

# StreamMeCo: Long-Term Agent Memory Compression for Efficient Streaming Video Understanding

Anonymous ACL submission

## Abstract

Vision agent memory has shown remarkable effectiveness in streaming video understanding. However, storing such memory for videos incurs substantial memory overhead, leading to high costs in both storage and computation. To address this issue, we propose **StreamMeCo**, an efficient **Stream Agent Memory Compression** framework. Specifically, based on the connectivity of the memory graph, StreamMeCo introduces edge-free minmax sampling for the isolated nodes and an edge-aware weight pruning for connected nodes, evicting the redundant memory nodes while maintaining the accuracy. In addition, we introduce a time-decay memory retrieval mechanism to further eliminate the performance degradation caused by memory compression. Extensive experiments on three challenging benchmark datasets (M3-Bench-robot, M3-Bench-web and Video-MME-Long) demonstrate that under **70%** memory graph compression, StreamMeCo achieves a **1.87×** speedup in memory retrieval while delivering an average accuracy improvement of **1.0%**. *Our code is available in the supplementary materials and will be released on GitHub.*

## 1 Introduction

With the rapidly growing demands of application scenarios such as live streaming (Lu et al., 2018), real-time surveillance (Rao et al., 2022), and autonomous driving (Yurtsever et al., 2020), research on streaming video understanding has become increasingly important. Unlike offline video understanding (Zuo et al., 2025; Tang et al., 2025; Kahatapitiya et al., 2025; Wang et al., 2025b), where models can access the complete video content and user questions prior to inference, streaming video understanding (Wang et al., 2025a; Di et al., 2025; Xiong et al., 2025; Qian et al., 2025; Zhang et al., 2024; Yang et al., 2025a) can only rely on the limited information observed before

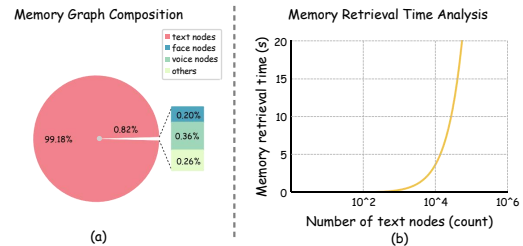


Figure 1: An example from M3-Bench-robot. (a) Composition of the memory graph. (b) Analysis of memory retrieval time with respect to the number of text nodes, indicating that the *increment of memory nodes makes the retrieval time not acceptable.*

the user’s question arrives, making how to efficiently process continuously incoming visual information a critical challenge. Current research on streaming video understanding primarily focuses on achieving long-video understanding by carefully handling the vision tokens and KV cache (Yang et al., 2025b; Yao et al., 2025a). However, for extremely long videos, these methods still suffer from the loss of critical visual information and struggle to maintain long-term consistency of entities such as human identities.

To tackle this challenge, agent memory-based methods are introduced to organize the information of videos as memory in plaintext. For example, M3-Agent (Long et al., 2025) leverages agent memory mechanism to model streaming video as a memory graph, representing entities with face and voice nodes and abstracting rich visual content into text nodes, thereby preserving information integrity and long-term entity consistency. Nevertheless, as the memory graph scales up, storage and retrieval efficiency become major bottlenecks. As shown in Figure 1, the increasing number of text nodes leads to a rapidly increasing memory retrieval latency, which hinders the real-time answering of questions.

To address this, we propose **StreamMeCo**, an efficient training-free Long-Term **Stream Agent**

*Memory Compression* framework, achieving the first memory compression of industrial-level streaming video agent and supporting efficient memory graph retrieval. Based on the connectivity of the memory graph, we design a dual-branch memory compression method. For isolated text nodes, we use the *Edge-free Minmax sampling (EMsampling)* module for compression. First, we cluster the text nodes into different groups, then for each cluster, we first select the node closest to the cluster center and add it to the selected set, then calculate the minimum distance between each unselected node and the selected set, choosing the farthest point. This process continues until the number of selected nodes meets the pre-defined requirement. For connected text nodes, we use the *Edge-aware Weighting pruning (EW-pruning)* module. By calculating the weight edge matrix between text nodes and face nodes or voice nodes, we determine the entity importance of each text node. We then combine the embedding similarity of the text nodes to select the most dissimilar text nodes, and ultimately determine the retained text nodes based on a comprehensive fusion score.

Additionally, to reduce the performance degradation caused by compressing the memory graph, inspired by the human memory mechanism, we design a *Time-decay Memory Retrieval (TMR)* mechanism for efficiently retrieving the memory graph. We calculate the total similarity between text nodes in the memory graph and the model’s query for each time segment. Based on the total similarity of each time segment, we dynamically allocate the number of text nodes to different segments. Furthermore, to simulate the memory decay mechanism of the human brain, we assign nonlinear time-varying decay to the text nodes of each time segment, thereby adaptively balancing long-term memory and recent critical information during the retrieval stage, and guaranteeing the retrieval of more memories from recent segments.

In summary, the main contributions of this paper are as follows:

- We propose *edge-free minmax sampling* module and *edge-aware weighting pruning* module, which efficiently compress the memory graph based on its structural properties.
- We introduce a novel *time-decay memory retrieval* mechanism, which has been experimentally proven to significantly outperform previous retrieval frameworks.

- We demonstrate the effectiveness of *Stream-MeCo* through extensive experiments. On three high-difficulty benchmark datasets, under approximately **70%** memory graph compression, the proposed method achieves a **1.87×** speedup in memory retrieval while delivering an average accuracy improvement of **1.0%**.

## 2 Related Work

### 2.1 Streaming Video Understanding

Streaming video understanding enables a model to continuously process incoming frames, and upon receiving a user question  $Q_t$  at time  $t$ , generate responses based on the accumulated content  $C_{1:t}$ . Existing methods fall into two main paradigms: visual token compression and KV cache compression. The former selects informative tokens either in pixel space or after feature encoding (Chen et al., 2025; Dorovatas et al., 2025; Zeng et al., 2025; Yao et al., 2025a; Li et al., 2025; Huang et al., 2025; Chen et al., 2024). While the latter maintains a short-term sliding window and a long-term memory stack at the LLM level to preserve information across temporal scales (Yang et al., 2025b; Kim et al., 2025; Ning et al., 2025; Di et al., 2025; Xu et al., 2025b). Recently, M3-Agent (Long et al., 2025) departs from these designs by integrating long-term agent memory into streaming video understanding, showing strong potential for modeling long-range temporal dependencies and maintaining entity consistency.

### 2.2 Long-Term Agent Memory

Long-term agent memory is designed to provide LLM-based agents with persistent memory beyond a single context window. Depending on the type of input it handles, it can be categorized into unimodal and multimodal forms. Unimodal memory (Yu et al., 2025a; Zhou et al., 2025; Chhikara et al., 2025; Packer et al., 2023) typically maintains consistency across multi-turn reasoning by constructing structured representations, latent embeddings, or summaries. In contrast, multimodal memory (Yu et al., 2025b; Zhang et al., 2025; Yao et al., 2025b; Fan et al., 2025) integrates visual, linguistic, and other signals into a unified space, updating them to support stable long-term reasoning and decision-making. Building on this foundation, M3-Agent is the first to introduce long-term agent memory into streaming video scenarios, enabling application in real-time settings.

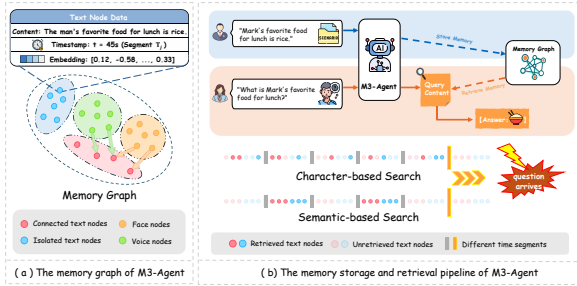


Figure 2: The overview of M3-Agent. (a) The memory graph of M3-Agent. (b) The memory storage and retrieval pipeline of M3-Agent, illustrating two different retrieval strategies.

