# WorldMM: Dynamic Multimodal Memory Agent for Long Video Reasoning

Woongyeong Yeo[1*]    Kangsan Kim[1*]    Jaehong Yoon[2†]    Sung Ju Hwang[1,3†]

KAIST[1]    Nanyang Technological University[2]    DeepAuto.ai[3]

https://worldmm.github.io

## Abstract

*Recent advances in video large language models have demonstrated strong capabilities in understanding short clips. However, scaling them to hours- or days-long videos remains highly challenging due to limited context capacity and the loss of critical visual details during abstraction. Existing memory-augmented methods mitigate this by leveraging textual summaries of video segments, yet they heavily rely on text and fail to utilize visual evidence when reasoning over complex scenes. Moreover, retrieving from fixed temporal scales further limits their flexibility in capturing events that span variable durations. To address this, we introduce WorldMM, a novel multimodal memory agent that constructs and retrieves from multiple complementary memories, encompassing both textual and visual representations. WorldMM comprises three types of memory: episodic memory indexes factual events across multiple temporal scales, semantic memory continuously updates high-level conceptual knowledge, and visual memory preserves detailed information about scenes. During inference, an adaptive retrieval agent iteratively selects the most relevant memory source and leverages multiple temporal granularities based on the query, continuing until it determines that sufficient information has been gathered. WorldMM significantly outperforms existing baselines across five long video question-answering benchmarks, achieving an average 8.4% performance gain over previous state-of-the-art methods, showing its effectiveness on long video reasoning.*

## 1. Introduction

With the increasing deployment of video large language models (video LLMs) [1, 3, 16, 37] in real-world applications, such as AI glasses and household robots, these models are now required to process and reason over extremely long videos from several hours to even days [5, 25, 33]. Recent works [4, 12, 13] have introduced memory-based

approaches that build external memories from abstracted video representations. Such methods allow the model to focus on essential information by retrieving a small number of relevant memories, thereby reducing the number of input tokens. This is a more efficient and effective strategy compared to processing all frames in the video, requiring high computational cost as illustrated in Fig. 1(a).

Despite their promise, most existing approaches remain highly dependent on textual representations. Typically, each detected event or clip is converted into captions, summaries, or structured text descriptions for downstream retrieval and reasoning [12, 23, 33]. Although Long et al. [13] incorporates visual inputs when building memory, its use of multimodal features is limited to entity recognition and is not fully exploited during inference (Fig. 1(b)). Moreover, existing models [13, 33] typically retrieve a fixed number of clips with predetermined durations, such as retrieving three 30-second clips. These rigid architecture designs in video memory agents face the following two major limitations.

First, they fail to adaptively leverage visual information from videos in conjunction with textual memory during retrieval and generation. Visual details are essential for many real-world tasks requiring attribute recognition, spatial reasoning, or precise scene understanding, while this knowledge cannot be fully represented in text. Meanwhile, as shown in Fig. 1(c), a fixed strategy that always includes both captions and frames during response generation yields suboptimal results since excessive visual context can even distract the model. Therefore, an adaptive mechanism for selecting multimodal memories is essential for retrieving the most informative references for a given query, which remains unexplored in previous works.

Second, retrieving a fixed number of clips limits the model's ability to handle queries that require varying temporal scopes. For instance, a question like *"Where did I leave my glasses?"* may require only a few seconds of video, whereas *"What happened in the second half of the soccer match?"* demands a much longer temporal context. Existing approaches retrieve a predetermined length of segments for simplicity [12, 13, 33], which inherently over-
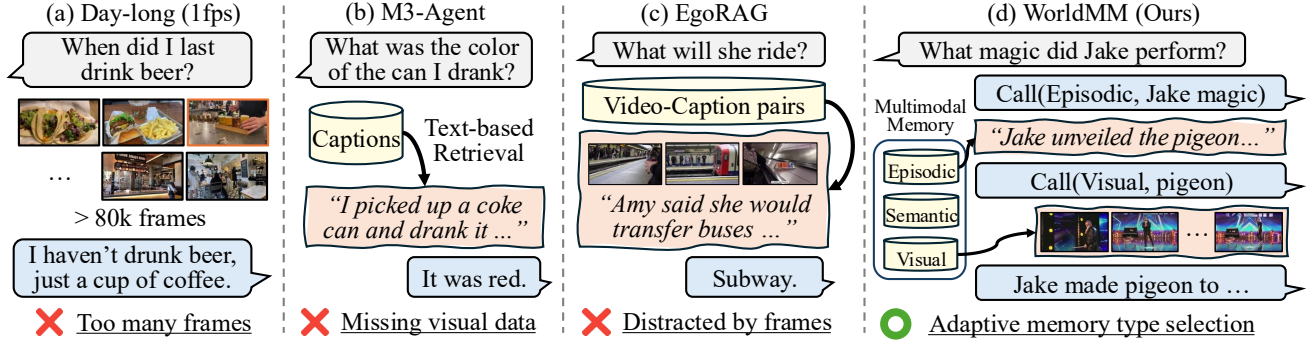
---

*Equal contribution; †Equal advising

**Figure 1. Concept Figure.** (a) A day-long video sampled at 1 fps has frames that exceed the context limits of video LLMs. (b) M3-Agent [13] relies on textual representation of video, which can underrepresent visual information. (c) EgoRAG [33] retrieves both captions and the corresponding visual frames, but irrelevant frames may distract model. (d) WorldMM (Ours) constructs multiple memories, incorporating both textual and visual representations, and uses adaptive memory retrieval to effectively leverage multimodal information.

looks the diverse temporal scales of real-world events, limiting flexibility. Instead, the retriever should dynamically gather information at multiple temporal scales, combining hour-level summaries with minute-level details as needed.

To fill this gap, we present WorldMM, a novel memory-based agent that constructs separate textual and visual memories and employs an adaptive retrieval agent to select the optimal memory modality and temporal granularity for each query. The textual memory comprises two components: episodic memory, which stores multiple events across different time scales, and semantic memory, which captures high-level, long-term knowledge such as relationships and habits, organized within knowledge graphs. The visual memory divides a long video into short-term segments indexed within a retrieval corpus, enabling the model to access visual information when required. The retrieval agent iteratively selects the most relevant memory across modalities and timescales, ensuring that the agent retrieves only the information necessary for each query. The proposed multimodal memory selection design therefore prevents the model from being forced to condition on paired yet unnecessary modality memories when retrieving data for a given query, minimizing potential distraction during reasoning.

In addition, WorldMM is able to retrieve information at varying levels of granularity over the appropriate time range by leveraging multiple graphs operating at different temporal scales, such as seconds, minutes, and hours. When episodic memory is selected for retrieval, the retrieval agent searches each memory to gather potentially relevant information from all temporal levels. The collected candidates are then jointly examined to determine which pieces of information should be used to answer the query. In the end, we dynamically access both short- and long-term video contexts to assemble only the necessary information for reasoning. Furthermore, the model performs retrieval in multiple turns by iteratively selecting memories and queries, thereby expanding the range of possible combinations and allowing

adaptive selection of the information for each query.

We evaluate WorldMM with five long video question-answering benchmarks from hour- to week-long durations. The proposed approach consistently outperforms strong baselines, including long video LLMs and memory-augmented models. Comprehensive ablation studies further demonstrate the effectiveness of our multi-memory, multi-scale design. Specifically, episodic memory enables reasoning over events at multiple timescales, visual memory improves performance on object- and action-centered queries, while semantic memory enhances reasoning over long-term contexts. When all the memories are adaptively integrated, the model achieves the best overall results. These results highlight that our multimodal memory system represents a promising direction toward robust long video reasoning.

## 2. Related Work

### 2.1. Long Video Understanding

Existing video LLMs demonstrate strong understanding capabilities for short videos, and recent research has shifted toward reasoning over longer videos. Current proprietary models, such as GPT-5 [16] and Gemini 2.5 [3], have advanced to minute- or hour-level video understanding [2, 5, 12, 25, 34] by utilizing extended context lengths. However, these models still incur high computational costs and uniformly sampling every frame is often suboptimal for questions focused on localized events [17, 24].

To address these challenges, several strategies have been explored. Visual token compression [8, 9, 11, 19, 20, 26, 35] improves efficiency by reducing token counts but often loses fine-grained details, limiting the capture of subtle or sparse events. Key frame selection [22, 30] retains only the most informative frames to reduce redundancy but fails to detect relevant frames when video streams are too long and may miss rare events. More recently, reasoning-centric training and inference [28, 29] have enhanced long-
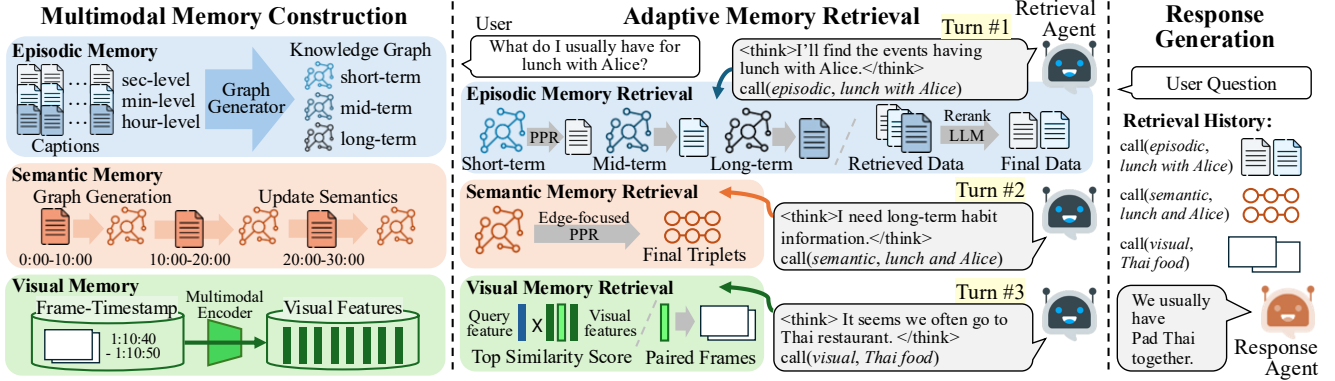
Figure 2. **Overview of WorldMM.** (Left) Multimodal Memory Construction: WorldMM builds three complementary memories (episodic, semantic, and visual memory) that capture temporal events, long-term relations, and visual details from video streams. (Middle) Adaptive Memory Retrieval: A retrieval agent iteratively selects and integrates relevant information from diverse memories for a given query. (Right) Response Generation: The retrieved content and reasoning history are used by a response agent to produce a grounded response.

range temporal grounding through reinforcement learning and adaptive test-time scaling, yet still face scalability limits on ultra-long videos over ten hours.

Beyond hour-level videos, emerging benchmarks push video understanding and reasoning toward day- or even week-long continuous recordings [23, 33]. The aforementioned strategies struggle in these settings due to the massive scale of frames and long-term temporal dependencies. This highlights the necessity for more efficient, context-aware, and scalable approaches to handle extremely long videos.

## 2.2. Memory-based Video LLMs

In order to effectively reason over long videos, retrieval-augmented generation (RAG) based methods that retrieve relevant frames or clips instead of sampled frames have been introduced. They typically retrieve query-relevant information using textual and visual cues, and allow the model to focus on crucial clips [10, 14, 31]. Some recent methods extend this approach by constructing graph structures to encode multimodal interactions across frames [18, 21, 32]. However, these models rely on textual representation or a simple similarity score between visual and query features, limiting their ability to perform holistic understanding and complex reasoning over long video sequences.

Beyond naive retrieval-based design, memory-based methods have emerged to construct structured knowledge over video streams. EgoRAG [33] organizes hierarchical textual memories that store events from egocentric video streams in a hierarchical manner, allowing reasoning throughout day-long activities. Ego-R1 [23] extends this by leveraging vision-centric tools with iterative reasoning to perform long-horizon reasoning. HippoMM [12] proposes a dual-process memory using semantic summaries with multimodal cues. M3-Agent [13] constructs entity-centric long-term memory by processing multimodal contexts and adopts iterative reasoning to retrieve relevant knowledge

from the memory. Despite these advances, existing works still struggle to fully integrate multimodal information and to dynamically retrieve knowledge across varying temporal scales to handle complex, long video scenarios.

## 3. WorldMM

We introduce WorldMM, a novel framework that leverages both textual and visual contexts of video streams to build a multimodal memory for comprehensive understanding and reasoning over long videos. As illustrated in Fig. 2, the model operates in three stages: multimodal memory construction, adaptive memory retrieval, and response generation. Given a long video stream, WorldMM first builds multiple memories, including two textual memories and one visual memory (Sec. 3.1). Next, a retrieval agent iteratively collects query-relevant information from different memories and timescales until sufficient knowledge is gathered to answer the question (Sec. 3.2). Finally, the query is fed into a response agent along with the retrieved contents and retrieval history to generate a response (Sec. 3.3).

