

ViSAGE: Constructing Self-Correcting Memories for Long-Form Video Understanding

Anonymous ACL submission

Abstract

Multimodal agents operating in long-horizon environments must build and continually update memory to support entity-consistent, temporally grounded reasoning. However, existing agentic memory approaches often discard fine-grained identity cues under aggressive compression and segment-wise processing. They also over-trust vector-similarity retrieval, surfacing semantically related yet identity-mismatched evidence and thereby causing identity errors and hallucinations. We propose ViSAGE, a multimodal agentic memory framework that constructs self-correcting, entity-centric memories. Specifically, ViSAGE anchors entity identity via cross-modal binding over long temporal ranges. It then applies bidirectional memory refinement to propagate delayed identity evidence, retroactively unifying historical records and stabilizing future reasoning. We also introduce multi-agent cross-verification to adjudicate retrieved evidence under an identity–evidence alignment constraint, enabling verified refusals instead of hallucinations when evidence is missing. Extensive results demonstrate that ViSAGE consistently outperforms the strongest baseline, achieving 5.9% higher accuracy.

1 Introduction

With the rapid progress of multimodal large language models (MLLMs) (Google DeepMind, 2025; Lin et al., 2024b; Li et al., 2023), agents can now perceive and communicate through vision, audio, and language (Fan et al., 2024; Fung et al., 2025). However, deploying such agents in long-horizon settings requires memory that persists beyond a single context window. Therefore, prior work typically uses retrieval-augmented memory to extend MLLMs beyond finite context windows (Park et al., 2023; Zhong et al., 2024; Chhikara et al., 2025; Liu et al., 2024; Hu et al., 2025). In this paradigm, the agent writes ob-

servations to external memory and retrieves relevant entries to guide generation and decisions. Many systems use vector-store retrieval (Zhong et al., 2024; Chhikara et al., 2025; Fan et al., 2024, 2025), while hierarchical designs organize memories across multiple time scales for long-horizon planning (Liu et al., 2024; Hu et al., 2025). Optimus-1 (Li et al., 2024) further explores hybrid multimodal memory, combining a hierarchical directed knowledge graph with an abstracted experience pool to encode world knowledge and past multimodal experience. More recently, M3-Agent structures memory as an entity-centric graph to better support long-horizon reasoning over entities (Long et al., 2025).

Although these efforts have shown promising performance in long-horizon multimodal agentic memory tasks, we identify several limitations in current paradigms:

- **Spatio-Temporal Detail Loss.** Existing systems face a trade-off between span and detail. To ingest extensive multimodal streams, they often rely on aggressive down-sampling or compression, which erases critical micro-dynamics and thus removes fine-grained identity cues needed for reliable entity disambiguation.
- **Segmentation Dilemma.** Many long-horizon pipelines process inputs in isolated chunks. Identity evidence often arrives late and must be linked back to earlier events. Chunk-wise processing blocks cross-segment revision, so late cues cannot fix earlier memories.
- **Similarity–Veracity Gap.** Retrieval-based memory conflates semantic relevance with factual correctness. Vector-similarity retrieval can surface contextually related yet identity-mismatched or incorrect evidence; without verification, agents are prone to answering based on noisy retrieval.