### 3 Preliminaries

Graph-based agent memory has demonstrated great potential across a variety of research areas. Our method exhibits strong generality and can be readily adapted to different graph-based agent memory frameworks, while this paper focuses on its application in streaming video understanding scenarios. Next, we provide a detailed introduction to M3-Agent (Long et al., 2025).

#### 3.1 Memory Graph of M3-Agent

As shown in Figure 2(a), the memory graph includes face nodes, voice nodes, text nodes, and other relevant information such as edge weights. Each text node contains content, embedding vectors, and a timestamp of the content. Among these nodes, only some voice nodes are connected to text nodes, and some face nodes are connected to text nodes, while there are no internal connections within face nodes, voice nodes, or text nodes, nor between voice and face nodes. Based on the graph’s connectivity, the content of isolated text nodes is similar to: "The man’s favorite food for lunch is rice." with no specific entity references. In contrast, the content of connected text nodes is like: "Mark’s favorite food for lunch is rice." which explicitly includes entity information.

#### 3.2 Retrieval of Memory Graph

As shown in Figure 2(b), after a question is input into M3-Agent, the model performs reasoning and outputs the query content to be retrieved. Based on this content, the model retrieves from the memory graph, which is primarily divided into two approaches: *character-based* and *semantic-based*. If the retrieval content contains characters like "character id", the model returns the  $k$  most similar content based on the similarity between the embedding of the query content and the

content embedding of the text nodes, as the retrieved memory. If there is no "character id" identifier, semantic retrieval is performed. In this case, the model segments all text nodes based on their timestamps, selecting the most similar node as the representative for each time segment. Finally, the two most similar text nodes are selected, and all text node contents at the timestamps of these two nodes are returned as the retrieved memory.

## 4 Methodology

### 4.1 Overview

The overview of StreamMeCo is illustrated in Figure 3. The framework consists of two components, achieving efficient compression of memory graph while supporting fast and accurate retrieval.

### 4.2 Text Memory Compression

Based on the connectivity of the memory graph, we categorize text nodes into isolated nodes and connected nodes, and accordingly propose a dual-branch memory compression method. Different compression strategies are applied to these two types of text nodes.

**Edge-free Minmax Sampling.** For isolated text nodes, we have carefully designed the EMSampling module to select the most representative text nodes. Specifically, given  $N_1$  text nodes, we use the *Spherical KMeans Algorithm* (Dhillon and Modha, 2001) to cluster them into  $N_1 \times a$  clusters based on the embeddings of the text node contents, as shown in the following formula:

$$\{C_j\} = KMeans(\{e_i\}, a), \quad (1)$$

where  $e_i$  ( $i = 1, 2, \dots, N_1$ ) represents the embedding vector of the text node contents, and clustering ratio  $a$  is used to control the number of clusters, while  $C_j$  ( $j = 1, 2, \dots, N_1 \times a$ ) represents the  $j$ -th cluster.

Then, we assume the retention ratio of isolated text nodes is  $\alpha$ , and distribute them evenly across each cluster. The number of text nodes retained in the  $j$ -th cluster is  $|C_j| \times \alpha$ , where  $|C_j|$  represents the number of text nodes in the  $j$ -th cluster.

Finally, we adopt an effective *minmax sampling* strategy to select the text nodes to be retained from each cluster. Specifically, for the  $j$ -th cluster, we first select the text node closest to the cluster center and add it to the selected node set  $S_j$ . Then, we compute the distance from each unselected text

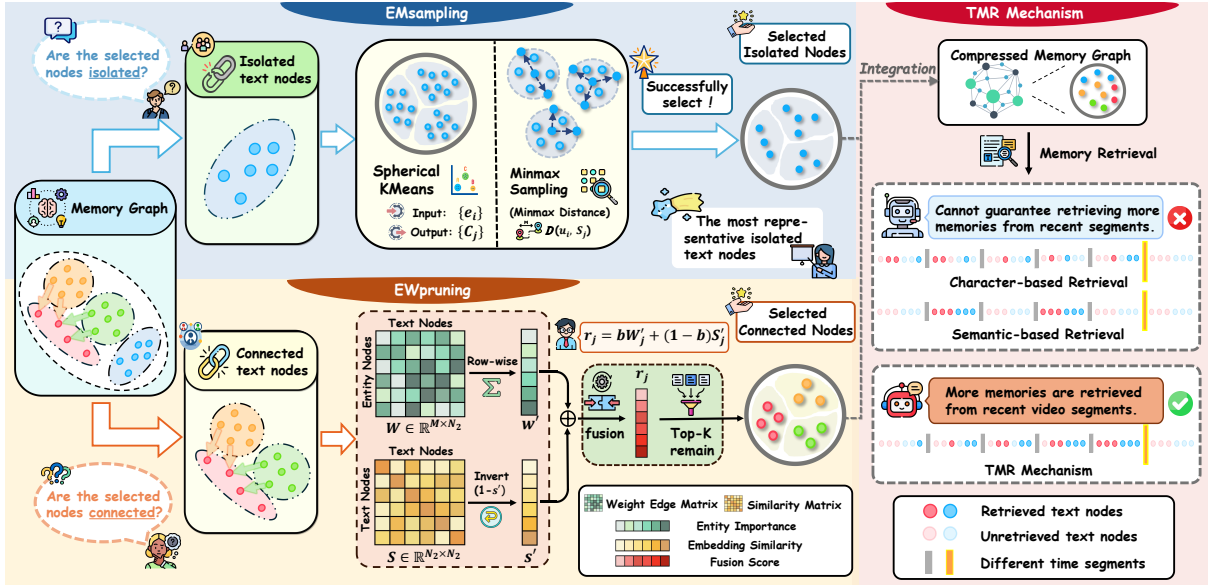


Figure 3: The overview of StreamMeCo. Efficient memory graph compression is achieved via the EMSampling and EWpruning modules, while the TMR mechanism prioritizes retrieving memories from more recent segments, enabling fast and accurate memory retrieval.

node in the unselected node set  $U_j$  to the selected node set  $S_j$ , and select the node with the largest distance to add to  $S_j$ . The distance from each unselected node to the selected node set  $S_j$  is defined as the minimum distance to any node in the selected set, as shown in the following formula:

$$D(u_i, S_j) = \min_{s_k \in S_j} \|e_{u_i} - e_{s_k}\|, \quad (2)$$

where  $u_i \in U_j$  is an unselected node, and  $s_k \in S_j$  is a selected node, while  $e_{u_i}$  and  $e_{s_k}$  represent the embedding vectors of the unselected and selected nodes, respectively.

We select unselected node  $u_{i^*}$  with the farthest distance and add it to the selected node set  $S_j$ :

$$u_{i^*} = \arg \max_{u_i \in U_j} D(u_i, S_j). \quad (3)$$

Until the number of selected nodes reaches  $|C_j| \times \alpha$ , the number of retained text nodes.

**Edge-aware Weighting Pruning.** For connected text nodes, we carefully design the EWpruning module, thereby constructing a more compact yet informative memory representation to enable efficient node pruning. Specifically, given  $N_2$  text nodes that are connected to face and voice nodes via edges, we construct an entity-based weight edge matrix  $W \in \mathbb{R}^{M \times N_2}$ , where  $M$  denotes the total number of face and voice nodes. Each entry  $w_{ij}$  represents the edge weight between the  $i$ -th face or voice node and the  $j$ -th

text node, reflecting the importance of the  $j$ -th text node to the  $i$ -th face or voice node.