## 3.1. Multimodal Memory Construction

As we described in Sec. 1, an effective memory agent for long-form video understanding must address two key requirements: 1) *adaptively leveraging visual information alongside text memory*, and 2) *retrieving knowledge across diverse temporal ranges*. To achieve this, WorldMM constructs three types of memory, each encoding complementary video knowledge across diverse modalities. Episodic memory captures diverse events over multiple dynamic timescales, semantic memory incrementally updates high-level relational knowledge, and visual memory preserves spatial and appearance details. Together, they form a comprehensive multimodal memory that supports episodic retrieval, semantic reasoning, and visually grounded understanding of long-form videos.

**Episodic Memory Construction** Episodic memory consists of multiple textual graphs, each of which encodes events at different temporal resolutions. Before constructing the graphs, we first perform fine-grained captioning on the unit temporal scale $t_0$. We divide the video into short segments of length $t_0$, each converted into a caption using a video LLM. Most existing approaches rely on a fixed temporal scale during memory construction [13, 33], overlooking the diverse spans of real-world events. In contrast, we introduce a multi-scale memory composed of multiple temporal resolutions that flexibly encodes information with different levels of density:

$$\mathcal{T} = \{t_0, t_1, \ldots, t_N\}, \quad t_0 < t_1 < \cdots < t_N. \quad (1)$$

For each temporal scale $t_i \in \mathcal{T}$, the video is partitioned into non-overlapping segments of length $t_i$. The segments are captioned and transformed into factual triplets (entity-action-entity) to construct a knowledge graph (KG) $G_{t_i}$. Finally, episodic memory is represented as a set of KGs:

$$\mathcal{M}_e = \{G_{t_0}, G_{t_1}, \ldots, G_{t_N}\}. \quad (2)$$

This multi-scale episodic memory enables temporally grounded reasoning that spans both fine-grained event details and long-range narrative understanding.

**Semantic Memory Construction** Semantic memory captures long-term, evolving knowledge about relationships and habits within a video. Since episodic graphs are constructed from independent events, they fail to preserve continuity across distant scenes and cannot capture high-level knowledge. Semantic memory, on the contrary, maintains an evolving graph that continuously integrates relational and habitual knowledge over time.

To build this continually updating memory, we first split the input video into coarse segments with a fixed timescale $t_s$. Textual captions are generated for each segment and converted into semantic triplets $T_{t_s}^k$, focusing on conceptual knowledge rather than event-specific details. These triplets are incrementally integrated into an evolving semantic graph through a consolidation process that merges new knowledge while preserving stable relationships. Specifically, embedding-based similarity is first used to identify overlapping or conflicting triplets between the current graph $G_{t_s}^k$ and the newly extracted triplets $T_{t_s}^{k+1}$. The matched triplets are then provided to an LLM along with the new triplets, which determines outdated or conflicting triplets $T_{\text{remove}}$ and triplets that should be revised or added $T_{\text{update}}$. Formally, the consolidation process can be represented as:

$$\text{Consolidate}(G_{t_s}^k, T_{t_s}^{k+1}) = (G_{t_s}^k \setminus T_{\text{remove}}) \cup T_{\text{update}}. \quad (3)$$

The resulting semantic memory is a continuously evolving KG $\mathcal{M}_s = G_{t_s}^M$, where $M$ denotes the final segment index, capturing the video's long-term knowledge.

**Visual Memory Construction** Visual memory captures rich visual details that text cannot fully convey, including detailed object appearances, scene dynamics, and precise spatial context. WorldMM explicitly constructs a visual memory to ground reasoning in visual evidence. We consider two scenarios in which visual memory is invoked, when the retrieval agent searches for scenes associated with a specific keyword, and when the agent has timestamps identified during preceding retrieval steps to inspect the corresponding frames. Therefore, we adopt two complementary strategies for building visual memory: feature-based retrieval via natural language query, and timestamp-based retrieval for precise temporal grounding.

Specifically, we partition each video into short, fixed-length segments of duration $t_v$, encoding each segment $V_{t_v}^k$ into a visual feature $f_v^k$ using a multimodal encoder, forming a feature-based visual memory as a set of embeddings:

$$\mathcal{M}_v^f = \{f_v^1, f_v^2, \ldots, f_v^L\}. \quad (4)$$

In parallel, to support timestamp-based retrieval, each frame is paired with its corresponding timestamp:

$$\mathcal{M}_v^I = \{(t_i, I_i) \mid I_i = V(t_i), \ t_i \in [0, \text{len}(V)]\}. \quad (5)$$

This allows direct access to visual evidence at specific moments in the video. Finally, the complete visual memory integrates both components $\mathcal{M}_v = \mathcal{M}_v^f \cup \mathcal{M}_v^I$ by combining feature-level embeddings and frame-level indices.

### 3.2. Adaptive Memory Retrieval

In this section, we present how WorldMM dynamically retrieves the most relevant multimodal memories from the appropriate temporal scope for a given query.

**Retrieval Agent** Reasoning over long-form videos requires integrating heterogeneous information from multiple memory sources. To handle this, the retrieval agent iteratively decides which memory to access and what query to issue, conditioned on the user question and retrieval history. Leveraging the distinct characteristics of each memory component, it adaptively selects the most relevant source and formulates modality-specific queries. Through successive iterations, the agent progressively refines its retrieval strategy and constructs better knowledge collection.

Formally, we define the retrieval agent $\mathcal{R}$ as a multimodal reasoning module that iteratively selects a memory source and formulates a corresponding query. At each iteration $i$, $\mathcal{R}$ takes an input the user query $q$ and the set of previous retrieval histories $r_{<i} = \{r_1, \ldots, r_{i-1}\}$, and outputs either a memory–query pair or a STOP signal:

$$\mathcal{R}(q, r_{<i}) = \begin{cases} (m_i, q_i) & \text{if } r_{<i} \text{ insufficient and } i \leq N, \\ \text{STOP} & \text{otherwise,} \end{cases} \quad (6)$$

where $m_i \in \{\mathcal{M}_e, \mathcal{M}_s, \mathcal{M}_v\}$ and $N$ denotes the maximum number of iterations. If the retriever outputs a memory–query pair $(m_i, q_i)$, it retrieves the relevant information from the memory $m_i$ with search query $q_i$ and proceeds to the next iteration with the updated context $r_{\leq i}$. When the retriever outputs STOP, it indicates that sufficient information has been collected. The iterative process then terminates, and all retrieved results $\{r_1, \ldots, r_n\}$ are passed to the response agent for the final response generation.

**Episodic Memory Retrieval**  Episodic memory retrieval is guided by a query $q$ provided by the retriever, which specifies the desired information from episodic memory. The main challenge lies in determining the appropriate temporal scope, as episodic memory contains multiple graphs covering different temporal ranges. WorldMM adopts a coarse-to-fine, multi-timescale retrieval strategy. Specifically, for each temporal graph $G_{t_i}$, the model first retrieves top-$k$ candidate captions using a graph-based retrieval framework guided by the Personalized PageRank (PPR) score and the query, following Gutiérrez et al. [7]. Subsequently, an LLM serves as a cross-scale reranker, jointly analyzing the query and retrieved candidates across all timescales. It then selects the most relevant temporal range and refines the retrieved content, producing the final top-$m$ captions as output. By retrieving from multi-scale memory, the model leverages both coarse temporal context and fine-grained details.

**Semantic Memory Retrieval**  The semantic memory, also represented as a graph, is queried using a PPR-based retrieval algorithm. In contrast to episodic memory retrieval which operates over nodes and their temporal structures, semantic retrieval focuses on relational knowledge encoded as edges between entities. Since the standard PPR score measures node-level relevance, we adapt it for edge-based reasoning by assigning each edge a score equal to the sum of the PPR values of its two connected nodes. The top-$k$ triplets corresponding to the highest-scoring edges are then selected as the final retrieved results.

**Visual Memory Retrieval**  Following Sec. 3.1, visual memory retrieval operates in two complementary modes: feature-based search and timestamp-based access. In feature-based mode, the retrieval agent formulates a query $q$, encodes it into a text feature $f_t$ using a multimodal encoder, and retrieves the top-$k$ relevant video segments from $\mathcal{M}_v^f$ based on the cosine similarity between $f_t$ and the visual features. In timestamp-based mode, when specific temporal ranges are identified, typically following episodic retrieval, the corresponding frames are directly fetched from $\mathcal{M}_v^I$. By combining these two modes, WorldMM enables flexible and effective access to visual information at both semantic and temporal levels.

### 3.3. Response Generation

Finally, once the retrieval agent determines that sufficient information has been gathered, the retrieval process is terminated. The retrieval history, including the selected memories, their corresponding queries, and the retrieved results, is then passed to the response agent along with the original user query. The response agent generates the final answer by grounding its response in the retrieved information. This clear separation between the retriever and the responder allows each component to focus on its respective objective, ensuring effective retrieval and response generation.

## 4. Experiment Results

### 4.1. Experimental Setup

**Datasets and Metrics**  We assess the performance of WorldMM across five benchmarks that require reasoning over long videos. EgoLifeQA [33] and Ego-R1 Bench [23] contain week-long videos, and HippoVlog [12] features vlog-style content, requiring comprehension of audio and visual streams. We also assess general video understanding on LVBench [25] and Video-MME (long) [5] with hour-level videos. All benchmarks consist of multiple-choice questions, with accuracy used as the evaluation metric. Please see additional dataset details in the Sec. A.

**Baselines**  We compare WorldMM against a comprehensive set of baselines spanning base video LLMs, long video understanding models, RAG systems, and memory-based models. Base video LLMs include GPT-5 [16], Gemini 2.5 Pro [3], and Qwen3-VL-8B-Instruct [1], while long video understanding models include VideoChat-Flash [11], Time-R1 [28], and Video-RTS [29], which all use uniformly sampled frames within their input capacity. We further evaluate RAG approaches, including text retrieval methods like LightRAG [6] and HippoRAG [7], which retrieve video captions, and Video-RAG [14], which retrieves relevant clips. Finally, we compare with memory-based frameworks for long video reasoning, including EgoRAG [33], Ego-R1 [23], HippoMM [12], and M3-Agent [13].

**Implementations Details**  We adopt VLM2Vec-V2 [15] as a multimodal encoder for visual memory retrieval. During the memory construction, GPT-5-mini is used for building episodic and semantic memories. We experiment ours with two video LLMs, GPT-5 and Qwen3-VL-8B-Instruct, which serve as the retrieval and response agent, respectively denoted as WorldMM-GPT and WorldMM-8B. For temporal segmentation in episodic memory, we apply timescales specific to each dataset. For example, we use 30-second, 3-minute, 10-minute, and 1-hour intervals for EgoLifeQA. Configurations for other benchmarks, more experimental details, and the prompts are provided in Sec. B.

Table 1. Performance of WorldMM with various baselines across long video QA benchmarks. "–" denotes a proprietary backbone.

| Model | | EgoLife QA | Ego-R1 Bench | Hippo Vlog | LV Bench | Video-MME (L) | Avg. |
|---|---|---|---|---|---|---|---|
| ***Base Models*** | | | | | | | |
| Qwen3-VL-8B [1] | 8B | 38.6 | 35.7 | 74.4 | 48.3 | 61.0 | 51.6 |
| Gemini 2.5 Pro [3] | – | 46.4 | 46.7 | 72.0 | 57.0 | 55.7 | 55.6 |
| GPT-5 [16] | – | 48.6 | 46.3 | 75.7 | 60.4 | 74.3 | 61.1 |
| ***Long Video LLMs*** | | | | | | | |
| VideoChat-Flash [11] | 7B | 34.2 | 42.7 | 58.0 | 33.2 | 44.1 | 42.4 |
| Time-R1 [28] | 3B | 48.8 | 48.0 | 54.6 | 31.1 | 37.6 | 44.0 |
| Video-RTS [29] | 7B | 48.2 | 47.4 | 59.0 | 39.8 | 47.9 | 48.6 |
| ***RAG-based Video LLMs*** | | | | | | | |
| LightRAG [6] | – | 48.8 | 52.3 | 47.4 | 30.4 | 46.6 | 45.1 |
| HippoRAG [7] | – | 59.6 | 56.0 | 63.2 | 54.0 | 52.1 | 57.0 |
| Video-RAG [14] | – | 55.4 | 49.7 | 65.1 | 33.1 | 55.4 | 51.7 |
| ***Memory-based Video LLMs*** | | | | | | | |
| EgoRAG [33] | – | 52.0 | 49.0 | 57.5 | 32.2 | 41.1 | 46.4 |
| Ego-R1 [23] | 3B | 53.0 | 52.0 | 58.8 | 34.1 | 42.7 | 48.1 |
| HippoMM [12] | – | 54.6 | 53.0 | 71.9 | 38.2 | 41.6 | 51.8 |
| M3-Agent [13] | 7B | 53.5 | 52.0 | 65.5 | 49.3 | 55.3 | 55.1 |
| ***WorldMM (Ours)*** | | | | | | | |
| **WorldMM-8B** | 8B | 56.4 | 52.0 | 69.7 | 55.4 | 66.0 | 59.9 |
| **WorldMM-GPT** | – | **65.6** | **65.3** | **78.3** | **61.9** | **76.6** | **69.5** |

## 4.2. Main Results

Tab. 1 presents the evaluation results of the proposed WorldMM and baseline models. WorldMM significantly outperforms all baselines across various long video understanding benchmarks. In particular, WorldMM-GPT achieves an average score of 69.5%, exceeding the strongest baseline by 8.4%. Compared with base models, both variants of our model surpass their corresponding baselines by more than 8% on average, highlighting the effectiveness of our framework in leveraging strong reasoning capabilities without requiring full video processing. On the other hand, models in the long video LLM category show the weakest performance, with all results falling below 50% on EgoLifeQA and Ego-R1 Bench, indicating that these approaches are not effective in days-long videos.