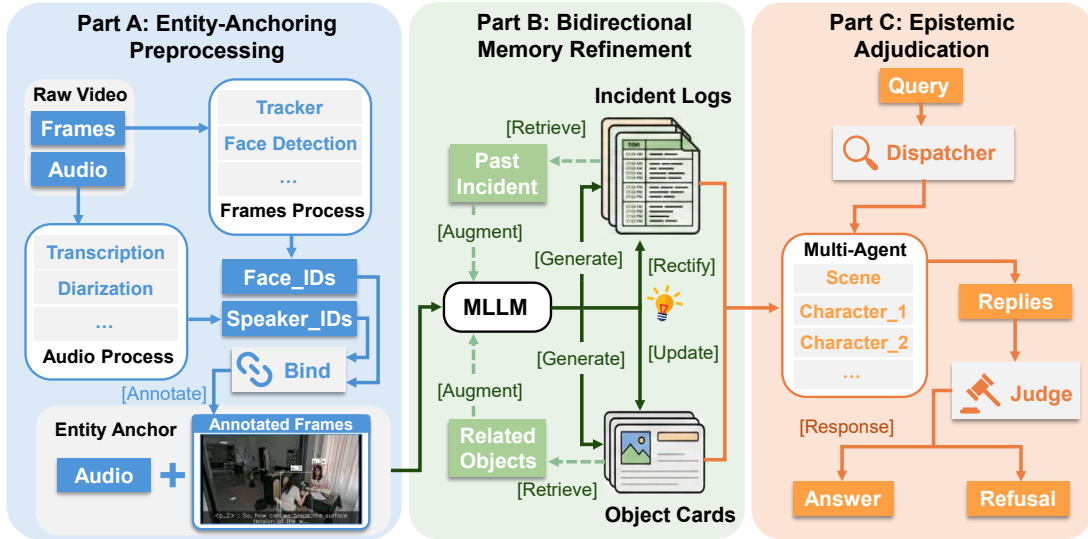


Figure 1: Architecture of ViSAGE. The framework proceeds in three stages: (A) Entity-Anchoring Preprocessing binds visual and audio signals to annotated frames; (B) Bidirectional Context Refinement maintains self-correcting memory (Incident Logs and Object Cards) via retrieval and rectification; and (C) Epistemic Adjudication employs multi-agent cross-verification to ensure reliable responses.

To address these challenges, we introduce ViSAGE, a **Visual Self-correcting AGentic** memory framework that constructs *entity-centric, self-correcting* memories beyond linear, chunk-wise processing. ViSAGE maintains a dual-structured memory: *Incident Logs* record sequential narrative events, while *Object Cards* maintain evolving entity profiles that accumulate and revise identity attributes. To align identity cues across spatial degradation and temporal sparsity, we perform **Cross-Modal Entity Binding** by binding long-range visual tracklets with audio signals. To overcome isolated segmentation, ViSAGE further introduces **Bidirectional Memory Refinement** to propagate delayed identity evidence, retroactively unifying historical records and stabilizing future reasoning. Finally, we propose **Multi-Agent Cross-Verification** to adjudicate retrieved evidence under an identity–evidence alignment constraint, enabling verified refusals instead of hallucinations when evidence is missing.

We evaluate ViSAGE through extensive experiments on long-horizon multimodal benchmarks, including M3-Bench-robot, M3-Bench-web (Long et al., 2025), and Video-MME-long (Fu et al., 2025). Experimental results show that ViSAGE consistently outperforms strong baselines across diverse reasoning categories by 5.9%, demonstrating robust gains on identity-critical queries. Moreover, our results verify that Multi-Agent Cross-Verification improves reliability by

enforcing identity-evidence alignment, enabling verified refusals instead of hallucinated answers when supporting evidence is missing or noisy. Our contributions are summarized as follows:

- **Framework Design:** We introduce *ViSAGE*, a visual self-correcting agentic memory framework that enables *retroactive, entity-centric* updates beyond chunk-wise processing. ViSAGE separates memory into sequential *Incident Logs* and structured *Object Cards*, allowing delayed identity evidence to be aligned with historical events for long-horizon entity consistency.
- **Algorithmic Novelty:** We propose **Bidirectional Memory Refinement** to resolve the Segmentation Dilemma. This mechanism transcends sequential processing via a feedback loop that retroactively unifies fragmented narratives using delayed identity evidence, ensuring long-term entity consistency.
- **Epistemic Safety:** To bridge the Similarity-Veracity Gap, we introduce **Multi-Agent Cross-Verification**. This framework enforces identity-constrained reasoning to mitigate hallucinations, prioritizing Verified Epistemic Refusals over plausible fabrications in high-stakes embodied scenarios.
- **Empirical Impact:** ViSAGE sets new state-of-the-art on authoritative benchmarks. Our framework consistently outperforms baselines,

142 achieving accuracy of 45.5% on M3-Bench-
143 robot, 58.4% on M3-Bench-web, and 79.1% on
144 Video-MME-long, validating its robustness in
145 long-form video understanding.

146 2 Related Work

147 2.1 Long Video Understanding

148 Recent advancements in Multimodal Large Lan-
149 guage Models (MLLMs), such as GPT-4o (Hurst
150 et al., 2024), Gemini-3-Pro (Google DeepMind,
151 2025), and the open-source Qwen3-VL (Bai
152 et al., 2025), have revolutionized visual percep-
153 tion. While methods like Video-LLaVA (Lin et al.,
154 2024a) and VILA (Lin et al., 2024b) excel in short
155 clips, extending them to hour-long videos remains
156 a challenge due to context limits. Existing solu-
157 tions typically follow two paradigms:

158 Early works, inspired by Socratic Models (Zeng
159 et al., 2022), decompose video reasoning into text-
160 based tasks. Methods like Video Recap (Islam
161 et al., 2024) and AutoAD (Han et al., 2023) gen-
162 erate hierarchical captions or summaries to fit tex-
163 tual RAG systems (Lewis et al., 2020). Recent
164 benchmarks like EgoSchema (Mangalam et al.,
165 2023) and LongVideoBench (Wu et al., 2024)
166 highlight that while efficient, these methods suf-
167 fer from severe information loss, discarding fine-
168 grained visual cues essential for reasoning.