For the weight matrix  $W$ , the row-wise summation is computed as  $W' = \sum_{i=1}^M w_{ij}$ , where  $W' \in \mathbb{R}^{1 \times N_2}$  represents the importance scores of the  $N_2$  connected text nodes with respect to the set of face and voice entity nodes. Subsequently, we normalize  $W'$  to the range  $[0, 1]$  to obtain the final *entity importance* for each text node.

Next, for the  $N_2$  text nodes, we compute their pairwise similarities to obtain a similarity matrix  $S \in \mathbb{R}^{N_2 \times N_2}$ , where  $s_{ij}$  denotes the embedding similarity between the  $i$ -th and  $j$ -th text nodes. We then perform a row-wise summation on the similarity matrix, computed as  $S' = \sum_{i=1}^{N_2} s_{ij}$ , where  $S' \in \mathbb{R}^{1 \times N_2}$  represents the total similarity of each text node to all other text nodes. Since we aim to select the text nodes with higher diversity, we normalize  $S'$  to the range  $[0, 1]$  and then invert it as  $1 - S'$ , resulting in the *embedding similarity* for each text node relative to the connected text node set. For brevity, we continue to denote the processed entity importance and embedding similarity as  $W'$  and  $S'$ .

We assume the retention ratio for connected text nodes is  $\beta$ , meaning that we need to keep a total of  $N_2 \times \beta$  text nodes. For the  $j$ -th connected text node, we compute its normalized entity importance  $W'_j$  and its embedding similarity  $S'_j$  through a weighted combination. The final fusion score is computed as:

$$r_j = bW'_j + (1 - b)S'_j, \quad (4)$$

where the balancing coefficient  $b$  controls the trade-off between entity importance and embedding similarity.

Finally, we retain the top  $N_2 \times \beta$  text nodes with the highest scores and prune the others.

### 4.3 Time-decay Memory Retrieval

To reduce the performance degradation caused by compressing the memory graph, inspired by the human memory mechanism, we propose a TMR framework. Specifically, given the query content generated by the model in Section 3.2, we first compute the embedding similarity between the query and the content embeddings of all text nodes. According to their timestamps, these text nodes are grouped into a series of chronologically ordered time segments  $T_j$  ( $j = 1, 2, \dots, t$ ). At time step  $t$ , when the system receives a query, the relevance score of each time segment is computed by aggregating the similarity scores of all text nodes within that segment as  $E_j = \sum_{i=1}^{|T_j|} s_{ij}$ , where  $s_{ij}$  denotes the embedding similarity between the query and the  $i$ -th text node in  $T_j$ ,  $E_j$  denotes the aggregated similarity within the time segment  $T_j$ , and  $s_{ij}$  denotes the embedding similarity between the query content and the  $i$ -th text node in time segment  $T_j$ , and  $|T_j|$  represents the number of text nodes in segment  $T_j$ .

Next, to simulate human memory decay, we first compute the average similarity of each time segment as  $\bar{E}_j = \frac{E_j}{|T_j|}$ . Based on this averaged relevance score, we then apply a temporal decay function to model the fading of older memories:

$$E'_j = \bar{E}_j \cdot e^{-\lambda(t-j)}, \quad (5)$$

where  $j = 1, 2, \dots, t$  indicates the temporal position of the segment, and the memory decay coefficient  $\lambda$  controls the rate of memory decay.

To be consistent with the number of retrieved text nodes in Section 3.2, we assume that the model needs to retrieve a total of  $k$  text nodes as the final memory output. Accordingly, we compute the number of nodes selected from each time segment as:

$$Num_j = \frac{E'_j}{\sum_{i=1}^t E'_i} \times k, \quad (6)$$

where  $Num_j$  denotes the number of nodes selected from the  $j$ -th segment.

Finally, within each time segment, we select the top-ranked text nodes according to their similarity to the query content until reaching  $Num_j$  nodes, which constitute the final retrieved memory.

## 5 Experiments

### 5.1 Experiments Setup

**Datasets.** We select three challenging open-source benchmark datasets, M3-Bench-robot, M3-Bench-web (Long et al., 2025) and Video-MME-Long (Fu et al., 2025), to evaluate the effectiveness of StreamMeCo. M3-Bench-robot is a streaming video dataset, while the others are offline video datasets. Consistent with the original M3-Agent setup, we use GPT-4o to assess the answer quality and employ text-embedding-3-large to encode the query content generated by M3-Agent to supporting subsequent memory graph retrieval. For more details, please refer to the Appendix A.1.

**Baselines.** Our baselines include closed-source MLLMs: Gemini-1.5-Pro (Team et al., 2024), GPT-4o (Hurst et al., 2024), Gemini-Agent and Gemini-GPT4o-Hybrid (Long et al., 2025). open-source MLLMs: Qwen2.5-VL-7B-Instruct (Bai et al., 2025), Qwen2.5-Omni-7B (Xu et al., 2025a). As well as streaming video models: MovieChat (Song et al., 2024), MA-LMM (He et al., 2024), VideoChat-Online (Huang et al., 2025), Flash-VStream (Zhang et al., 2024), TimeChat-Online (Yao et al., 2025a), StreamForest (Zeng et al., 2025), and M3-Agent (Long et al., 2025). For more introductions about the baselines, please refer to the Appendix A.2.

**Implementation Details.** All experiments are conducted on two NVIDIA A100 (80G) GPUs. To reduce error, each experiment is executed three times and the average result is reported. In our setup, the clustering ratio  $a$  is set to 0.05, the balancing coefficient  $b$  is set to 0.1, and the memory decay coefficient  $\lambda$  is set to 0.1. For the M3-Bench-robot and M3-Bench-web, we use the memory graphs provided by the M3-Agent paper, for Video-MME-Long, we replace the Gemini-1.5-Pro API with Gemini-2.5-Pro, while keeping all other settings consistent with M3-Agent. We apply compression only to the memory graph produced by M3-Agent, while keeping all other components identical to the original model.

Dataset	M3-Bench-robot						M3-Bench-web						Video-MME-Long	Avg.
	ME	MH	CM	PU	GK	All	ME	MH	CM	PU	GK	All		
<i>Closed-source MLLMs</i>														
Gemini-1.5-Pro	6.5	7.5	8.0	9.7	7.6	8.0	18.0	17.9	23.8	23.1	28.7	23.2	38.0	23.1
GPT-4o	9.3	9.0	8.4	10.2	7.3	8.5	21.3	21.9	30.9	27.1	39.6	28.7	38.8	25.3
Gemini-Agent	15.8	17.1	15.3	20.0	15.5	16.9	29.3	20.9	33.8	34.6	45.0	34.1	55.1	35.4
Gemini-GPT4o-Hybrid	21.3	25.5	22.7	28.8	23.1	24.0	35.9	26.2	37.6	43.8	52.2	41.2	56.5	40.6
<i>Open-source MLLMs</i>														
Qwen2.5-VL-7b-Instruct	2.9	3.8	3.6	4.6	3.4	3.4	11.9	10.5	13.4	14.0	20.9	14.9	46.9	21.7
Qwen2.5-Omni-7b	2.1	1.4	1.5	1.5	2.1	2.0	8.9	8.8	13.7	10.8	14.1	11.3	42.2	18.5
<i>Streaming Video Models</i>														
MovieChat [CVPR24]	13.3	9.8	12.2	15.7	7.0	11.2	12.2	6.6	12.5	17.4	11.1	12.6	19.4	14.4
MA-LMM [CVPR24]	25.6	23.4	22.7	39.1	14.4	24.4	26.8	10.5	22.4	39.3	15.8	24.3	17.3	22.0
VideoChat-Online [CVPR25]	31.7	24.7	30.5	43.1	17.1	29.9	34.3	18.9	34.1	48.3	23.1	32.7	45.9	36.2
Flash-VStream [ICCV25]	21.6	19.4	19.3	24.3	14.1	19.4	24.5	10.3	24.6	32.5	20.2	23.6	25.0	22.7
TimeChat-Online [MM25]	35.6	29.4	31.5	44.3	22.3	33.6	38.8	22.8	38.4	51.6	28.0	36.6	49.2	39.8
StreamForest [NIPS25]	33.4	29.4	31.3	44.3	19.0	32.1	37.9	22.3	39.9	51.7	27.6	36.5	51.9	40.2
M3-Agent [Arxiv25]	30.9	29.4	29.6	41.1	22.6	30.3	44.9	25.6	44.8	58.6	53.7	47.9	<b>56.0</b>	44.7
+Random ( $\downarrow$ 30%)	28.9	29.4	27.7	40.5	18.0	28.5	42.4	24.7	40.1	54.8	50.8	44.6	54.2	42.4
+Clustering ( $\downarrow$ 30%)	30.0	32.9	30.0	42.5	20.8	29.6	41.2	26.3	41.7	55.9	50.9	45.3	54.6	43.2
+DART ( $\downarrow$ 30%)	29.5	24.7	25.8	40.1	21.7	29.1	40.4	26.0	41.3	55.6	51.2	45.4	54.8	43.1
+StreamMeCo ( $\downarrow$ 30%)	35.7	32.9	35.5	44.0	27.5	<b>34.6</b>	46.4	32.2	46.5	61.8	55.8	<b>50.7</b>	<b>55.2</b>	<b>46.8</b>
+StreamMeCo ( $\downarrow$ 50%)	34.8	34.1	32.4	43.1	24.5	<b>33.4</b>	45.4	31.3	45.0	59.2	52.7	<b>48.9</b>	<b>56.6</b>	<b>46.3</b>
+StreamMeCo ( $\downarrow$ 70%)	34.9	31.8	32.1	42.7	24.8	<b>34.2</b>	43.8	28.9	43.9	57.9	52.9	<b>47.9</b>	54.9	<b>45.7</b>