Meanwhile, retrieval- and memory-based approaches, including ours, achieve scores mostly above 52% on Ego-LifeQA and Ego-R1 Bench, suggesting that selective retrieval of relevant segments is more effective for long video understanding than processing full video sequences. Compared with other retrieval-based models such as HippoRAG and HippoMM, which also rely on GPT backbones, our model achieves markedly higher accuracy on average (69.5% vs. 57.0% and 51.8%). These demonstrate that integrating textual and visual memory and adaptively selecting temporal scopes are crucial for effective video reasoning.

## 4.3. Efficacy of Multimodal Memory

To examine the contribution of each memory in our framework, we perform an ablation study that varies the composition of available memories. The evaluation results, summarized in Tab. 2, show a consistent improvement in performance as additional memories are incorporated. This
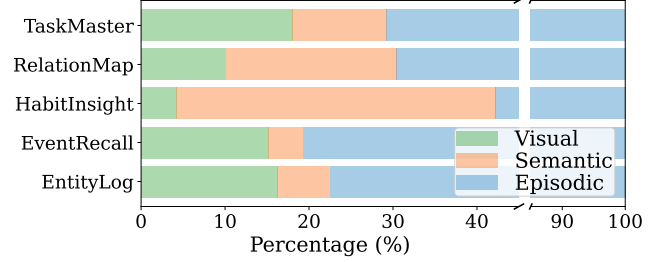


Figure 3. Memory type utilization of WorldMM on five distinctive categories in EgoLifeQA.

finding confirms that different memories capture complementary forms of knowledge. All following experiments in this section are conducted using WorldMM-GPT.

**Effect of Episodic Memory** To examine the performance differences arising from the retrieved data modality in WorldMM, we evaluate models using only episodic memory (E) and only visual memory (V), and report the results in Tab. 2. Using only episodic memory shows 20% higher performance than using only visual memory on average. This is mostly because textual information can be more readily organized into a graph, which enables effective retrieval, while indexing visual frames into a structured representation remains challenging.

**Effect of Visual Memory** Visual memory plays a particularly important role in categories that demand perceptual understanding, such as object recognition or action interpretation. In Tab. 2, visual memory significantly enhances accuracy in categories like EntityLog and EventRecall of EgoLifeQA and Ego-R1 Bench, as well as Visual and Audio+Visual of HippoVlog. The full configuration (E+S+V) surpasses the non-visual configuration (E+S) by an average margin of 4.2%. This improvement arises because visual information preserves spatial and perceptual details that are difficult to represent in text. As shown in Fig. 4(a), when relying solely on episodic memory, the model fails to capture object details such as the type of baked item, leading to an incorrect response. In contrast, visual memory provides access to corresponding frames that contain a complete scene, enabling accurate interpretation of objects and activities.

**Effect of Semantic Memory** Semantic memory proves the most beneficial for categories that require reasoning over long-term dependencies and abstract relationships. This effect is evident in the HabitInsight and Relation-Map categories of EgoLifeQA and Ego-R1 Bench. The model equipped with full memory achieves 76.9% accuracy in HabitInsight, representing a 23% improvement over the setting without semantic memory (E+V). This substantial gain indicates that semantic memory serves as a structured

Table 2. Performance of WorldMM across multiple benchmarks using different memory types. E, S, and V denote episodic, semantic, and visual memories, respectively. Combinations with "+" indicate multiple memory types are used.

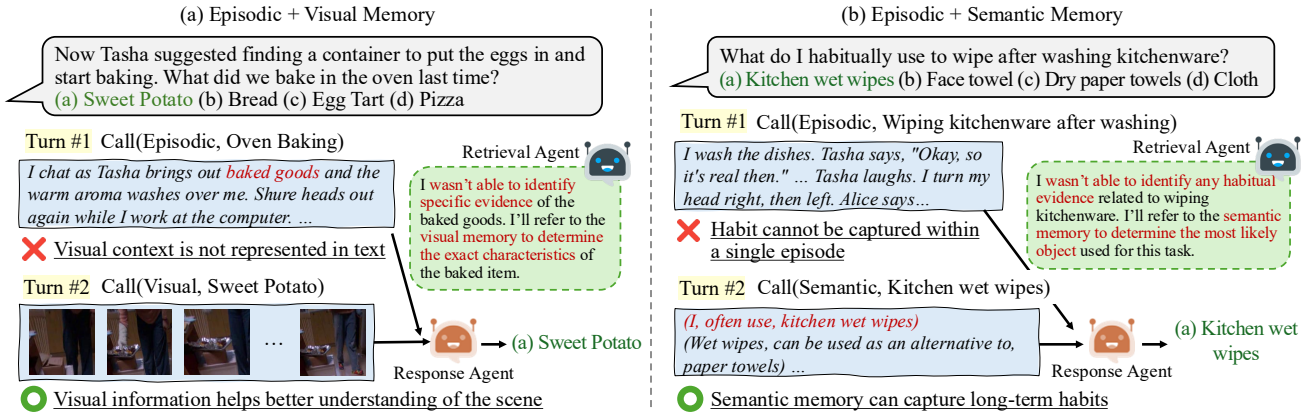| Model | EgoLifeQA | | | | | | Ego-R1 Bench | | | | | | HippoVlog | | | | | LVBench | Video-MME (L) | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Aud. | Vis. | A+V | Summ. | Avg. | | | |
| E | 57.6 | 61.1 | 70.5 | 61.6 | 69.8 | 62.6 | 54.5 | 70.7 | 53.9 | 52.6 | 57.9 | 57.0 | 72.4 | 73.2 | 68.4 | 80.4 | 73.6 | 60.6 | 72.7 | 64.9 |
| V | 40.8 | 35.7 | 36.1 | 34.4 | 39.7 | 37.2 | 36.5 | 34.1 | 23.1 | 31.6 | 28.2 | 34.2 | 35.2 | 66.4 | 54.8 | 48.8 | 51.3 | 47.4 | 64.2 | 44.9 |
| E+S | 56.8 | 61.9 | 73.8 | 62.4 | 71.4 | 63.4 | 59.3 | 68.3 | 69.2 | 57.9 | 60.5 | 61.0 | 70.8 | 75.2 | 68.8 | 80.4 | 73.8 | 58.8 | 74.1 | 66.8 |
| E+V | 59.2 | 63.5 | 70.5 | 60.8 | 68.8 | 63.3 | 65.1 | 68.3 | 53.9 | 47.4 | 57.9 | 63.0 | 73.2 | 77.2 | 70.8 | 79.6 | 75.2 | 59.8 | 76.0 | 66.9 |
| **E+S+V** | **62.4** | **64.3** | **75.4** | **62.4** | **71.4** | **65.6** | **64.6** | **70.7** | **76.9** | **57.9** | **63.2** | **65.3** | **75.6** | **81.6** | **73.2** | **82.8** | **78.3** | **61.9** | **76.6** | **69.5** |



Figure 4. **Qualitative results.** (a) Episodic memory alone cannot capture detailed visual context. The retrieval agent dynamically retrieves from visual memory, enabling access to fine-grained visual details. (b) To address the limitations of episodic memory in representing relationships or habitual behaviors, the retrieval agent proactively accesses semantic memory, allowing it to incorporate habitual knowledge.

knowledge base that captures relational or habitual knowledge accumulated over time. The qualitative example in Fig. 4(b) illustrates this behavior, where episodic memory alone fails to infer habitual actions that extend beyond a single event. By contrast, semantic memory captures the repeated use of kitchen wet wipes, allowing the model to infer the correct answer through long-term reasoning.

**Adaptive Retrieval on Multimodal Memory** To further analyze how categories differ in their reliance on distinct memory modalities in WorldMM, we quantify the usage proportion of each memory across all retrieval iterations per category. As shown in Fig. 3, while episodic memory plays a foundational role across all tasks, certain categories tend to select it more frequently than other memory types: HabitInsight and RelationMap depend primarily on semantic memory, reflecting their reliance on reasoning over long-term patterns. In contrast, EntityLog and EventRecall benefit more from visual memory, which provides fine-grained perceptual details not fully captured by text. This selective utilization suggests that the model dynamically emphasizes the most relevant memory type for each category, leveraging the strength of each type in a context-dependent manner. These results confirm that different types of memory contribute distinct yet complementary strengths.

### 4.4. Dynamic Temporal Scope Retrieval

We evaluate episodic memory retrieval performance across diverse temporal scales of events using temporal intersection over union (tIoU), which measures the overlap between retrieved and ground truth segments as the ratio of their intersection to their union duration. We compare WorldMM with various models in temporal grounding, single-modality retrieval, long-form egocentric video retrieval, and keyframe selection. Details about baselines are given in Sec. C.1. As shown in Tab. 3, WorldMM significantly superior tIoU scores than strong baselines. Notably, reasoning-based retrieval and keyframe selection methods exhibit lower tIoU values, indicating difficulty in handling long input contexts. Moreover, Fig. 5 demonstrates that the superior tIoU is directly correlated with higher overall accuracy, particularly in understanding long videos.

### 4.5. Efficacy of Multi-turn Retrieval

We validate the effectiveness of our model's multi-turn approach by limiting the maximum number of retrieval steps. The results in Fig. 7 show that performance consistently improves as the number of iterations increases across all benchmarks. Notably, on the EgoLifeQA benchmark, allowing a maximum of five steps yields a 9.3% improvement

Table 3. Average tIoU (%) across three benchmarks.

| Model | EgoLifeQA | Ego-R1 Bench | LVBench |
|---|---|---|---|
| Time-R1 [28] | 0.58 | 0.59 | 2.70 |
| Qwen3 Emb. [36] | 4.35 | 2.87 | 4.54 |
| HippoRAG [7] | 4.00 | 3.28 | 4.30 |
| InternVideo2 [27] | 3.36 | 2.60 | 3.55 |
| EgoRAG [33] | 3.60 | 2.73 | 3.50 |
| Ego-R1 [23] | 3.70 | 2.89 | 3.60 |
| AKS [22] | 2.75 | 2.30 | 3.52 |
| **WorldMM (Ours)** | **10.09** | **9.17** | **9.57** |

Table 4. Comparison of model variants by changing each module.

