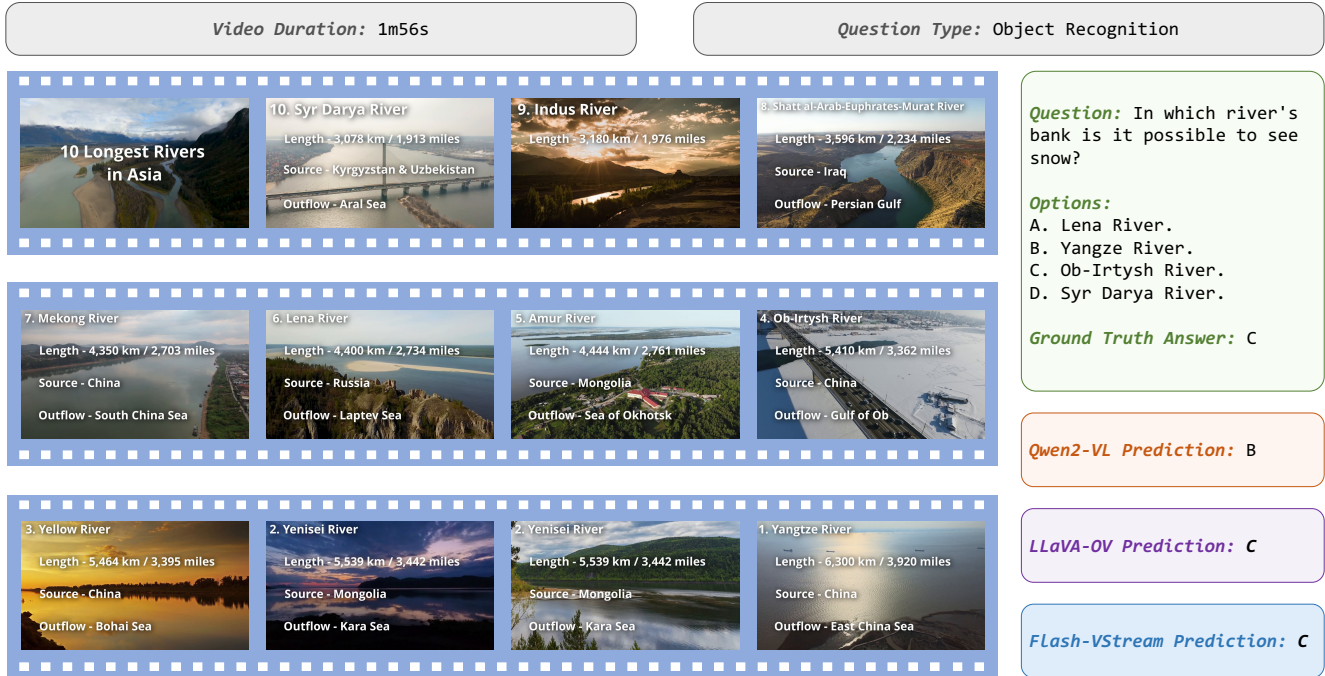
**Video Duration:** 1m56s     **Question Type:** Object Recognition