Table 1: Performance comparison of different models on M3-Bench-robot, M3-Bench-web and Video-MME-Long. The best three results for the streaming video models are highlighted in **bold black font**.

## 5.2 Main Results

To comprehensively evaluate the performance of our method, we conduct comparisons with 13 competitive baselines. For more details about the baselines, please refer to the Appendix A.3. As shown in Table 1, the overall performance of **StreamMeCo** significantly surpasses that of existing closed-source and open-source MLLMs, as well as other streaming video models, demonstrating the effectiveness of the agent memory mechanism for streaming video understanding. Even with **70%** of text nodes compressed, StreamMeCo still achieves an average accuracy improvement of **1.0%** across all datasets compared to the uncompressed M3-Agent.

In addition, we have constructed three memory graph compression methods and compared them with StreamMeCo, with the compression ratio controlled at 30%. In the random method, we randomly selected nodes from two types of text nodes for compression. For the clustering method, we used the KMeans clustering algorithm to evaluate the differences between our method and tra-

ditional clustering approaches. Specifically, we divided the two types of text nodes into a series of smaller clusters, each of which retains the text node closest to the cluster center. For the DART method, we were inspired by the DART (Wen et al., 2025). Specifically, for the two types of text nodes, we first randomly selected 2% of the pivot nodes, then calculated the total similarity of the remaining nodes with these pivot nodes, prioritizing the compression of nodes with higher similarity to the pivot nodes, until the target compression ratio was reached. As shown in the Table 1, StreamMeCo significantly outperforms these methods in terms of performance.

For the M3-Bench-robot dataset, all three compression ratios lead to varying degrees of performance improvement. In contrast, the performance gains on the M3-Bench-web and Video-MME-Long datasets are relatively limited. This is because M3-Bench-robot consists of real-world filming captured from a first-person robotic perspective, involving daily environments and human-robot interactions, which naturally contain

Components			M3-Bench-Robot	M3-Bench-Web
A	B	C	Accuracy	Accuracy
-	-	-	30.3 (100.0%)	47.9 (100.0%)
			28.5 (94.1%)	44.6 (93.1%)
✓			30.3 (100.0%)	45.7 (95.4%)
	✓	✓	30.1 (99.3%)	46.1 (96.2%)
✓	✓		30.6 (101.0%)	46.1 (96.2%)
✓		✓	31.4 (103.6%)	46.9 (97.9%)
✓	✓	✓	<b>31.5 (104.0%)</b>	<b>47.3 (98.7%)</b>

Table 2: Ablation study of the compression components. The first row reports the accuracy of the M3-Agent without compression. **A**: EMsampling; **B**: entity importance; **C**: embedding similarity.

a higher level of redundancy. In comparison, both M3-Bench-web and Video-MME-Long are sourced from YouTube, where videos typically exhibit stronger narrative structures or scripted production, resulting in more concise content and relatively less redundancy, leaving limited room for effective compression.

### 5.3 Ablation Study

Due to the increased overlap between the nodes compressed by StreamMeCo and those retained through random compression at higher text-node compression ratios, the performance differences between modules become less distinguishable. To highlight the contribution of each component in our compression, we conduct the components ablation analysis under a 30% compression ratios.

#### Compression Components Ablation Analysis.

This component analysis focuses solely on the compression method and does not involve the TMR mechanism. As shown in Table 2, we evaluate the performance of EMsampling (A), entity importance (B), and embedding similarity (C) under different combinations. The results indicate that enabling any single component leads to varying degrees of performance improvement. Specifically, the EMsampling module effectively selects the most representative text nodes from the isolated-node subset, while for connected text nodes, our method jointly considers both entity importance and embedding similarity, yielding a more compact yet informative memory representation. When all three components are used together, the model achieves its best performance on both benchmark datasets, reaching accuracies of **31.5%** and **47.3%**, corresponding to **104.0%** and **98.7%** of the original M3-Agent, respectively.

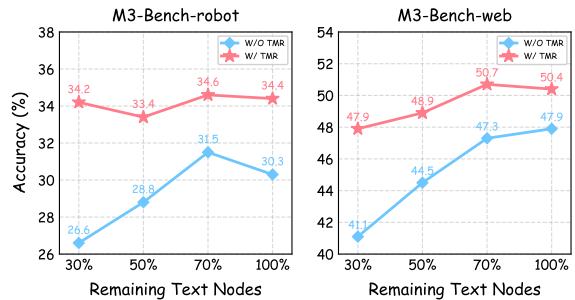


Figure 4: Impact of the TMR Mechanism on M3-Bench-robot and M3-Bench-web at different compression ratios.

In contrast, random compression achieves only 28.5% and 44.6% accuracy on the two benchmark datasets, resulting in substantial performance loss.

We observe that compression causes greater performance degradation on M3-Bench-web than on M3-Bench-robot, further supporting our earlier hypothesis that the text nodes in M3-Bench-robot contain higher redundancy.

#### TMR Impact Analysis.

As shown in Figure 4, when the TMR mechanism is not applied, the model’s performance drops significantly on both benchmark datasets once the compression ratio exceeds 50%. In contrast, after introducing TMR, the model accuracy improves by **3.3%** over the original M3-Agent even without any compression. Compressing 30% of the text nodes further enhances the performance. More importantly, under a 70% compression setting, the model does not suffer from noticeable performance degradation and performs comparably to or even better than the uncompressed M3-Agent. These observations demonstrate that TMR mechanism can significantly reduce the accuracy loss caused by memory graph compression and effectively retrieve memories from more recent segments for reasoning.