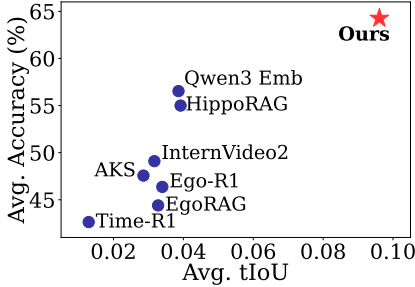| Model | EgoLifeQA | | | | | | LVBench |
|---|---|---|---|---|---|---|---|
| | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Acc. |
| *Episodic Memory* | | | | | | | |
| Fixed Timescale | 44.8 | 51.6 | 60.7 | 51.2 | 58.7 | 51.8 | 47.9 |
| Embedding Retrieval | 45.6 | 52.4 | 59.0 | 54.4 | 52.4 | 52.0 | 50.9 |
| *Semantic Memory* | | | | | | | |
| w/o Consolidation | 48.8 | 53.2 | 57.4 | 51.2 | **60.7** | 53.0 | 54.2 |
| *Visual Memory* | | | | | | | |
| Feature Retrieval | 45.6 | 51.6 | 62.3 | **58.4** | 55.6 | 53.6 | 52.4 |
| Timestamp Retrieval | 41.6 | 50.0 | **63.9** | 56.8 | 54.0 | 51.8 | 52.9 |
| **WorldMM (Ours)** | **49.6** | **56.4** | 63.9 | **58.4** | 58.7 | **56.4** | **55.4** |



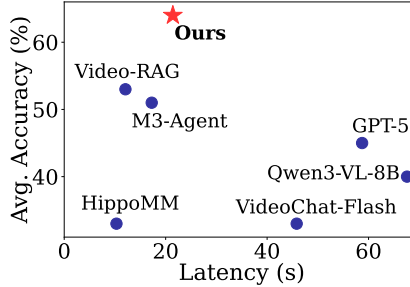Figure 5. Average tIoU and performance of WorldMM and baselines.



Figure 6. Average latency and performance of WorldMM and baselines.
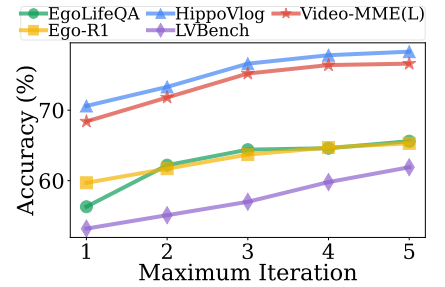


Figure 7. Accuracy of WorldMM with different maximum retrieval steps.

over single-step retrieval. This gain arises because multiple iterations enable the retrieval agent to gather additional relevant information and refine its retrieval strategy when earlier attempts are suboptimal. An example of this refinement is shown in Sec. E.2, where the model corrects an initially irrelevant retrieval to produce a more accurate and contextually grounded response.

### 4.6. Analysis on Efficiency

To assess the efficiency of our framework, we measure the end-to-end latency of WorldMM on 100 randomly sampled queries from EgoLifeQA. As shown in Fig. 6, our method achieves a superior latency–accuracy trade-off compared with baselines. Long-video LLMs incur significantly higher inference latency, while still exhibiting relatively low performance. Although RAG- or memory-based approaches offer better latency, they often require substantial preprocessing and show a significant performance gap. In contrast, by allowing the retrieval agent to adaptively finish iterations and by retrieving only the relevant segments, WorldMM achieves low latency and substantially higher accuracy.

### 4.7. Efficacy of Memory Modules

To enable effective long video reasoning, WorldMM applies different strategies for each memory. Tab. 4 reports the results of WorldMM-8B under various module configurations with details about each method in Sec. C.2. For episodic memory, using a fixed single timescale or embed-

dings instead of graphs results in a 6.1% and 4.4% drop in average accuracy, respectively, highlighting the importance of multi-scale structured knowledge. For semantic memory, removing the consolidation process results in approximately 7% drops for the category that requires long-term reasoning, demonstrating the need for continuous integration of knowledge to support long-term reasoning. Finally, for visual memory, disabling its dual-mode retrieval leads to an accuracy drop of about 3%, indicating that each mode contributes complementary benefits for retrieving particular scenes or accessing broader temporal ranges.

## 5. Conclusion

We propose WorldMM, a novel memory agent designed to perceive and remember the world as represented in long video streams. To address the challenges of long video reasoning, we introduce a multimodal, multi-scale memory that integrates textual and visual information through adaptive retrieval. By constructing separate memories across different modalities and timescales, together with a retrieval agent that iteratively identifies relevant information, our approach enables effective and flexible reasoning over long videos. We validate our model on multiple benchmarks ranging from hour- to week-long videos, demonstrating that WorldMM provides a promising solution capable of robust performance across various long video reasoning tasks.

# References

[1] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025. 1, 5, 6, 15

[2] Keshigeyan Chandrasegaran, Agrim Gupta, Lea M. Hadzic, Taran Kota, Jimming He, Cristóbal Eyzaguirre, Zane Durante, Manling Li, Jiajun Wu, and Li Fei-Fei. Hourvideo: 1-hour video-language understanding. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. 2

[3] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 2, 5, 6, 15

[4] Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi Li, Zhi Gao, and Qing Li. Videoagent: A memory-augmented multimodal agent for video understanding. In *European Conference on Computer Vision*, pages 75–92, 2024. 1

[5] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118, 2025. 1, 2, 5, 11, 12

[6] Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*, 2024. 5, 6, 12, 15

[7] Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. From rag to memory: Non-parametric continual learning for large language models. *arXiv preprint arXiv:2502.14802*, 2025. 5, 6, 8, 12, 15, 16

[8] Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. Ma-lmm: Memory-augmented large multimodal model for long-term video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13504–13514, 2024. 2

[9] Sunil Hwang, Jaehong Yoon, Youngwan Lee, and Sung Ju Hwang. Everest: Efficient masked video autoencoder by removing redundant spatiotemporal tokens. In *International Conference on Machine Learning*, 2024. 2

[10] Soyeong Jeong, Kangsan Kim, Jinheon Baek, and Sung Ju Hwang. VideoRAG: Retrieval-augmented generation over video corpus. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 21278–21298, Vienna, Austria, 2025. Association for Computational Linguistics. 3

[11] Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhan Zhu, Haian Huang, Jianfei Gao, Kunchang Li, Yinan He, Chenting Wang, et al. Videochat-flash: Hierarchical com-

[12] Yueqian Lin, Qinsi Wang, Hancheng Ye, Yuzhe Fu, Hai Li, Yiran Chen, et al. Hippomm: Hippocampal-inspired multimodal memory for long audiovisual event understanding. *arXiv preprint arXiv:2504.10739*, 2025. 1, 2, 3, 5, 6, 11, 15

[13] Lin Long, Yichen He, Wentao Ye, Yiyuan Pan, Yuan Lin, Hang Li, Junbo Zhao, and Wei Li. Seeing, listening, remembering, and reasoning: A multimodal agent with long-term memory. *arXiv preprint arXiv:2508.09736*, 2025. 1, 2, 3, 4, 5, 6, 12, 14, 15

[14] Yongdong Luo, Xiawu Zheng, Xiao Yang, Guilin Li, Haojia Lin, Jinfa Huang, Jiayi Ji, Fei Chao, Jiebo Luo, and Rongrong Ji. Video-rag: Visually-aligned retrieval-augmented long video comprehension. *arXiv preprint arXiv:2411.13093*, 2024. 3, 5, 6, 12, 15

[15] Rui Meng, Ziyan Jiang, Ye Liu, Mingyi Su, Xinyi Yang, Yuepeng Fu, Can Qin, Zeyuan Chen, Ran Xu, Caiming Xiong, et al. Vlm2vec-v2: Advancing multimodal embedding for videos, images, and visual documents. *arXiv preprint arXiv:2507.04590*, 2025. 5

[16] OpenAI. Gpt-5 system card, 2025. 1, 2, 5, 6, 15

[17] Jongwoo Park, Kanchana Ranasinghe, Kumara Kahatapitiya, Wonjeong Ryu, Donghyun Kim, and Michael S Ryoo. Too many frames, not all useful: Efficient strategies for long-form video qa. *arXiv preprint arXiv:2406.09396*, 2024. 2

[18] Xubin Ren, Lingrui Xu, Long Xia, Shuaiqiang Wang, Dawei Yin, and Chao Huang. Videorag: Retrieval-augmented generation with extreme long-context videos. *arXiv preprint arXiv:2502.01549*, 2025. 3

[19] Saul Santos, António Farinhas, Daniel C McNamee, and André FT Martins. ∞-video: A training-free approach to long video understanding via continuous-time memory consolidation. *arXiv preprint arXiv:2501.19098*, 2025. 2

[20] Xiaoqian Shen, Yunyang Xiong, Changsheng Zhao, Lemeng Wu, Jun Chen, Chenchen Zhu, Zechun Liu, Fanyi Xiao, Balakrishnan Varadarajan, Florian Bordes, Zhuang Liu, Hu Xu, Hyunwoo J. Kim, Bilge Soran, Raghuraman Krishnamoorthi, Mohamed Elhoseiny, and Vikas Chandra. Longvu: Spatiotemporal adaptive compression for long video-language understanding. *arXiv preprint arXiv:2410.17434*, 2024. 2

[21] Xiaoqian Shen, Wenxuan Zhang, Jun Chen, and Mohamed Elhoseiny. Vgent: Graph-based retrieval-reasoning-augmented generation for long video understanding, 2025. 3

[22] Xi Tang, Jihao Qiu, Lingxi Xie, Yunjie Tian, Jianbin Jiao, and Qixiang Ye. Adaptive keyframe sampling for long video understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29118–29128, 2025. 2, 8, 13, 16

[23] Shulin Tian, Ruiqi Wang, Hongming Guo, Penghao Wu, Yuhao Dong, Xiuying Wang, Jingkang Yang, Hao Zhang, Hongyuan Zhu, and Ziwei Liu. Ego-r1: Chain-of-toolthought for ultra-long egocentric video reasoning. *arXiv preprint arXiv:2506.13654*, 2025. 1, 3, 5, 6, 8, 11, 12, 15, 16

[24] Shaoguang Wang, Ziyang Chen, Yijie Xu, Weiyu Guo, and Hui Xiong. Less is more: Token-efficient video-qa via adap-

pression for long-context video modeling. *arXiv preprint arXiv:2501.00574*, 2024. 2, 5, 6, 15

9

tive frame-pruning and semantic graph integration. *arXiv preprint arXiv:2508.03337*, 2025. 2

[25] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Ming Ding, Xiaotao Gu, Shiyu Huang, Bin Xu, et al. Lvbench: An extreme long video understanding benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22958–22967, 2025. 1, 2, 5, 11

[26] Xidong Wang, Dingjie Song, Shunian Chen, Junyin Chen, Zhenyang Cai, Chen Zhang, Lichao Sun, and Benyou Wang. Longllava: Scaling multi-modal llms to 1000 images efficiently via a hybrid architecture. *arXiv preprint arXiv:2409.02889*, 2024. 2

[27] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Zun Wang, Yansong Shi, et al. Internvideo2: Scaling foundation models for multimodal video understanding. In *European Conference on Computer Vision*, pages 396–416, 2024. 8, 12, 16

[28] Ye Wang, Ziheng Wang, Boshen Xu, Yang Du, Kejun Lin, Zihan Xiao, Zihao Yue, Jianzhong Ju, Liang Zhang, Dingyi Yang, et al. Time-r1: Post-training large vision language model for temporal video grounding. *arXiv preprint arXiv:2503.13377*, 2025. 2, 5, 6, 8, 12, 15, 16

[29] Ziyang Wang, Jaehong Yoon, Shoubin Yu, Md Mohaiminul Islam, Gedas Bertasius, and Mohit Bansal. Video-rts: Rethinking reinforcement learning and test-time scaling for efficient and enhanced video reasoning. *arXiv preprint arXiv:2507.06485*, 2025. 2, 5, 6, 15

[30] Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3272–3283, 2025. 2

[31] Zeyu Xu, Junkang Zhang, Qiang Wang, and Yi Liu. E-vrag: Enhancing long video understanding with resource-efficient retrieval augmented generation, 2025. 3

[32] Zhucun Xue, Jiangning Zhang, Xurong Xie, Yuxuan Cai, Yong Liu, Xiangtai Li, and Dacheng Tao. Adavideorag: Omni-contextual adaptive retrieval-augmented efficient long video understanding. *arXiv preprint arXiv:2506.13589*, 2025. 3

[33] Jingkang Yang, Shuai Liu, Hongming Guo, Yuhao Dong, Xiamengwei Zhang, Sicheng Zhang, Pengyun Wang, Zitang Zhou, Binzhu Xie, Ziyue Wang, et al. Egolife: Towards egocentric life assistant. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28885–28900, 2025. 1, 2, 3, 4, 5, 6, 8, 11, 12, 15, 16

[34] Emmanouil Zaranis, António Farinhas, Saul Santos, Beatriz Canaverde, Miguel Moura Ramos, Aditya K Surikuchi, André Viveiros, Baohao Liao, Elena Bueno-Benito, Nithin Sivakumaran, Pavlo Vasylenko, Shoubin Yu, Sonal Sannigrahi, Wafaa Mohammed, Ben Peters, Danae Sánchez Villegas, Elias Stengel-Eskin, Giuseppe Attanasio, Jaehong Yoon, Stella Frank, Alessandro Suglia, Chrysoula Zerva, Desmond Elliott, Mariella Dimiccoli, Mohit Bansal, Oswald Lanz, Raffaella Bernardi, Raquel Fernández, Sandro

Pezzelle, Vlad Niculae, and André F. T. Martins. Movie facts and fibs (mf2): A benchmark for long movie understanding. *arXiv preprint arXiv:2506.06275*, 2025. 2

[35] Pan Zhang, Xiaoyi Dong, Yuhang Cao, Yuhang Zang, Rui Qian, Xilin Wei, Lin Chen, Yifei Li, Junbo Niu, Shuangrui Ding, et al. Internlm-xcomposer2.5-omnilive: A comprehensive multimodal system for long-term streaming video and audio interactions. *arXiv preprint arXiv:2412.09596*, 2024. 2

[36] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, et al. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025. 8, 12, 16

[37] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun MA, Ziwei Liu, and Chunyuan Li. LLaVA-video: Video instruction tuning with synthetic data. *Transactions on Machine Learning Research*, 2025. 1

# WorldMM: Dynamic Multimodal Memory Agent for Long Video Reasoning

## Supplementary Material

## A. Additional Details on Dataset

In this section, we provide additional details for each dataset used in our experiments. Tab. 5 summarizes the datasets, including the number of queries, domain categories, and the average video duration.