**Question:** In which river's bank is it possible to see snow?

**Options:**
A. Lena River.
B. Yangze River.
C. Ob-Irtysh River.
D. Syr Darya River.

**Ground Truth Answer:** C

**Qwen2-VL Prediction:** B

**LLaVA-OV Prediction:** C

**Flash-VStream Prediction:** C

Figure 7. **Case Study.** This figure presents a case study on documentary video about the 10 longest rivers in Asia, highlighting their lengths, sources, and outflows. The study includes a question regarding the possibility of seeing snow on the banks of these rivers, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.
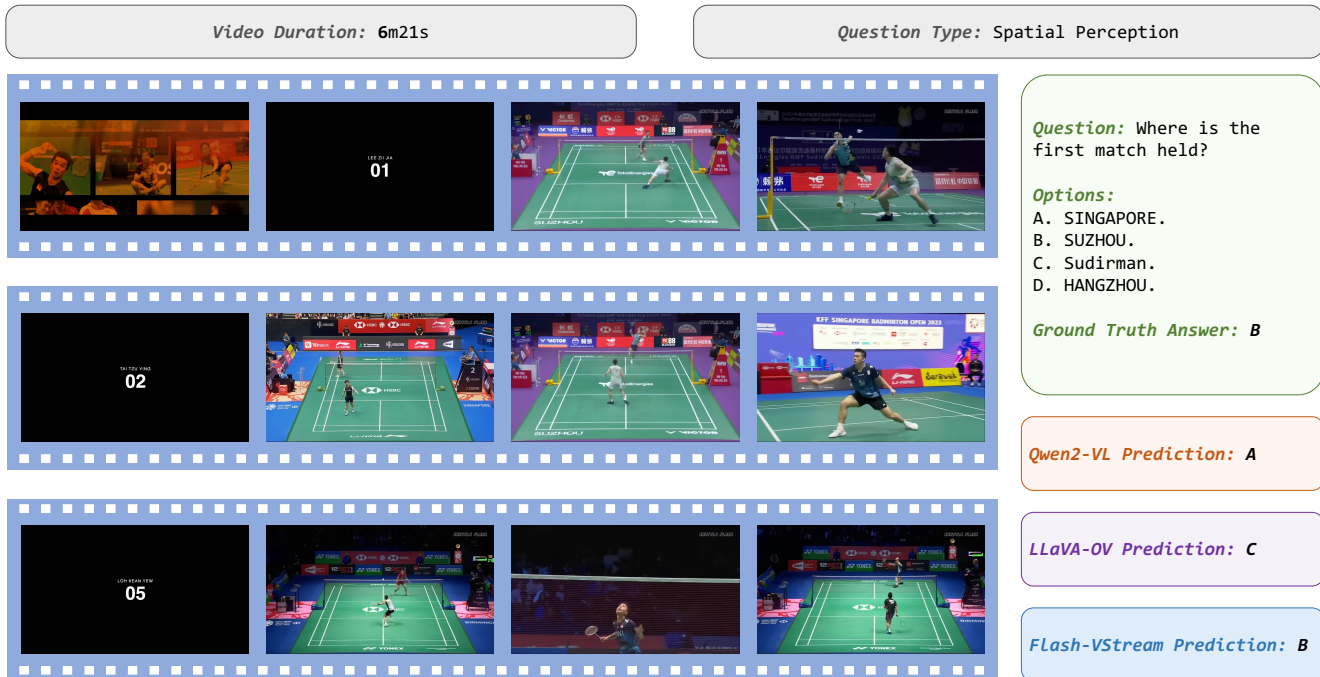


**Video Duration:** 12m6s     **Question Type:** Action Reasoning

**Question:** Why does the mother bird bring a fish to the fox?
**Options:**
A. To thank the fox for raising its own offspring.
B. Because the fox is extremely hungry.
C. To fulfill its baby's wish to help the fox.
D. Because the fox is a friend of its child.
**Ground Truth Answer:** A
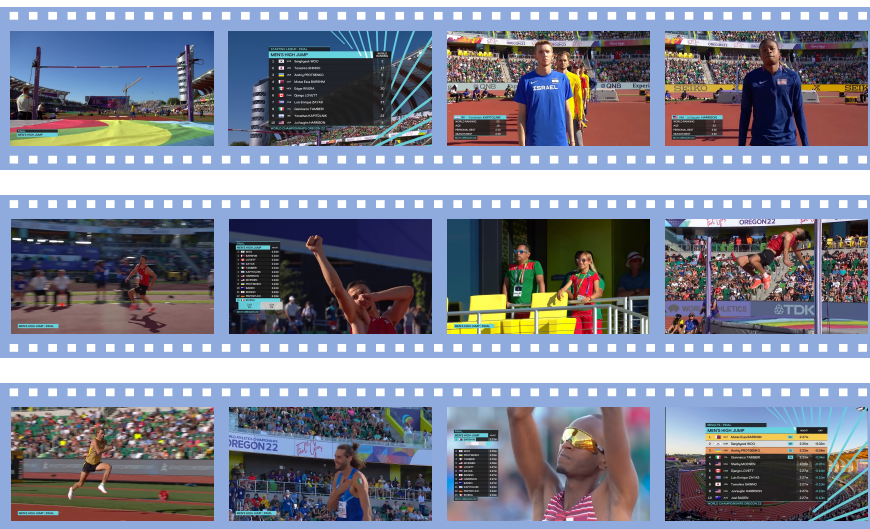
**Qwen2-VL Prediction:** C