#### Time Efficiency Analysis.

In this section, we analyze the impact of the EMsampling module, EWpruning module, and TMR mechanism on the time efficiency of M3-Agent, with a focus on the memory retrieval time after converting the model query into embeddings. As shown in Figure 5, when only applying the EMsampling and EWpruning modules for compression, SS-Time decreases as the compression ratio increases, while Avg-Retrieval increases correspondingly (M3-Agent triggers up to five memory retrievals per question by default). We attribute this phenomenon to the reduced memory capacity

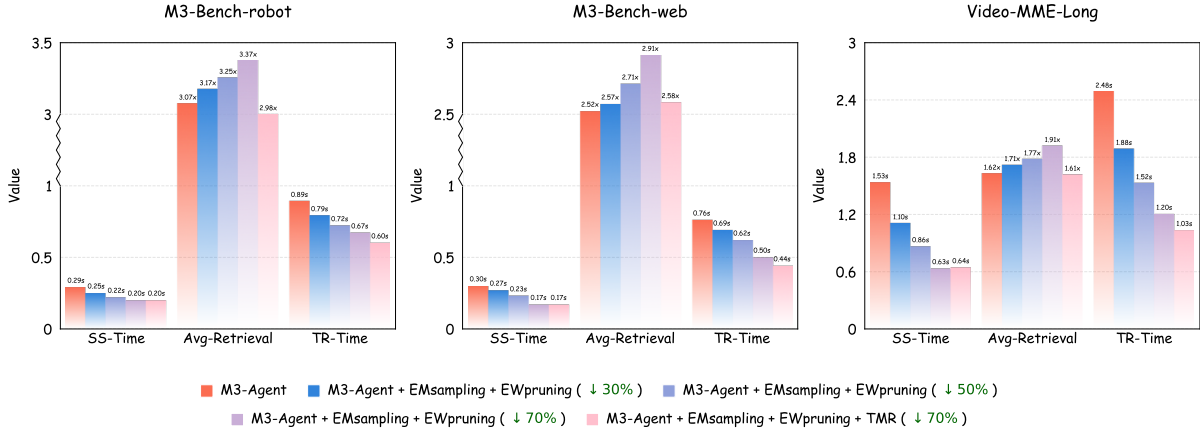


Figure 5: Time Efficiency Analysis on M3-Bench-Robot, M3-Bench-Web and Video-MME-Long. **SS-Time:** Time to compute and sort the similarity between the query embedding and each text node’s content embedding; **Avg-Retrieval:** Average number of memory graph queries per question; **TR-Time:** Average total time spent querying the memory graph for each question after obtaining the model’s query embedding.

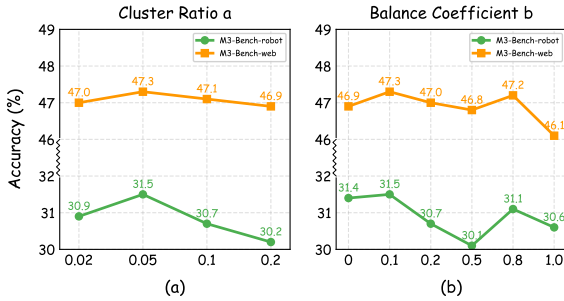


Figure 6: Impact of cluster ratio  $a$  and balance coefficient  $b$  on M3-Bench-robot and M3-Bench-web.

after compression, which limits the amount of effective information the model can retrieve, lowers the confidence in the retrieval results, and requires more retrieval rounds to finalize the answer.

In contrast, with the introduction of the TMR mechanism, Avg-Retrieval significantly decreases, as TMR mechanism prioritizes more reliable and higher-confidence memory entries, enabling the model to obtain sufficient information with fewer retrieval iterations, thus effectively reducing the overall processing time.

When compressing **70%** of the memory, our method achieves an average **1.87×** speedup in memory retrieval time across all datasets.

**Effect of Cluster Ratio  $a$ .** As shown in Figure 6(a), the model performs best on both datasets when the cluster ratio  $a = 0.05$ . When  $a$  is smaller, the number of clusters decreases, and each cluster contains more samples, which makes it more likely that text nodes with large semantic differences are assigned to the same cluster, thereby affecting the preservation of semantic distinctions. In contrast, when  $a$  is larger, the number of clus-

ters increases, and each cluster contains fewer samples, making it difficult to form stable and representative semantic structures, ultimately degrading the performance of the compressed model.

**Effect of Balance Coefficient  $b$ .** As shown in Figure 6(b), the model achieves the best performance on both datasets when the parameter balance coefficient  $b = 0.1$ . When  $b$  is further decreased or increased, the performance drops to varying degrees. This observation indicates that an excessively small balance coefficient prevents the model from effectively leveraging entity-importance information, whereas an excessively large balance coefficient causes the model to overly rely on entity importance, thereby ignoring embedding similarity information and ultimately degrading the overall performance.

In addition, we analyzed the impact of the memory decay coefficient  $\lambda$ , and conducted other ablation experiments on the Video-MME-Long. Detailed results can be found in Appendix B.

## 6 Conclusion

We propose **StreamMeCo**, the first memory compression framework specifically designed for industrial-scale streaming video agent. It achieves efficient compression and accurate retrieval of memory graphs through the integration of our EMSampling and EWpruning modules, along with the TMR mechanism. Experimental results show that even with a **70%** reduction in memory graph size, StreamMeCo achieves a **1.87×** speedup in memory retrieval, while delivering an average accuracy improvement of **1.0%**.

## 583 Limitations

584 Due to the frequent need to call the Gemini-  
585 2.5-Pro and text-embedding-3-large APIs during  
586 memory graph generation, as well as the substan-  
587 tial time required to generate the memory graphs,  
588 and since M3-Agent currently only constructs  
589 and provides memory graphs for M3-Bench-robot  
590 and M3-Bench-web, we were limited by budget  
591 and time constraints and were only able to test  
592 the Video-MME-Long dataset additionally. We  
593 hope that future work can test more benchmark  
594 datasets to more comprehensively validate the per-  
595 formance of our method.

## 596 References

597 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang,  
598 Wenbin Ge, Siboz Song, Kai Dang, Peng Wang, Shi-  
599 jie Wang, Jun Tang, and 1 others. 2025. Qwen2.5-vl  
600 technical report. *arXiv preprint arXiv:2502.13923*.

601 Joya Chen, Zhaoyang Lv, Shiwei Wu, Kevin Qinghong  
602 Lin, Chenan Song, Difei Gao, Jia-Wei Liu, Ziteng  
603 Gao, Dongxing Mao, and Mike Zheng Shou.  
604 2024. Videollm-online: Online video large lan-  
605 guage model for streaming video. In *Proceedings*  
606 *of the IEEE/CVF Conference on Computer Vision*  
607 *and Pattern Recognition*, pages 18407–18418.

608 Xueyi Chen, Keda Tao, Kele Shao, and Huan Wang.  
609 2025. Streamingtom: Streaming token compression  
610 for efficient video understanding. *arXiv preprint*  
611 *arXiv:2510.18269*.

612 Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet  
613 Singh, and Deshraj Yadav. 2025. Mem0: Building  
614 production-ready ai agents with scalable long-term  
615 memory. *arXiv preprint arXiv:2504.19413*.

616 Inderjit S Dhillon and Dharmendra S Modha. 2001.  
617 Concept decompositions for large sparse text data  
618 using clustering. *Machine learning*, 42(1):143–175.

619 Shangzhe Di, Zhelun Yu, Guanghao Zhang, Haoyuan  
620 Li, Tao Zhong, Hao Cheng, Bolin Li, Wangui He,  
621 Fangxun Shu, and Hao Jiang. 2025. Streaming  
622 video question-answering with in-context video kv-  
623 cache retrieval. *arXiv preprint arXiv:2503.00540*.

624 Vaggelis Dorovatas, Soroush Seifi, Gunshi Gupta, and  
625 Rahaf Aljundi. 2025. Recurrent attention-based  
626 token selection for efficient streaming video-llms.  
627 *arXiv preprint arXiv:2510.17364*.