Table 5. Summary of benchmark datasets used in experiments.

| Dataset | # Queries | Domain | Avg. Video Length |
|---|---|---|---|
| EgoLifeQA [33] | 500 | Egocentric | 44.3h |
| Ego-R1 Bench [23] | 300 | Egocentric | 44.3h |
| HippoVlog [12] | 1,000 | Vlog | 0.45h |
| LVBench [25] | 1,534 | General | 1.14h |
| Video-MME (L) [5] | 900 | General | 0.69h |

### A.1. EgoLifeQA

EgoLifeQA [33] is a set of questions designed to test the capability of models to understand and remember everyday life from week-long video recordings. It includes questions that require recalling past events, tracking object locations, and reasoning over long-term activities. In our experiments, we use questions from the perspective of a single participant (A1: JAKE), along with his corresponding video stream, which spans 44.3 hours. The benchmark is organized into five distinct categories as follows.

**EntityLog (Ent.)** Questions that require recalling information about objects, such as their locations, states, or interactions. (Example: "Who used the screwdriver first?")

**EventRecall (EvR.)** Questions that ask about specific past events, including what happened, when it occurred, and relevant context. (Example: "Shure mentioned Tiramisu, when was the last time we discussed making Tiramisu?")

**HabitInsight (Hab.)** Questions aimed at identifying a person's recurring behaviors or long-term activity patterns. (Example: "What food does Alice love to eat?")

**RelationMap (Rel.)** Questions involving understanding social relationships and interactions between people. (Example: "Who usually sings when Shure plays the guitar?")

**TaskMaster (Task)** Questions focused on ongoing or pending tasks that require reasoning about what actions still need to be completed. (Example: "What are we planning to do in the afternoon?")

### A.2. Ego-R1 Bench

Ego-R1 Bench [23] is designed as a complementary evaluation to EgoLifeQA, but with a distinct focus on model reasoning. While both benchmarks focus on the same week-long egocentric video, Ego-R1 Bench targets multi-step, tool-augmented reasoning over ultra-long video. We reorganize query types of Ego-R1 Bench to the category adopted by EgoLifeQA, as shown in Tab. 6.

Table 6. Classification of queries under the EgoLifeQA category.

| Category | Ego-R1 Category |
|---|---|
| EntityLog | EntityLog, FoodLog, HealthLog, TechLog |
| EventRecall | EventRecall, Event Recollection, Event Memory |
| HabitInsight | HabitInsight, Behavior Habit(s) |
| RelationMap | RelationMap, Interpersonal Relationships |
| TaskMaster | TaskMaster, Future Plan(s) |

### A.3. HippoVlog

HippoVlog [12] contains 25 daily vlog videos with 1,000 multiple-choice questions for continuous audiovisual event understanding. The benchmark evaluates a model's ability to handle modality-specific information, with **Auditory (Aud.)** questions requiring reasoning over the audio stream (or transcript) and **Visual (Vis.)** questions focusing on the visual content. **Auditory+Visual (A+V)** queries test the model's ability to integrate information across both modalities, while **Summarization (Summ.)** questions assess higher-level reasoning over long temporal spans, requiring synthesis of events and semantic understanding from the continuous video.

### A.4. LVBench

LVBench [25] consists of 103 long videos, typically longer than an hour, with 1,549 multiple-choice questions for extreme long video understanding. The videos cover a general and diverse set of domains. Questions include both visual perception for recognizing entities or events in short segments and summarization for higher-level reasoning across extended sequences, evaluating models' ability to integrate information over both local and long-horizon contexts. In our experiments, we categorize questions into three groups based on their segment length, defined as the duration of video required to answer the question: **Short** (<30s), **Medium (Med.)** (30s~5min), and **Long** (>5min). We excluded 15 questions without segment tags, leaving 1,534 questions in total for evaluation.

### A.5. Video-MME

Video-MME [5] is a comprehensive video understanding benchmark with 2,700 questions and varying video durations. In this experiment, we use only the long subset (>30min), containing 900 questions, to assess the model's capability on long video reasoning. We adopt the categories provided by the benchmark, with acronyms as follows: Action Reasoning (ARES), Action Recognition (AREC), Attribute Perception (ATTR), Counting Problem (CNT), Information Synopsis (ISYN), OCR Problems (OCR), Object Reasoning (ORES), Object Recognition (OREC), Spatial Perception (SPER), Spatial Reasoning (SRES), Temporal Perception (TPER), and Temporal Reasoning (TRES).

## B. Additional Implementation Details

We provide additional details on the baseline setup (Sec. B.1), the configuration of our proposed WorldMM (Sec. B.2), and the prompts used (Sec. B.3).

### B.1. Baseline Setup

**Base Models & Long Video LLMs** For all base models and long video LLMs, the video input is uniformly sampled at 0.5 fps and capped at 768 frames, since we cannot process all frames due to context limit, as mentioned in Sec. 1. For Time-R1 [28], we employ the 7B checkpoint[1].

**RAG-based Video LLMs** For text-based RAG video models, we construct a knowledge base from video captions. Specifically, each video is segmented into 30 second chunks, and set of captions from these segments serve as retrieval pool. LightRAG [6] performs dual-level retrieval, selecting either fine-grained (low-level) or abstracted (high-level) information from the knowledge graph generated from set of captions depending on the query. HippoRAG [7], in contrast, retrieves raw captions ranked by their PPR scores, treating each caption as a separate document. For Video-RAG [14] model, retrieval is performed directly on the raw video using tools such as optical character recognition (OCR) and automatic speech recognition (ASR) to extract textual signals. Unless otherwise stated, we follow the retrieval specifications described in each model's corresponding paper or implementation.

**Memory-based Video LLMs** Memory-based video LLMs construct explicit memories from the video stream. For EgoRAG [33] and Ego-R1 [23], which build hierarchical textual memories, we use the same temporal granularity applied when constructing WorldMM's memory. For models that perform iterative reasoning, including Ego-R1 [23] and M3-Agent [13], we evaluate the checkpoints released

---

[1] https://huggingface.co/Boshenxx/Time-R1-7B

by authors and set the maximum number of reasoning iterations to 5 to ensure consistent evaluation across all systems. All other implementation details follow the official specifications provided by the respective authors.

### B.2. WorldMM

To construct multi-scale episodic memory, we tailor the temporal resolutions to each dataset's duration. For Ego-LifeQA and Ego-R1 Bench, which contain week-long videos, we use four broad timescales: 30 seconds, 3 minutes, 10 minutes, and 1 hour. For HippoVlog, LVBench, and Video-MME, which contain shorter recordings averaging about an hour, we adopt shorter timescales of 10 seconds, 30 seconds, 3 minutes, and 10 minutes to better match their temporal structure. For semantic memory, triplets with a similarity score above 0.6 are consolidated using an LLM, and the top 10 triplets are retrieved at query time. The retrieval agent is limited to a maximum of five iterations, consistent with the baseline evaluation setting.

### B.3. Prompts

To construct and retrieve memory, and to generate the final response of WorldMM, we employ carefully optimized prompts for use with an LLM. In particular, we use prompts for episodic triple extraction (Figs. 9 and 10) and multi-scale memory construction (Fig. 11), adapted from Yang et al. [33]. Furthermore, we utilize prompts for multi-scale memory retrieval (Fig. 12), semantic triple extraction (Fig. 13), semantic consolidation (Fig. 14), iterative reasoning by the retrieval agent (Fig. 15), and final response generation (Fig. 16).

## C. Additional Description on Experiments

In this section, we provide additional description of the settings used in our ablation experiments.

### C.1. Dynamic Temporal Scope Retrieval (Sec. 4.4)

To evaluate performance on dynamic temporal reasoning with WorldMM, we employ several approaches, including temporal grounding model, embedding-based retrieval models, hierarchical retrieval models, and keyframe selection method. For each method, we measure tIoU using either the returned timestamps or the timestamps of the selected content. For the temporal grounding model, we use Time-R1 [28], with a slightly modified prompt that enables it to return both the evidence timestamps and the corresponding grounded responses. We sample videos at 0.5 fps and provide up to 768 frames. For embedding-based and hierarchical retrieval models, we follow the configurations described in Sec. B.1. Additionally, we include Qwen3 Emb., which applies the Qwen3-Embedding-4B [36] text encoder for caption retrieval, and InternVideo2, which encodes each segment using InternVideo2 [27] as an video en-

coder with uniform 16 frame averaging to enable segment-level retrieval. Both methods retrieve 30 second segments based on similarity search. For key frame selection, we apply AKS [22], which selects keyframes from the 0.5 fps sampled sequence. For tIoU evaluation, we interpret frames as representing their corresponding 30 second segments.

## C.2. Efficacy of Memory Modules (Sec. 4.7)

To assess the contribution of each component within WorldMM's multimodal memory system, we evaluate several ablated variants in Sec. 4.7. In this section, we detail each variant of WorldMM created by selectively disabling a specific component. For episodic memory variants, we first construct a **fixed timescale** variant by replacing hierarchical episodic memory with a single fixed timescale memory. Specifically, we use the episodic memory of the finest granularity timescale. We also experiment an **embedding retrieval** variant in which the model's graph-based episodic retrieval is replaced with an embedding-based similarity search using Qwen-Embedding-4B. To examine the effect of semantic consolidation, we use a **w/o consolidation** version that bypasses the consolidation procedure to update the memory and instead store the raw extracted triplets without any update to existing memory. Finally, for visual memory, we ablate components of dual-retrieval mechanism by evaluating systems that rely exclusively on either **feature retrieval** through natural-language keyword search or **timestamp retrieval** based purely on temporal indices.

## D. Detailed Experimental Results

**Main results** Tabs. 7 and 8 present the category-wise performance breakdown of WorldMM and baseline methods. Beyond overall benchmark averages, WorldMM consistently outperforms existing approaches across most categories. Notably, the gains are especially larger in categories that rely on visual information. For instance, in the EntityRecall category of EgoLifeQA, where visual cues can help answering, WorldMM exceeds the previous best method, Ego-R1, by a substantial 11.2%. Similarly, on HippoVlog, our model achieves a 4% improvement in the Aud. and A+V categories, both of which require visual reasoning. These margins are greater than those observed in categories that do not explicitly depend on visual content, highlighting the strong advantage of our multimodal multi-memory architecture.

**Efficacy of multimodal memory** Fig. 8 shows memory type utilization of our model on HippoVlog benchmark, where categories are grouped by their modality requirements. The Audio category requires reasoning over spoken content and therefore is expected to depend primarily on textual memory derived from caption transcripts, while the

Visual category focuses on visual understanding and correspondingly is designed to rely more on visual memory. Our results clearly support these expectations, showing that the Audio category predominantly activates textual memory while the Visual category relies heavily on visual memory, indicating that each category effectively leverages the required memory. Moreover, the Summarization category, which requires long-term reasoning, utilizes semantic memory more than any other category, demonstrating the complementary roles and effectiveness of each memory module in handling different reasoning demands. Together with this distribution of memory usage and the demonstrated performance gains in Tab. 2, these underscore the effectiveness of our multimodal multi-memory framework.

**Dynamic temporal scope retrieval** Tabs. 9 and 10 detail the per-category tIoU and accuracy results for WorldMM and baseline methods. While WorldMM significantly outperforms existing baselines on average, the results on LVBench particularly highlight the effectiveness of our dynamic episodic memory. In LVBench's Long category, where answering requires reasoning over more than five minutes of video, WorldMM outperforms the baselines by a notably larger margin than in categories that require shorter timescale, underscoring its ability to flexibly retrieve and integrate information over diverse temporal spans.