**LLaVA-OV Prediction:** C

**Flash-VStream Prediction:** A

Figure 8. **Case Study.** This figure presents a case study involving a cartoon video depicting a mother bird bringing a fish to a fox. The study includes a question about the reason behind this action, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.

Figure 9. **Case Study.** This figure presents a case study involving an advertising video, depicting various scenes including people by the pool, on the beach, and along a coastal hillside. The study includes a question about the number of people on the staircase at the end of the video, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.



Figure 10. **Case Study.** This figure presents a case study involving a sports documentary video of badminton tournaments, depicting various matches and players. The study includes a question about the location of the first match, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.

Figure 11. **Case Study.** This figure presents a case study involving a tutorial video depicting various magic tricks. The study includes a question about the order of events in the video, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.



Figure 12. **Case Study.** This figure presents a case study involving a sports video from a high jump competition, depicting various athletes and their performances. The video frames capture moments of intense competition, showcasing the athletes' skills and determination as they strive to achieve their best performances. The analysis aims to evaluate the models' ability to accurately interpret and predict the outcomes based on visual and contextual cues from the video. The study includes a question about the countries of the top three athletes in the competition, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.
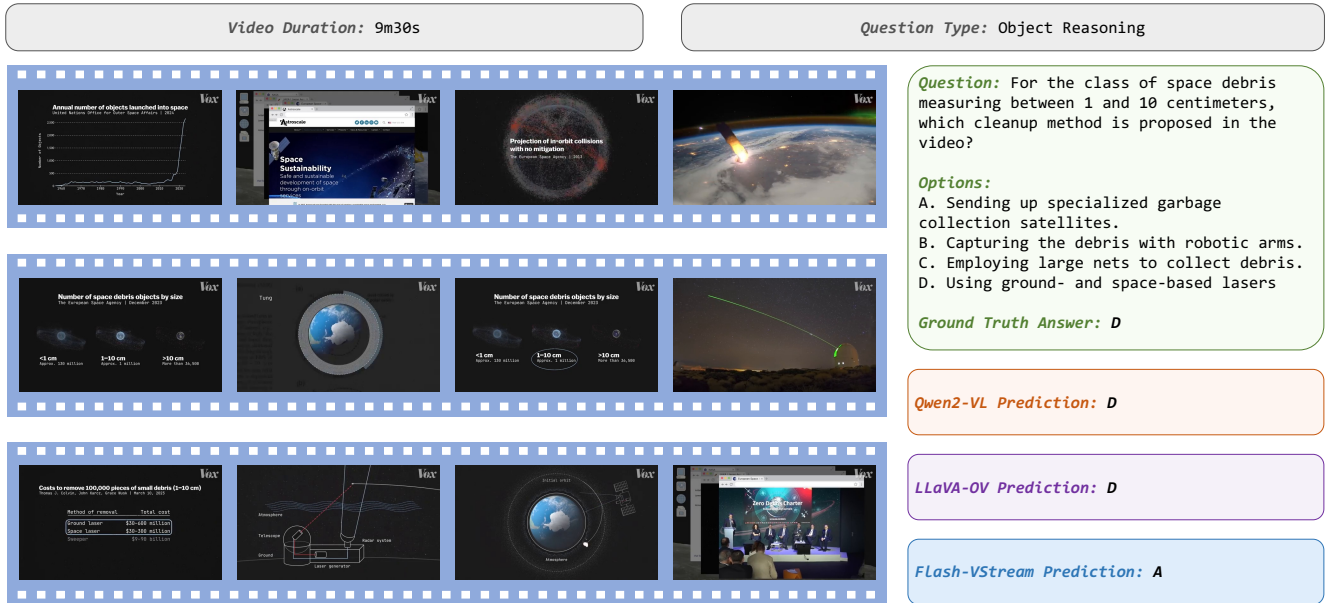
Figure 13. **Fail Case Analysis.** This figure presents a case study involving a video on space debris and proposed cleanup methods. The video frames illustrate various statistics and methods related to space debris, highlighting the challenges and potential solutions for mitigating the growing problem of space junk. The study includes a question about the recommended method for cleaning up space debris measuring between 1 and 10 centimeters, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.
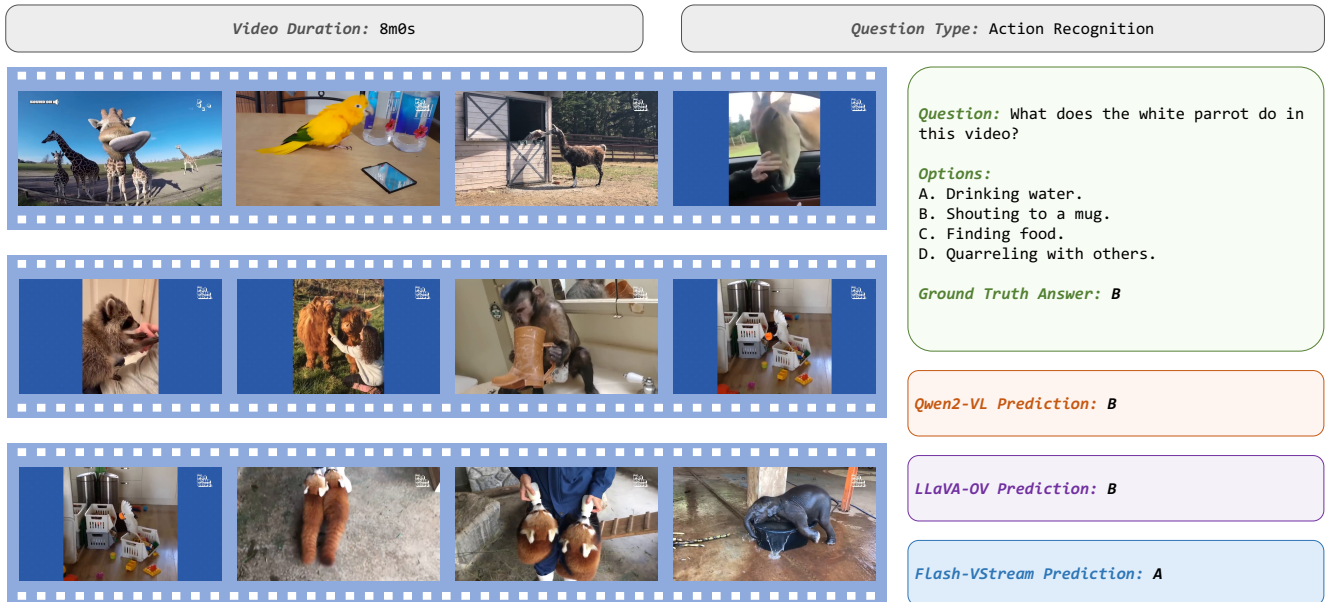


Figure 14. **Fail Case Analysis.** This figure presents a case study involving a video showing various animals and their behaviors. The video frames capture different moments of animal interactions and activities, highlighting the diverse behaviors exhibited by the animals. The study includes a question about the specific action of a white parrot in the video, with multiple-choice options provided. The ground truth answer is indicated, along with the predictions from three different models: Qwen2-VL, LLaVA-OV, and Flash-VStream.