628 Yue Fan, Xiaojian Ma, Rongpeng Su, Jun Guo, Rujie  
629 Wu, Xi Chen, and Qing Li. 2025. Embodied videoa-  
630 gent: Persistent memory from egocentric videos  
631 and embodied sensors enables dynamic scene un-  
632 derstanding. In *Proceedings of the IEEE/CVF In-*  
633 *ternational Conference on Computer Vision*, pages  
634 6342–6352.

635 Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li,  
636 Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
637 Zhou, Yunhang Shen, Mengdan Zhang, and 1 oth-  
638 ers. 2025. Video-mme: The first-ever compre-  
639 hensive evaluation benchmark of multi-modal llms in  
640 video analysis. In *Proceedings of the Computer*  
641 *Vision and Pattern Recognition Conference*, pages  
642 24108–24118.

643 Bo He, Hengduo Li, Young Kyun Jang, Menglin  
644 Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava,  
645 and Ser-Nam Lim. 2024. Ma-lmm: Memory-  
646 augmented large multimodal model for long-term  
647 video understanding. In *Proceedings of the*  
648 *IEEE/CVF Conference on Computer Vision and Pat-*  
649 *tern Recognition*, pages 13504–13514.

650 Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xi-  
651 angyu Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang  
652 Li, and Limin Wang. 2025. Online video under-  
653 standing: Ovbench and videochat-online. In *Pro-*  
654 *ceedings of the Computer Vision and Pattern Recog-*  
655 *nition Conference*, pages 3328–3338.

656 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam  
657 Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,  
658 Akila Welihinda, Alan Hayes, Alec Radford, and 1  
659 others. 2024. Gpt-4o system card. *arXiv preprint*  
660 *arXiv:2410.21276*.

661 Kumara Kahatapitiya, Kanchana Ranasinghe, Jong-  
662 woo Park, and Michael S Ryoo. 2025. Language  
663 repository for long video understanding. In *Find-*  
664 *ings of the Association for Computational Linguis-*  
665 *tics: ACL 2025*, pages 5627–5646.

666 Minsoo Kim, Kyuhong Shim, Jungwook Choi,  
667 and Simyung Chang. 2025. Infinipot-v:  
668 Memory-constrained kv cache compression for  
669 streaming video understanding. *arXiv preprint*  
670 *arXiv:2506.15745*.

671 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio  
672 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-  
673 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-  
674 täschel, and 1 others. 2020. Retrieval-augmented  
675 generation for knowledge-intensive nlp tasks. *Ad-*  
676 *vances in neural information processing systems*,  
677 33:9459–9474.

678 Ruanjun Li, Yuedong Tan, Yuanming Shi, and Jiawei  
679 Shao. 2025. Videoscan: Enabling efficient stream-  
680 ing video understanding via frame-level semantic  
681 carriers. *arXiv preprint arXiv:2503.09387*.

682 Lin Long, Yichen He, Wentao Ye, Yiyuan Pan, Yuan  
683 Lin, Hang Li, Junbo Zhao, and Wei Li. 2025.  
684 Seeing, listening, remembering, and reasoning: A  
685 multimodal agent with long-term memory. *arXiv*  
686 *preprint arXiv:2508.09736*.

687 Zhicong Lu, Haijun Xia, Seongkook Heo, and Daniel  
688 Wigdor. 2018. You watch, you give, and you en-  
689 gage: a study of live streaming practices in china.  
690 In *Proceedings of the 2018 CHI conference on hu-*  
691 *man factors in computing systems*, pages 1–13.

692	Zhenyu Ning, Guangda Liu, Qihao Jin, Wenchao Ding, Minyi Guo, and Jieru Zhao. 2025. Livevlm: Efficient online video understanding via streaming-oriented kv cache and retrieval. <i>arXiv preprint arXiv:2505.15269</i> .	Haomiao Xiong, Zongxin Yang, Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Jiawen Zhu, and Huchuan Lu. 2025. Streaming video understanding and multi-round interaction with memory-enhanced knowledge. <i>arXiv preprint arXiv:2501.13468</i> .	748 749 750 751 752
697	Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. 2023. Memgpt: Towards llms as operating systems. <i>arXiv preprint arXiv:2310.08560</i> .	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, and 1 others. 2025a. Qwen2. 5-omni technical report. <i>arXiv preprint arXiv:2503.20215</i> .	753 754 755 756
701	Rui Qian, Shuangrui Ding, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. 2025. Dispider: Enabling video llms with active real-time interaction via disentangled perception, decision, and reaction. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 24045–24055.	Ruyi Xu, Guangxuan Xiao, Yukang Chen, Liuning He, Kelly Peng, Yao Lu, and Song Han. 2025b. Streamingvlm: Real-time understanding for infinite video streams. <i>arXiv preprint arXiv:2510.09608</i> .	757 758 759 760
708	Aravinda S Rao, Marko Radanovic, Yuguang Liu, Songbo Hu, Yihai Fang, Kourosh Khoshelham, Marimuthu Palaniswami, and Tuan Ngo. 2022. Real-time monitoring of construction sites: Sensors, methods, and applications. <i>Automation in construction</i> , 136:104099.	Haolin Yang, Feilong Tang, Lingxiao Zhao, Xiang An, Ming Hu, Huifa Li, Xinlin Zhuang, Yifan Lu, Xiaofeng Zhang, Abdalla Swikir, and 1 others. 2025a. Streamagent: Towards anticipatory agents for streaming video understanding. <i>arXiv preprint arXiv:2508.01875</i> .	761 762 763 764 765 766
714	Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, and 1 others. 2024. Moviechat: From dense token to sparse memory for long video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 18221–18232.	Yanlai Yang, Zhuokai Zhao, Satya Narayan Shukla, Aashu Singh, Shlok Kumar Mishra, Lizhu Zhang, and Mengye Ren. 2025b. Streammem: Query-agnostic kv cache memory for streaming video understanding. <i>arXiv preprint arXiv:2508.15717</i> .	767 768 769 770 771
721	Xi Tang, Jihao Qiu, Lingxi Xie, Yunjie Tian, Jianbin Jiao, and Qixiang Ye. 2025. Adaptive keyframe sampling for long video understanding. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 29118–29128.	Linli Yao, Yicheng Li, Yuancheng Wei, Lei Li, Shuhuai Ren, Yuanxin Liu, Kun Ouyang, Lean Wang, Shicheng Li, Sida Li, and 1 others. 2025a. Timechat-online: 80% visual tokens are naturally redundant in streaming videos. In <i>Proceedings of the 33rd ACM International Conference on Multimedia</i> , pages 10807–10816.	772 773 774 775 776 777 778
726	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	Yiqun Yao, Naitong Yu, Xiang Li, Xin Jiang, Xuezhi Fang, Wenjia Ma, Xuying Meng, Jing Li, Aixin Sun, and Yequan Wang. 2025b. Egomem: Lifelong memory agent for full-duplex omnimodal models. <i>arXiv preprint arXiv:2509.11914</i> .	779 780 781 782 783
732	Haibo Wang, Bo Feng, Zhengfeng Lai, Mingze Xu, Shiyu Li, Weifeng Ge, Afshin Dehghan, Meng Cao, and Ping Huang. 2025a. Streambridge: Turning your offline video large language model into a proactive streaming assistant. <i>arXiv preprint arXiv:2505.05467</i> .	Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiying Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, and 1 others. 2025a. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. <i>arXiv preprint arXiv:2507.02259</i> .	784 785 786 787 788 789
738	Mengyue Wang, Shuo Chen, Kristian Kersting, Volker Tresp, and Yunpu Ma. 2025b. Metok: Multi-stage event-based token compression for efficient long video understanding. <i>arXiv preprint arXiv:2506.02850</i> .	Xinlei Yu, Chengming Xu, Guibin Zhang, Zhangquan Chen, Yudong Zhang, Yongbo He, Peng-Tao Jiang, Jiangning Zhang, Xiaobin Hu, and Shuicheng Yan. 2025b. Vismem: Latent vision memory unlocks potential of vision-language models. <i>arXiv preprint arXiv:2511.11007</i> .	790 791 792 793 794 795
743	Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang, Qintong Zhang, Weijia Li, Conghui He, and Linfeng Zhang. 2025. Stop looking for important tokens in multimodal language models: Duplication matters more. <i>arXiv preprint arXiv:2502.11494</i> .	Ekim Yurtsever, Jacob Lambert, Alexander Carballo, and Kazuya Takeda. 2020. A survey of autonomous driving: Common practices and emerging technologies. <i>IEEE access</i> , 8:58443–58469.	796 797 798 799
747		Xiangyu Zeng, Kefan Qiu, Qingyu Zhang, Xinhao Li, Jing Wang, Jiaxin Li, Ziang Yan, Kun Tian, Meng Tian, Xinhai Zhao, and 1 others. 2025.	800 801 802

803 Streamforest: Efficient online video understand-  
804 ing with persistent event memory. *arXiv preprint*  
805 *arXiv:2509.24871*.