## E. Qualitative Results

### E.1. Memory Construction

Tab. 11 presents an example of episodic triplet extraction. Given a caption generated from sampled frames of a segment along with its corresponding transcript, an LLM is prompted (using the prompt in Fig. 10) to extract episodic triplets. Semantic triplets are extracted using a different prompt (Fig. 13), designed to focus on long-term dependencies and capture more abstract relationships across the segments, as shown in Tab. 12. To better capture persistent knowledge across segments, we introduce semantic consolidation, which incrementally updates the semantic graph by integrating new triplets and resolving conflicts. Using embedding-based matching and an LLM, duplicated or conflicting triplets are removed, and new or revised ones are added, generating an evolving semantic memory, as shown in Tab. 13. For instance, the new triplet "[I, uses WeChat for, money transfers]" is merged with the existing triplet to consolidate redundant information, and conflicting triplets, such as "[Lucia, dislikes, overly sweet food]" versus "[Lucia, likes, sweet desserts]", are removed to ensure consistency in the semantic memory.

### E.2. Multi-turn Refinement

WorldMM demonstrates the effectiveness of multi-turn reasoning by progressively refining its retrieval strategy to answer questions, as shown in Tab. 14. In this example, the first round retrieves episodic memory using a narrow keyword focused on the "discussion" of the air conditioning, but it provides insufficient detail about the activity. In the second round, the model expands to a more general keyword, "air conditioning", which enables retrieval of every scene where the air conditioning is involved to obtain sufficient textual evidence. Moreover, in the third round, since the textual evidence fails to capture specific visual details of the scene, WorldMM refines its strategy to retrieve video frames corresponding to the relevant timestamp. Through this stepwise process, WorldMM effectively refines its search strategy with different keyword strategies and memory types to respond to the question.

## F. Limitation and Broader Impact

While WorldMM serves as an effective multimodal memory agent for long video reasoning, it still requires careful preprocessing, including video captioning, triplet extraction, and semantic consolidation. Yet, this limitation is not unique to our approach but a broader constraint shared by existing memory-based video LLMs. For example, M3-Agent [13] incurs even heavier preprocessing due to its reliance on entity recognition, and many other approaches operate with offline preprocessing. In contrast, WorldMM is designed for online operation. Memories are updated at fixed intervals (e.g., every 10 seconds), and the required preprocessing for each segment can be performed within these windows. Moreover, new information can be seamlessly integrated into the knowledge graph, and our consolidation mechanism efficiently refines the knowledge base without requiring the reconstruction of memory from scratch.

With strong long-term reasoning capabilities and support for real-time updates, WorldMM serves as a practical solution for streaming scenarios such as egocentric assistants and embodied agents. This foundation enables richer and more persistent assistance for everyday tasks and accessibility. However, the continuous accumulation of structured knowledge over periods of time raises serious privacy and security concerns. Real-world deployments must therefore enforce safeguard policies, including strict access controls, secure data handling, and privacy protections.

Table 7. Category-wise performance breakdown of WorldMM and baselines on EgoLifeQA, Ego-R1 Bench, HippoVlog, and LVBench.

| Model | EgoLifeQA | | | | | | Ego-R1 Bench | | | | | | HippoVlog | | | | | LVBench | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Aud. | Vis. | A+V | Summ. | Avg. | Short | Med. | Long | Avg. |
| *Base Models* | | | | | | | | | | | | | | | | | | | | | |
| Qwen3-VL-8B [1] | 35.2 | 30.2 | 39.3 | 46.4 | 46.0 | 38.6 | 31.8 | 41.5 | 38.5 | 42.1 | 44.7 | 35.7 | 73.6 | 74.0 | 69.2 | 80.8 | 74.4 | 48.8 | 44.4 | 53.4 | 48.3 |
| Gemini 2.5 Pro [3] | 43.2 | 40.5 | 41.0 | 55.2 | 52.4 | 46.4 | 43.9 | 56.1 | 53.9 | 47.4 | 47.4 | 46.7 | 69.2 | 75.2 | 63.6 | 80.0 | 72.0 | 57.1 | 52.2 | 65.2 | 57.0 |
| GPT-5 [16] | 47.2 | 42.1 | 47.5 | 53.6 | 55.6 | 48.6 | 41.8 | 58.5 | 53.9 | 52.6 | 50.0 | 46.3 | 73.6 | 75.6 | 69.2 | **84.4** | 75.7 | **59.1** | 59.1 | 69.1 | 60.4 |
| *Long Video LLMs* | | | | | | | | | | | | | | | | | | | | | |
| VideoChat-Flash [11] | 28.8 | 32.5 | 37.7 | 37.6 | 38.1 | 34.2 | 43.4 | 43.9 | 38.5 | 31.6 | 44.7 | 42.7 | 60.8 | 59.2 | 56.4 | 55.6 | 58.0 | 34.9 | 23.1 | 44.6 | 33.2 |
| Time-R1 [28] | 39.2 | 50.8 | 65.6 | 48.8 | 47.6 | 48.8 | 49.2 | 48.8 | 46.2 | 42.1 | 44.7 | 48.0 | 58.2 | 58.2 | 49.4 | 52.4 | 54.6 | 32.1 | 23.6 | 40.2 | 31.1 |
| Video-RTS [29] | 40.8 | 48.4 | 62.3 | 48.8 | 47.6 | 48.2 | 47.6 | 46.3 | 53.9 | 52.6 | 47.4 | 48.0 | 58.8 | 62.0 | 56.8 | 58.4 | 59.0 | 43.4 | 25.7 | 49.5 | 39.8 |
| *RAG-based Video LLMs* | | | | | | | | | | | | | | | | | | | | | |
| LightRAG [6] | 40.8 | 48.4 | 67.2 | 50.4 | 44.4 | 48.8 | 54.0 | 61.0 | 46.2 | 42.1 | 42.1 | 52.3 | 51.6 | 46.0 | 44.8 | 47.2 | 47.4 | 30.2 | 28.6 | 34.3 | 30.4 |
| HippoRAG [7] | 48.8 | 60.3 | 70.5 | 60.8 | 66.7 | 59.6 | 54.5 | 65.9 | 69.2 | 52.6 | 50.0 | 56.0 | 72.4 | 53.2 | 54.0 | 73.2 | 63.2 | 54.9 | 47.5 | 62.3 | 54.0 |
| Video-RAG [14] | 49.6 | 56.3 | 67.2 | 55.2 | 54.0 | 55.4 | 48.7 | 58.5 | 53.9 | 47.4 | 44.7 | 49.7 | 63.2 | 64.8 | 63.6 | 68.8 | 65.1 | 32.9 | 30.2 | 39.7 | 33.1 |
| *Memory-based Video LLMs* | | | | | | | | | | | | | | | | | | | | | |
| EgoRAG [33] | 40.0 | 56.3 | 62.3 | 54.4 | 52.4 | 52.0 | 46.6 | 56.1 | 46.2 | 47.4 | 55.3 | 49.0 | 64.8 | 53.2 | 47.6 | 64.4 | 57.5 | 32.4 | 32.0 | 31.9 | 32.2 |
| Ego-R1 [23] | 51.2 | 53.2 | 63.9 | 50.4 | 50.8 | 53.0 | 50.8 | 63.4 | 38.5 | 36.8 | 57.9 | 52.0 | 57.2 | 58.8 | 52.0 | 67.2 | 58.8 | 32.5 | 36.5 | 37.3 | 34.1 |
| HippoMM [12] | 45.6 | 53.2 | 70.5 | 55.2 | 58.7 | 54.6 | 51.9 | 56.1 | 46.2 | 52.6 | 57.9 | 53.0 | 68.8 | 77.6 | 59.2 | 82.0 | 71.9 | 40.7 | 33.3 | 35.8 | 38.2 |
| M3-Agent [13] | 44.4 | 54.8 | 62.3 | 56.8 | 54.0 | 53.5 | 52.4 | 58.5 | 38.5 | 42.1 | 52.6 | 52.0 | 68.4 | 72.4 | 50.8 | 70.4 | 65.5 | 53.0 | 40.7 | 48.5 | 49.3 |
| *WorldMM (Ours)* | | | | | | | | | | | | | | | | | | | | | |
| **WorldMM-8B** | 49.6 | 56.4 | 63.9 | 58.4 | 58.7 | 56.4 | 48.2 | 63.4 | 53.9 | 52.6 | 57.9 | 52.0 | 69.6 | 73.6 | 65.2 | 70.4 | 69.7 | 55.0 | 54.1 | 59.8 | 55.4 |
| **WorldMM-GPT** | **62.4** | **64.3** | **75.4** | **62.4** | **71.4** | **65.6** | **64.6** | **70.7** | **76.9** | **57.9** | **63.2** | **65.3** | **75.6** | **81.6** | **73.2** | 82.8 | **78.3** | 58.3 | **65.4** | **72.1** | **61.9** |

Table 8. Category-wise performance breakdown of WorldMM and baselines on Video-MME (L).

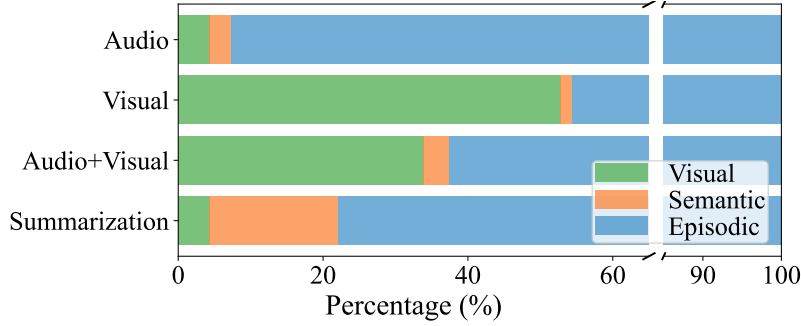| Model | ARES | AREC | ATTR | CNT | ISYN | OCR | ORES | OREC | SPER | SRES | TPER | TRES | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Base Models* | | | | | | | | | | | | | |
| Qwen3-VL-8B [1] | 62.2 | 54.0 | 51.9 | 43.8 | 68.1 | 42.9 | 62.9 | 57.4 | 33.3 | 45.5 | 33.3 | 67.0 | 61.0 |
| Gemini 2.5 Pro [3] | 56.9 | 47.6 | 66.7 | 41.7 | 71.8 | **57.1** | 53.3 | 40.7 | 0.0 | 72.7 | **66.7** | 48.4 | 55.7 |
| GPT-5 [16] | 71.1 | 69.8 | 70.4 | 47.9 | **88.3** | **57.1** | **75.8** | 74.1 | 33.3 | 72.7 | 50.0 | 75.8 | 74.3 |
| *Long Video LLMs* | | | | | | | | | | | | | |
| VideoChat-Flash [11] | 35.0 | 42.9 | 37.0 | 31.3 | 34.4 | 42.9 | 60.0 | 46.3 | 33.3 | 54.5 | 33.3 | 46.2 | 44.1 |
| Time-R1 [28] | 20.6 | 28.6 | 25.9 | 35.4 | 31.9 | 35.7 | 53.3 | 48.2 | 33.3 | 36.4 | 50.0 | 44.0 | 37.6 |
| Video-RTS [29] | 43.3 | 52.4 | 40.7 | 39.6 | 33.7 | 42.9 | 60.8 | 53.7 | 33.3 | 45.5 | 50.0 | 49.5 | 47.9 |
| *RAG-based Video LLMs* | | | | | | | | | | | | | |
| LightRAG [6] | 41.7 | 30.2 | 40.7 | 35.4 | 54.0 | 50.0 | 46.7 | 61.1 | 33.3 | 45.5 | 50.0 | 52.8 | 46.6 |
| HippoRAG [7] | 45.6 | 47.6 | 40.7 | 37.5 | 52.2 | 42.9 | 52.9 | 64.8 | **66.7** | 54.5 | 50.0 | 70.3 | 52.1 |
| Video-RAG [14] | 51.7 | 47.6 | 37.0 | 39.6 | 49.7 | 57.1 | 62.1 | 68.5 | **66.7** | 45.5 | 50.0 | 68.1 | 55.4 |
| *Memory-based Video LLMs* | | | | | | | | | | | | | |
| EgoRAG [33] | 31.1 | 55.6 | 33.3 | 22.9 | 41.1 | 28.6 | 44.6 | 48.2 | 33.3 | 54.5 | **66.7** | 48.4 | 41.1 |
| Ego-R1 [23] | 37.2 | 52.4 | 40.7 | 35.4 | 38.0 | 35.7 | 42.1 | 51.9 | **66.7** | 63.6 | 50.0 | 52.8 | 42.7 |
| HippoMM [12] | 41.1 | 42.9 | 55.6 | 35.4 | 38.7 | 35.7 | 37.9 | 53.7 | 33.3 | 54.5 | 50.0 | 47.3 | 41.6 |
| M3-Agent [13] | 52.2 | 57.1 | 59.3 | 45.8 | 51.5 | 42.9 | 54.6 | 64.8 | 33.3 | 45.5 | 50.0 | 71.4 | 55.3 |
| *WorldMM (Ours)* | | | | | | | | | | | | | |
| **WorldMM-8B** | 65.0 | 66.7 | 59.3 | 41.7 | 72.4 | 42.9 | 67.5 | 72.2 | 33.3 | 54.5 | **66.7** | 69.2 | 66.0 |
| **WorldMM-GPT** | **81.1** | **73.0** | **70.4** | **54.2** | 85.3 | 42.9 | 75.0 | **77.8** | 33.3 | **72.7** | **66.7** | **79.1** | **76.6** |

Figure 8. Memory type utilization of WorldMM on four distinctive categories in HippoVlog.

Table 9. Category-wise average tIoU (%) breakdown of WorldMM and dynamic temporal scope retrieval baselines.

| Model | EgoLifeQA | | | | | | Ego-R1 Bench | | | | | | LVBench | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Short | Med. | Long | Avg. |
| Time-R1 [28] | 0.34 | 0.72 | 1.07 | 0.52 | 0.41 | 0.58 | 0.27 | 0.84 | 0.71 | 1.15 | 1.58 | 0.59 | 3.10 | 2.60 | 1.00 | 2.70 |
| Qwen3 Emb. [36] | 2.87 | 4.31 | 5.58 | 2.98 | 8.91 | 4.35 | 2.68 | 2.74 | 3.85 | 2.74 | 3.70 | 2.87 | 4.48 | 6.20 | 1.75 | 4.54 |
| HippoRAG [7] | 3.02 | 4.19 | 4.99 | 2.12 | 8.36 | 4.00 | 3.32 | 2.85 | 3.28 | 2.23 | 4.07 | 3.28 | 4.23 | 5.76 | 1.88 | 4.30 |
| InternVideo2 [27] | 2.09 | 4.42 | 6.04 | 2.00 | 3.88 | 3.36 | 2.71 | 2.55 | 3.09 | 1.85 | 2.32 | 2.60 | 3.66 | 4.71 | 0.87 | 3.55 |
| EgoRAG [33] | 3.20 | 3.38 | 4.62 | 3.10 | 4.82 | 3.60 | 2.40 | 3.07 | 4.08 | 2.19 | 3.78 | 2.73 | 4.10 | 3.38 | 0.91 | 3.50 |
| Ego-R1 [23] | 3.31 | 3.52 | 5.03 | 2.87 | 5.18 | 3.70 | 2.57 | 2.83 | 4.13 | 2.83 | 4.12 | 2.89 | 4.08 | 3.72 | 1.14 | 3.60 |
| AKS [22] | 2.42 | 2.77 | 3.08 | 2.93 | 2.67 | 2.75 | 2.03 | 2.48 | 2.99 | 2.58 | 3.04 | 2.30 | 3.81 | 4.11 | 1.10 | 3.52 |
| **WorldMM (Ours)** | **9.79** | **10.43** | **11.85** | **7.73** | **12.97** | **10.09** | **8.91** | **9.85** | **8.86** | **9.63** | **9.58** | **9.17** | **7.53** | **14.41** | **10.02** | **9.57** |

Table 10. Category-wise performance breakdown of WorldMM and dynamic temporal scope retrieval baselines.

| Model | EgoLifeQA | | | | | | Ego-R1 Bench | | | | | | LVBench | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Ent. | EvR. | Hab. | Rel. | Task | Avg. | Short | Med. | Long | Avg. |
| Time-R1 [28] | 39.2 | 50.8 | 65.6 | 48.8 | 47.6 | 48.8 | 49.2 | 48.8 | 46.2 | 42.1 | 44.7 | 48.0 | 32.1 | 23.6 | 40.2 | 31.1 |
| Qwen3 Emb. [36] | 44.0 | 59.5 | 70.5 | 58.4 | 68.3 | 57.8 | 51.9 | 65.9 | 61.5 | 57.9 | 47.4 | 54.0 | 52.9 | 49.1 | 62.3 | 53.2 |
| HippoRAG [7] | 48.8 | 60.3 | 70.5 | 60.8 | 66.7 | 59.6 | 54.5 | 65.9 | 69.2 | 52.6 | 50.0 | 56.0 | 54.9 | 47.5 | 62.3 | 54.0 |
| InternVideo2 [27] | 40.8 | 54.0 | 60.7 | 51.2 | 52.4 | 50.6 | 50.3 | 56.1 | 46.2 | 47.4 | 52.6 | 51.0 | 47.4 | 37.3 | 53.4 | 45.7 |
| EgoRAG [33] | 40.0 | 56.3 | 62.3 | 54.4 | 52.4 | 52.0 | 46.6 | 56.1 | 46.2 | 47.4 | 55.3 | 49.0 | 32.4 | 32.0 | 31.9 | 32.2 |
| Ego-R1 [23] | 51.2 | 53.2 | 63.9 | 50.4 | 50.8 | 53.0 | 50.8 | 63.4 | 38.5 | 36.8 | 57.9 | 52.0 | 32.5 | 36.5 | 37.3 | 34.1 |
| AKS [22] | 41.6 | 51.6 | 63.9 | 51.2 | 52.4 | 50.6 | 51.3 | 63.4 | 46.2 | 36.8 | 50.0 | 51.7 | 43.3 | 33.9 | 39.2 | 40.4 |
| **WorldMM (Ours)** | **62.4** | **64.3** | **75.4** | **62.4** | **71.4** | **65.6** | **64.6** | **70.7** | **76.9** | **57.9** | **63.2** | **65.3** | **58.3** | **65.4** | **72.1** | **61.9** |

Table 11. Example of episodic triplet extraction.

| | |
|---|---|
| **Caption** | I stand and walk to the other side of the dining table. Katrina asks, "Is this for tomorrow's game?" "Yes—let's think about what to do tomorrow," I say. I raise my right hand as Katrina walks toward me. Lucia asks, "Using ancient poems? Or what else?" Katrina says, "I'm not good with ancient poems." Tasha asks, "Then what else to use?" Katrina says, "I'll be out in the first round. My room is already cleaned up." "Okay," I say. I turn toward the stairs, put down my phone, look back at the living room door, and walk into the second-floor living room. Lucia adds, "For example, not coming out." Katrina says, "Let me check that place we're going to." Tasha asks, "I just want to ask which fields it has expanded into." Lucia says, "Okay." |
| **Extracted Triplets** | [I, stand at, dining table] <br> [I, walk to, other side of the dining table] <br> [Katrina, asks about, tomorrow] <br> [I, confirm, tomorrow] <br> [I, raise, right hand] <br> [Katrina, walks toward, I] <br> [Lucia, asks about, using ancient poems] <br> [Katrina, says, not good with ancient poems] <br> [Tasha, asks, what else to use] <br> [Katrina, says, I will be out in the first round] <br> [Katrina, has, room already cleaned up] <br> [I, turn toward, stairs] <br> [I, put down, phone] <br> [I, look back at, living room door] <br> [I, walk into, second-floor living room] <br> [Lucia, adds, not coming out as an example] <br> [Katrina, says, let me check that place we're going to] <br> [Lucia, says, Okay] |

17

Table 12. Example of semantic triplet extraction.

| | |
|---|---|
| **Caption** | I got up, moved my phone, and checked it before turning it off. Alice expressed her feelings towards me, and I responded by checking my phone's chat interface. Alice then questioned her appearance, and I turned off the phone, looking around at the snacks and utensils on the table. I stood up, grabbed a pack of snacks, and proceeded to my room to enjoy them. Alice asked about something being fancy, and I fetched my glasses, placing them on the table. ... I managed my phone, swiping through pages, and interacted with others as I went about my tasks. I observed Alice and Tasha, discussing what to feed a cat, and continued interacting with my phone. As the environment darkened, I engaged with the surroundings, noting the layout and structures. Finally, I moved towards a house with blue-green walls, managing my power bank and surveying the area. |
| **Extracted Triplets** | [I, assigns tasks to, Katrina]<br>[I, handles reimbursements for, Alice]<br>[I, uses WeChat for, money transfers]<br>[I, often eats, snacks]<br>[I, wears, glasses]<br>[Lucia, dislikes, overly sweet food]<br>[Alice, expresses romantic feelings toward, I]<br>[Katrina, helps with, expense tracking]<br>[I, requires PDFs for, reimbursement]<br>[Tasha, participates in, house demolition tasks]<br>[Lucia, participates in, house demolition tasks] |

Table 13. Example of semantic consolidation.

| | | |
|---|---|---|
| **Original Triplets** | [I, uses WeChat to send money]<br>[I, wears, glasses]<br>[I, often eats, fruits]<br>[Lucia, likes, sweet desserts]<br>[Tasha, participates in, household projects] | |
| **New Triplets** | [I, assigns tasks to, Katrina]<br>[I, handles reimbursements for, Alice]<br>[I, uses WeChat for, money transfers]<br>[I, often eats, snacks]<br>[I, wears, glasses]<br>[Lucia, dislikes, overly sweet food]<br>[Alice, expresses romantic feelings toward, I]<br>[Katrina, helps with, expense tracking]<br>[I, requires PDFs for, reimbursement]<br>[Tasha, participates in, house demolition tasks]<br>[Lucia, participates in, house demolition tasks] | *% conflicts with existing "likes sweet desserts"* |
| **Consolidated Triplets** | [I, assigns tasks to, Katrina]<br>[I, handles reimbursements for, Alice]<br>[I, uses, WeChat to send money]<br>[I, often eats, fruits and snacks]<br>[I, wears, glasses]<br>[Alice, expresses romantic feelings toward, I]<br>[Katrina, helps with, expense tracking]<br>[I, requires PDFs for, reimbursement]<br>[Tasha, participates in, house demolition tasks]<br>[Lucia, participates in, house demolition tasks] | *% merged with "uses WeChat for money transfers"*<br>*% merged and generalized "fruits" and "snacks"*<br>*% merged with duplicate in new triples*<br><br><br>*% refined general triple "household projects"* |

Table 14. Example of the multi-turn refinement of WorldMM.

| | |
|---|---|
| **Question** | What were we doing last time we discussed the air conditioning temperature? (A) Eating hot pot (B) Shopping (C) Playing piano (D) Picking up a package |
| **Answer** | A |
| **Round 1** | Decision: Search // Memory: Episodic<br>Search Query: discussing the air conditioning temperature<br>*Retrieved:*<br>*[DAY2 13:36:00-13:39:00]*<br>*I joined my colleagues at the dining table and engaged in a discussion about the food while preparing to eat. ... During the meal, we discussed the air conditioner situation in the rooms. Tasha and Katrina talked about where the remotes were and who had access, with Shure mentioning that his room's AC had been on continuously due to a borrowed remote from me. I listened to their conversation, ...*<br>*[DAY2 17:48:00-17:51:00]*<br>*I coordinated with Shure and Lucia to find a suitable spot for my power bank, initially placing it on a flower pot as Shure suggested, but realizing it was unstable. ... We discussed the room environment, including turning off the air conditioner that was blowing cold air. Shure clarified that the remote we had was for the TV, not the air conditioner, prompting me to inquire about the correct remote. ...* |
| **Round 2** | Decision: Search // Memory: Episodic<br>Search Query: air conditioning<br>*Retrieved:*<br>*[DAY2 17:48:30-17:49:00]*<br>*"Flip it and put it underneath," I say. ... "Wow, this, this…" Lucia murmurs. "How do you turn off this air conditioner? It's blowing so cold up here." I adjust the power bank for a better position. Shure says, "Just tell it to turn off, it'll turn off." "Really? It's voice-controlled?" Lucia asks. I secure the power bank and turn to Shure and Lucia. Shure adds, "No, this is the TV remote." "Where is it?" I ask, and I chat with them.*<br>*[DAY2 18:34:01-18:34:29]*<br>*I watch Shure set the air conditioner to 26 degrees. "Save the Earth," I say. "Everyone's responsibility," I laugh. I kneel, then sit on the sofa. Shure asks, "Seriously, you think it's good to record this?" I look at my friends chatting and watch the movie. Shure says, "We're still using the AC." Shure says, "He really can't hear a thing." I watch Shure take a few steps. I laugh. Shure says, "Hahaha, it's like we're watching a silent movie." I smile and keep watching the movie. "Whoa, this is classic," I say. I stand up. "Come on, let's have a drink," I say.*<br>*[DAY2 18:33:00-18:33:30]*<br>*I walk in front of the fill light, pass behind Choiszt, and move toward the curtain. ... "The popularity is always high," I say. Shure: "Ha." Choiszt: "Turn it on." "Turn on the air conditioner?" I ask. Lucia: "I'll just, I suggest I drink water directly from a bowl." Lucia: "So later I'll eat with this bowl."* |
| **Round 3** | Decision: Search // Memory: Visual<br>Search Query: DAY2 18:34:01-18:34:29<br>*Retrieved:*<br><br>*(Author Note: They are gathered around a table cooking and sharing a hot pot meal while watching a projection on the wall.)* |
| **Round 4** | Decision: Answer |
| **Response** | A |

Your task is to extract named entities from the given paragraph. Respond with a JSON list of entities.