806 Guibin Zhang, Muxin Fu, and Shuicheng Yan.  
807 2025. Memgen: Weaving generative latent mem-  
808 ory for self-evolving agents. *arXiv preprint*  
809 *arXiv:2509.24704*.

810 Haoji Zhang, Yiqin Wang, Yansong Tang, Yong  
811 Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin.  
812 2024. Flash-vstream: Memory-based real-time un-  
813 derstanding for long video streams. *arXiv preprint*  
814 *arXiv:2406.08085*.

815 Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim,  
816 Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan  
817 Kian Hsiang Low, and Paul Pu Liang. 2025.  
818 Mem1: Learning to synergize memory and reason-  
819 ing for efficient long-horizon agents. *arXiv preprint*  
820 *arXiv:2506.15841*.

821 Jialong Zuo, Yongtai Deng, Lingdong Kong, Jingkan  
822 Yang, Rui Jin, Yiwei Zhang, Nong Sang, Liang Pan,  
823 Ziwei Liu, and Changxin Gao. 2025. Videolucy:  
824 Deep memory backtracking for long video under-  
825 standing. *arXiv preprint arXiv:2510.12422*.

## Appendix

### A Other Experimental Details 12

#### A.1 Dataset . . . . . 12

#### A.2 Baselines . . . . . 12

##### A.2.1 Closed-source MLLMs . . 12

##### A.2.2 Open-source MLLMs . . . 13

##### A.2.3 Streaming Video Models . 13

#### A.3 Baselines setup . . . . . 14

#### A.4 Closed-source MLLMs . . . . . 14

##### A.4.1 Open-source MLLMs . . . 14

#### A.5 Streaming Video Models . . . . . 14

### B Additional Ablation Study 14

#### B.1 Effect of Memory Decay Coefficient $\lambda$ . . . . . 14

#### B.2 TMR Impact on Video-MME-Long 14

### C Future Works 15

## A Other Experimental Details

### A.1 Dataset

The detailed statistics of the three benchmark datasets are provided in Table 3, and the detailed information of the memory graph is shown in Figure 7. Next, we provide a detailed description of each dataset.

**M3-Bench-robot.** M3-Bench-robot contains 100 real-world videos recorded from a robot’s first-person perspective, covering seven everyday environments, such as living rooms, kitchens, bedrooms, study rooms, offices, meeting rooms, and gyms. Each video depicts interactions between the robot and two to four different humans, where the robot must build and reason based on these interactions. The dataset includes five reasoning types: ME (multi-evidence reasoning), MH (multi-hop reasoning), CM (cross-modal reasoning), PU (person understanding), and GK (general knowledge extraction). M3-Bench-robot aims to evaluate how robots use long-term memory to reason, particularly focusing on complex spatial relationships, object tracking, and dynamic human interactions.

**M3-Bench-web.** M3-Bench-web contains 920 videos sourced from YouTube, covering 46 different video types such as documentaries, street interactions, variety shows, travel, food, and sports

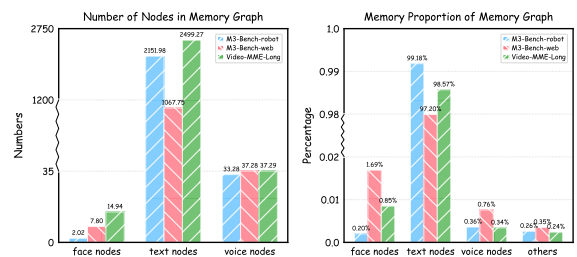


Figure 7: Numbers of different types of nodes and their memory proportion in the memory graphs of the three benchmark datasets.

events. The videos provide a wide range of content, designed to challenge the agent’s ability to reason in complex and dynamic environments. Like M3-Bench-robot, M3-Bench-web also includes five reasoning types, assessing how the agent integrates cross-modal information, understands relationships between people, and extracts general knowledge. This dataset is ideal for evaluating the agent’s memory reasoning capabilities in diverse tasks and complex scenarios, especially for cross-modal reasoning and memory retrieval in varied contexts.

**Video-MME-Long.** Video-MME-Long is the long-video subset of the Video-MME benchmark, consisting of 300 real-world long videos and 900 multiple-choice questions. The videos have an average duration of approximately 41 minutes and may extend up to 1 hour, covering 6 major visual domains and 30 fine-grained categories. The dataset provides multimodal inputs including video frames, audio, and subtitles, while the questions emphasize cross-temporal event association and complex spatiotemporal reasoning. This makes Video-MME-Long a highly challenging benchmark for evaluating multimodal large language models in long-term temporal understanding and multimodal information integration.

### A.2 Baselines

#### A.2.1 Closed-source MLLMs

**Gemini-1.5-Pro.** Gemini 1.5 Pro is an advanced multimodal model capable of handling long-contexts of up to 10 million tokens. It is based on a sparse mixture-of-experts (MoE) Transformer architecture, efficiently processing various data types such as text, video, and audio, and performs exceptionally well across multiple benchmarks.

**GPT-4o.** GPT-4o is an advanced multimodal model capable of processing a variety of input

Dataset	Avg. Duration(s)	#Videos	#QA Pairs	Avg. Text Nodes	#Source
<b>M3-Bench-robot</b>	2039.9	100	1276	2151.98	Real-world filming
<b>M3-Bench-web</b>	1630.7	920	3214	1067.75	YouTube
<b>Video-MME-Long</b>	2466.7	300	900	2499.27	YouTube

Table 3: The statistical information of M3-Bench-robot, M3-Bench-web and Video-MME-Long.

types, including text, images, and videos. It excels in real-time response, efficiently handling inputs and generating multimodal outputs. The system is particularly optimized for visual understanding, enabling fast and accurate interactions in multimodal tasks. This design allows GPT-4o to perform robustly across diverse challenging multimodal benchmarks.

### A.2.2 Open-source MLLMs

**Qwen2.5-VL.** Qwen2.5-VL is a powerful multimodal vision-language model designed for understanding and reasoning across various input types, including images, videos, and text. It excels in tasks such as precise object localization, document parsing, and long-video comprehension, with breakthroughs in video understanding and object localization. Qwen2.5-VL also supports efficient reasoning across diverse domains, handling complex input data and generating accurate multimodal outputs.

**Qwen2.5-Omni.** Qwen2.5-Omni is a unified multimodal model capable of processing various inputs such as text, images, audio, and video, while simultaneously generating text and natural speech responses in a streaming manner. It utilizes an innovative Thinker–Talker architecture, enabling real-time speech responses, and excels in tasks like audio understanding and video reasoning, making it suitable for real-time multimodal interaction scenarios.

### A.2.3 Streaming Video Models

**MovieChat.** MovieChat is a video understanding system that integrates video models and LLM. It uses a sliding window to extract frame-level features and stores them in hybrid memory, with the LLM performing QA based on this memory. Supporting long videos over 10K frames, the system effectively addresses computational complexity and memory overhead, enhancing overall long video comprehension.