Example:

Radio City is India's first private FM radio station and was started on 3 July 2001. It plays Hindi, English and regional songs. Radio City recently forayed into New Media in May 2008 with the launch of a music portal - PlanetRadiocity.com that offers music related news, videos, songs, and other music-related features.

{ "named_entities":
    ["Radio City", "India", "3 July 2001", "Hindi", "English", "May 2008", "PlanetRadiocity.com"]
}

Figure 9. Prompt for named entity recognition (NER). Recognized named entities are used to extract episodic triplets as shown in Fig. 10.

Your task is to construct an RDF (Resource Description Framework) graph from the given passages and named entity lists. Respond with a JSON list of triples, with each triple representing a relationship in the RDF graph.

Pay attention to the following requirements:
- Each triple should contain at least one, but preferably two, of the named entities in the list for each passage.
- When resolving pronouns, if the pronoun refers to the first-person (e.g., I, me, my), keep it as "I" instead of replacing with terms like "speaker" or "narrator". For other pronouns, clearly resolve them to their specific names to maintain clarity.

Convert the paragraph into a JSON dict, it has a named entity list and a triple list.

Example:

Radio City is India's first private FM radio station and was started on 3 July 2001. It plays Hindi, English and regional songs. Radio City recently forayed into New Media in May 2008 with the launch of a music portal - PlanetRadiocity.com that offers music related news, videos, songs, and other music-related features.

{ "named_entities":
    ["Radio City", "India", "3 July 2001", "Hindi", "English", "May 2008", "PlanetRadiocity.com"]
}

{ "triples": [
    ["Radio City", "located in", "India"],
    ["Radio City", "is", "private FM radio station"],
    ["Radio City", "started on", "3 July 2001"],
    ["Radio City", "plays songs in", "Hindi"],
    ["Radio City", "plays songs in", "English"],
    ["Radio City", "forayed into", "New Media"],
    ["Radio City", "launched", "PlanetRadiocity.com"],
    ["PlanetRadiocity.com", "launched in", "May 2008"],
    ["PlanetRadiocity.com", "is", "music portal"],
    ["PlanetRadiocity.com", "offers", "news"],
    ["PlanetRadiocity.com", "offers", "videos"],
    ["PlanetRadiocity.com", "offers", "songs"]
]}

Figure 10. Prompt for episodic triplet extraction.

As an Event Summary Documentation Specialist, your role is to systematically structure and summarize event information, ensuring that all key actions of major characters are captured while maintaining clear event logic and completeness. Your focus is on concise and factual summarization rather than detailed transcription.

# Specific Requirements

1. Structure the Events Clearly
- Merge related events: Consolidate similar content into major events and arrange them in chronological order to ensure a smooth logical flow.
- Logical segmentation: Events can be grouped based on location, task, or theme. Each event should have a clear starting point, progression, and key turning points without any jumps or fragmentation in the information.

2. Retain Key Information
- The primary character's ("I") decisions and actions must be fully presented, including all critical first-person activities. Transitions between different parts, such as moving between floors or starting/ending a task, should be seamless.
- Any discussions, decisions, and task execution involving the primary character and other key individuals that impact the main storyline must be reflected. This includes recording, planning, and confirming matters, but in a concise manner.
- The purpose and method of key actions must be recorded, such as "ordering takeout using a phone" or "documenting a plan on a whiteboard."

3. Concise Expression, Remove Redundancies
- Keep the facts clear, avoiding descriptions of atmosphere, emotions, or abstract content.
- Remove trivial conversations and extract only the core topics and conclusions of discussions. If a discussion is lengthy, summarize it into task arrangements, decision points, and specific execution details.

4. Strictly Adhere to Facts, No Assumptions
- Do not make assumptions or add interpretations—strictly organize content based on available information, ensuring accuracy. Every summarized point must have a basis in the original information, with no unnecessary additions.
- Maintain the correct chronological order of events. The sequence of developments must strictly follow their actual occurrence without any inconsistencies.

# Output Format
Each paragraph should represent one major event, structured in a summary-detail-summary format. Strictly output below 500 words in total. Do not report the word count in the output.

Figure 11. Prompt for episodic memory construction to generate coarser-level caption.

You are an expert assistant that helps filter and select relevant video captions based on a given query. Your task is to analyze the retrieved video captions and determine which ones are most relevant to answer the question.

Given the following question and retrieved video captions, select and rank the most relevant captions that should be used to answer the question.

Instructions:
1. Consider the nature of the question when selecting captions:
   - e.g., for queries about specific events, focus on finer granularities; for habitual, relationship, or general queries, consider coarser granularities.
   - Note that coarser granularity captions may provide broader context, but finer granularity captions often contain more specific details.
2. Each caption shows its time range (start_time to end_time)
3. Analyze each caption for relevance to the question
4. Select captions that directly help answer the question
5. Return the IDs in ranked order (most relevant first)
6. Only include captions that are truly relevant

Return ONLY a JSON array of caption IDs in order of relevance (most relevant first), without additional justification.

Figure 12. Prompt for episodic memory retrieval to select from multiple timescales.

You are tasked with extracting semantic knowledge from episodic triples. Your goal is to infer generalizable information that extends beyond the specific episode. Focus on capturing valid semantic triples that can guide reasoning about behavior, relationships, or preferences.

# What to Extract
1. Relationships: social bonds or roles between entities that persist over time
   (e.g., "Alice is a friend with Bob", "Jason is a teacher of Alice").
2. Attributes & Preferences: tendencies, likes/dislikes, personality-like traits, or behavioral habits
   (e.g., "Alice prefers not having dessert", "Bob enjoys music").
3. Habits & Capabilities: actions or patterns that suggest what an entity often does, can do, or tends to do
   (e.g., "Alice often helps friends", "Jason can give advice").
4. Conceptual Knowledge: directly useful facts that support reasoning, but avoid overly broad taxonomic statements
   (e.g., "Alice's office is near Cafe X", "Bob's gym is closed on Sundays").

# What to Avoid
- One-off events or transient states (e.g., "ate pizza yesterday", "was late once") unless explicitly declared as a preference/role
- Broad taxonomy or trivia unrelated to behavior (e.g., "a laptop is electronics", "Paris is in France")
- Speculative or mind-reading inferences without textual support (e.g., motives, beliefs not evidenced)

# Important Notes
- Prefer to base semantic triples on multiple supporting episodes.
- BUT if a single episode clearly reflects a role, preference, habit, or capability, it is valid to include it.
- Each semantic triple MUST have at least one supporting episodic triple.
- Reduce duplication. If multiple episodic triples support the same or very similar semantic knowledge, merge them into one semantic triple rather than repeating.
- The 'episodic_evidence[i]' list must always point to the indices that support 'semantic_triples[i]'.
- Aim for broad coverage: extract as many valid semantic triples as reasonably supported by the input.

# Output Format
- Return ONLY a JSON object with the following two keys:
  - 'semantic_triples' (List[List[str]]): Each item is a triple [subject, predicate, object].
  - 'episodic_evidence' (List[List[int]]) : Each item is a list of 0-based indices pointing to the input episodic triples that support the corresponding semantic triple at the same position.
- The two lists MUST have the same length and aligned order.
- If no semantic knowledge is inferable, return: {"semantic_triples": [], "episodic_evidence": []}

Example:

Episodic triples:
0. ["Alice", "talks to", "Bob"],
1. ["Alice", "laughs with", "Bob"],
2. ["Alice", "doesn't eat cake", "at restaurant"],
3. ["Alice", "shares personal stories with", "Bob"],
4. ["Alice", "brings coffee to", "Bob"],
5. ["Jason", "talks to", "Alice"],
6. ["Alice", "declines dessert", "at friend's house"]

Output:
{
    "semantic_triples": [
        ["Alice", "is a friend with", "Bob"],
        ["Alice", "prefers", "not having dessert"]
    ],
    "episodic_evidence": [
        [0, 1, 3],
        [2, 6]
    ]
}

Figure 13. Prompt for semantic triplet extraction.

You are tasked with consolidating semantic knowledge by processing a new semantic triple against relevant existing knowledge from previous timestamps.

Your job is to make two decisions:
1. Which existing triples to remove/pop — those that should be merged with the new triple or conflict with it
2. How to update the new triple — to capture merged information or resolve conflicts

# Consolidation Rules
1. Merge Similar Information: If existing triples express very similar information to the new triple, remove them and update the new triple to contain the most complete/accurate form.
2. Resolve Conflicts: If the new triple conflicts with existing ones, decide which is more accurate/recent and remove outdated ones.
3. Update with Context: Use information from existing triples to make the new triple more specific or more accurate.
4. Preserve Unique Information: Only remove existing triples when they are redundant or conflicting.

# Output Format
Return ONLY a JSON object with the following two keys:
- 'updated_triple' (List[str]): The new triple, possibly updated [subject, predicate, object].
- 'triples_to_remove' (List[int]): Indices of existing triples to remove (empty list if none).

Example:

New triple: ["Alice", "enjoys", "coffee"]

Existing triples:
0. ["Alice", "likes", "beverages"]
1. ["Alice", "favors", "to have coffee after dinner"]
2. ["Alice", "prefers", "hot drinks"]
3. ["Alice", "likes to drink", "coffee"]

Output:
{
    "updated_triple": ["Alice", "likes", "coffee"],
    "triples_to_remove": [1, 3]
}

Figure 14. Prompt for semantic memory consolidation.

You are a reasoning agent for a video memory retrieval system. Your job is to decide whether to stop and answer, or to search memory for more evidence. When searching, you must select exactly one memory type and form a query.

# Decision Modes
1. search: Retrieve memory to begin, continue, or extend progress toward the answer
   - Choose one memory type and form a keyword(phrase)-style search query.
2. answer: Stop searching because the accumulated results are sufficient.
   - No memory type selection is needed.

# Memory Types
1. Episodic: Specific events/actions. Stores memories of past events and actions. Query by EVENT/ACTION.
2. Semantic: Entities/relationships. Stores factual knowledge about entities and their relationships, roles, and habits. Query by ENTITY/CONCEPT.
3. Visual: Scene/setting snapshots. Stores visual snapshots of scenes and settings. Query by SCENE/SETTING or TIMESTAMP RANGE.
   - For timestamp range queries, return in the format: DAY X HH:MM:SS - DAY Y HH:MM:SS.
# Context Inputs
- Current Query
- Round History: Log of past retrieval rounds. Each round is written in this format:

    ### Round N
    Decision: <search|answer>
    Memory: <episodic|semantic|visual>
    Search Query: <query text>
    Retrieved: <retrieved items>

# Strict Output Rules
- If decision = "search": Must include "selected_memory" with exactly one memory type and one query.
- If decision = "answer": Do NOT include "selected_memory".
- Always output in valid JSON only, no extra commentary.

# Output Format
{
    "decision": "search" | "answer",
    "selected_memory": {
        "memory_type": "episodic" | "semantic" | "visual",
        "search_query": <str>
    } # Omit if decision = "answer"
}

(Few-shot examples given)

Figure 15. Prompt for retrieval agent to decide retrieval strategy.

You are an AI assistant that answers questions about egocentric video experiences using retrieved memory context. Your task is to answer multiple choice questions based on this accumulated context. Always choose the most relevant answer from the given choices based on the evidence provided.

# Guidelines
- Analyze all provided context carefully.
- Choose the answer that best matches the evidence.
- If evidence is unclear, make the most reasonable inference.

# Output Format
Provide your answer as a single letter (A, B, C, or D) based on the evidence.

Figure 16. Prompt for response agent to generate response based on retrieved results.