**MA-LMM.** MA-LMM is a memory-augmented large multimodal model designed for long-term

video understanding. It processes video frames sequentially and stores information in a long-term memory bank, avoiding the context length and GPU memory limitations faced by traditional models when handling long videos. The model efficiently captures temporal information in videos and achieves excellent performance in tasks such as video question answering and video captioning.

**VideoChat-Online.** VideoChat-Online is an efficient video LLM designed for streaming video understanding, utilizing a Pyramid Memory Bank (PMB) architecture to balance spatial and temporal details. The model is optimized through an offline-to-online learning paradigm, combined with an interleaved dialogue format for video data, enabling real-time video stream processing and outstanding performance across multiple benchmarks. It outperforms many existing models in online video tasks and demonstrates excellent cross-platform capabilities.

**Flash-VStream.** Flash-VStream is a real-time video stream understanding model that employs the STAR (Spatial-Temporal-Abstract-Retrieved) memory mechanism, enabling efficient processing of extremely long video streams and real-time user query responses. By compressing visual information and dynamically updating memory, the model significantly reduces inference latency and GPU memory consumption, supporting online video question answering. Flash-VStream outperforms existing methods in multiple benchmarks, especially in streaming video scenarios.

**TimeChat-Online.** TimeChat-Online is an advanced online VideoLLM designed for real-time video interaction. Its core innovation lies in the Differential Token Drop (DTD) module, which reduces visual redundancy in streaming videos by selectively preserving significant temporal changes while discarding static content between frames. This approach helps achieve an 82.8% reduction in video tokens, maintaining high performance in real-time video question answer-

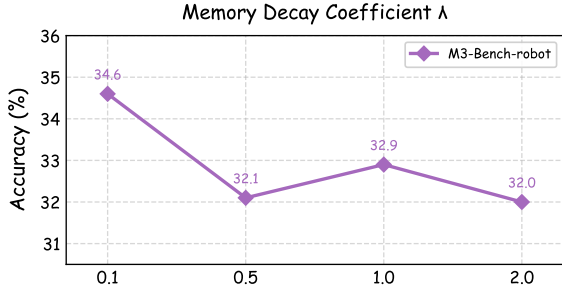


Figure 8: Impact of memory decay coefficient  $\lambda$  on M3-Bench-robot.

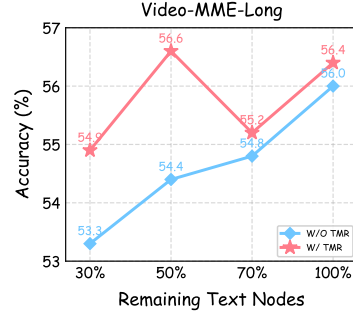


Figure 9: Impact of the TMR Mechanism on Video-MME-Long at different compression ratios.

ing (QA). Additionally, it integrates proactive responding capabilities, triggering answers based on future visual cues from video scene transitions.

**StreamForest.** StreamForest is an online video understanding model specifically designed for streaming video scenarios. Its core innovation lies in the Persistent Event Memory Forest, which dynamically organizes video frames into event-level tree structures, ensuring effective storage and processing of long-term visual information. The model also introduces a Fine-grained Spatiotemporal Window, enhancing real-time scene perception. It performs exceptionally well across multiple streaming video understanding benchmarks, particularly in real-time video question answering and autonomous driving applications.

**M3-Agent.** M3-Agent is the first industrial-grade streaming video agent, integrated with a long-term Agent Memory mechanism, capable of processing video and audio inputs in real-time. It generates episodic and semantic memories to gradually build world knowledge and completes tasks through multi-turn reasoning.

### A.3 Baselines setup

### A.4 Closed-source MLLMs

M3-Agent consists of two models, memorization and control, and is trained using reinforcement learning. Following the approach in the M3-Agent paper, we use the closed-source MLLMs Gemini-1.5-Pro and GPT-4o to replace the memorization model, generating descriptions for video clips and storing them as long-term memory. Subsequently, through retrieval-augmented generation (RAG) (Lewis et al., 2020), we use GPT-4o to perform memory retrieval as a replacement for the control model.

Following the approach in the M3-Agent paper, for Gemini-Agent, we use Gemini-1.5-Pro

to replace the memorization and control models in M3-Agent. For the generated memory graph, we use M3-Agent’s search function to replace RAG. Similarly, for Gemini-GPT4o-Hybrid, we use Gemini-1.5-Pro to generate memory as the memory model, use GPT-4o for control, and still use M3-Agent’s search function to replace RAG.

### A.4.1 Open-source MLLMs

Following the approach in the M3-Agent paper, we use the open-source MLLMs Qwen2.5-VL and Qwen2.5-Omni to replace the closed-source MLLMs Gemini-1.5-Pro and GPT-4o in Section A.4, to ensure a fair comparison, with the rest remaining consistent as before.

### A.5 Streaming Video Models

For the streaming video models, we use the pre-trained models and default parameter settings from the official papers.

## B Additional Ablation Study

### B.1 Effect of Memory Decay Coefficient $\lambda$ .

As shown in Figure 8, the model achieves its best performance when the memory decay coefficient is set to memory decay coefficient  $\lambda = 0.1$ . As  $\lambda$  increases further, the performance drops to varying degrees. This indicates that an excessively large decay coefficient leads to an overly strong temporal decay effect, causing the model to diminish the importance of historical memory too early. Consequently, the model fails to effectively leverage long-term contextual information, ultimately degrading its reasoning performance.

### B.2 TMR Impact on Video-MME-Long

As shown in Figure 9, consistent with the trends observed in M3-Bench-robot and M3-Bench-web, the model’s performance on the Video-MME-Long dataset significantly decreases when the

1066 compression ratio exceeds 50%. After introduc-  
1067 ing the TMR mechanism, the performance degrada-  
1068 tion is effectively mitigated, and even the per-  
1069 formance at a 50% compression ratio surpasses  
1070 that at 30%. These results are in strong agreement  
1071 with the previous experiments, indicating that the  
1072 TMR mechanism can significantly reduce the ac-  
1073 curacy loss caused by memory graph compression  
1074 and retrieve more critical memories for reasoning.

## 1075 **C Future Works**

1076 Future work can further extend the ideas pro-  
1077 posed in this study. First, this work adopts a  
1078 unified compression ratio for isolated text nodes  
1079 and connected text nodes; however, future re-  
1080 search could explore differentiated and adaptive  
1081 compression strategies for different node types,  
1082 so as to better exploit their structural character-  
1083 istics and information distributions. Second, fu-  
1084 ture studies may investigate the differential com-  
1085 pression of episodic nodes and semantic nodes,  
1086 where episodic nodes primarily focus on specific  
1087 events occurring in video segments, while seman-  
1088 tic nodes emphasize general and relatively stable  
1089 knowledge distilled from these segments. Effec-  
1090 tively preserving highly discriminative event-level  
1091 descriptions and critical stable semantic memories  
1092 during compression is an important direction for  
1093 further enhancing the expressive power and rea-  
1094 soning support of the memory graph. In addi-  
1095 tion, we observe a noticeable redundancy among  
1096 text nodes in the M3-Bench-robot dataset. Fu-  
1097 ture work may systematically study the identifi-  
1098 cation and resolution of redundant or even con-  
1099 flicting memories in the memory graph, as well  
1100 as potential memory poisoning issues. By selec-  
1101 tively removing or correcting unreliable memo-  
1102 ries, the robustness and reliability of the model in  
1103 complex streaming video scenarios can be further  
1104 improved. Finally, reducing the time overhead  
1105 required for memory graph construction in the  
1106 early stage is also an important research direction.  
1107 Future work may explore more efficient memory  
1108 construction mechanisms to lower the cost of gen-  
1109 erating memory graphs.