169 Recently, another line of work has fo-
170 cused on expanding the model to retain
171 visual features within limited context win-
172 dows. Specifically, token compression methods
173 such as LongVILA (Chen et al., 2024b) and
174 LongVU (Shen et al., 2024), utilize spatiotem-
175 poral pooling to reduce redundancy, while
176 VidCompress (Lan et al., 2024) and Video-
177 XL (Liu et al., 2025) employ adaptive selection to
178 fit hour-scale videos. In contrast, memory-centric
179 models like MovieChat (Song et al., 2024), MA-
180 LMM (He et al., 2024), LifelongMemory (Wang
181 et al., 2023) and MeMViT (Wu et al., 2022)
182 introduce explicit external banks or buffers, with
183 Flash-VStream (Zhang et al., 2024) specifically
184 optimizing for real-time updates.

185 2.2 Memory Systems in AI Agents

186 Recently, AI Agents have evolved from simple in-
187 struction followers to autonomous entities capable
188 of long-horizon planning, heavily relying on ro-
189 bust memory architectures.

190 Foundational works like Generative

191 Agents (Park et al., 2023) established the impor-
192 tance of retrieving past experiences. Recent frame-
193 works have formalized this: Mem0 (Chhikara
194 et al., 2025) and MemoryBank (Zhong et al.,
195 2024) offer production-ready vector memory lay-
196 ers. HiAgent (Hu et al., 2025) and AgentLite (Liu
197 et al., 2024) extend this with hierarchical memory
198 for complex planning, while AIOS (Mei et al.,
199 2024) explores operating-system-level memory
200 management for resource optimization.

201 Extending agents to the visual domain, VideoA-
202 gent (Fan et al., 2024) and its embodied variant
203 Embodied VideoAgent (Fan et al., 2025) utilize
204 tools for iterative video retrieval. Jarvis-1 (Wang
205 et al., 2024) integrates memory with multimodal
206 planning. A-Mem (Xu et al., 2025c) explores
207 agentic memory for evolving environments. Most
208 notably, M3-Agent (Long et al., 2025) structures
209 memory into an entity-centric graph.

210 3 Method

211 As shown in Fig. 1, to solve the challenges of
212 identity fragmentation and hallucination, ViSAGE
213 introduces an integrated framework composed of
214 three distinct specialized phases:

215 3.1 Memory Bank

216 Inspired by Tulving’s multiple memory systems
217 theory (Tulving, 2002), which distinguishes time-
218 and context-bound experience from decontextual-
219 ized knowledge about entities and emphasizes a
220 mechanistic link “from mind to brain”, we design
221 ViSAGE’s memory bank as an incident–entity pair
222 (Table 1). Specifically, *Incident Logs* store an
223 append-only record of *what happened*, including
224 timestamps, actions, and dialogue. This event con-
225 tent is immutable, but the *who* field remains refin-
226 able. In parallel, *Object Cards* maintain contin-
227 uously updated entity knowledge. This separation
228 enables retroactive identity unification without cor-
229 rupting the historical event record.

230 3.2 Entity-Anchored Preprocessing

231 **Countering Spatial Degradation via Sequence-
232 Level Entity Anchoring.** The representative
233 agent-based approach (Long et al., 2025) adopts
234 a partition-and-process paradigm that slices long
235 videos into independent segments. This confines
236 identity recognition to each local window, making
237 it brittle under occlusion, pose changes, or com-
238 pression artifacts. As a result, long videos often at-
239 tenuate fine-grained spatio-temporal cues, leading

Memory Track	Schema Definition	Instance Example (Snapshot)
Incident Logs	<i>A chronological sequence of narrative events:</i> $\mathcal{L} = \{(t, \mathcal{P}_{ids}, \text{Act}, \text{Dial})\}$	[00:07] <p_1=Adam=Adan> crouches down to examine the details of the bright yellow wing structure. [00:28] <p_1=Adam=Adan> gestures towards the mechanism and asks: “What does this platform represent?”
Object Cards	<i>Entity-centric profile and state containers:</i> $\mathcal{O} = \{(\text{ID}, \mathcal{N}, \mathcal{V}, \mathcal{P}, \mathcal{S})\}$	ID: 1 Names: [Adam, Adan] Vectors: { $v_{\text{face}}, v_{\text{voice}}, v_{\text{name}}$ } Profile: {Role: Presenter, Gender: Male} State: {Loc: Workshop, Action: Gesturing}

Table 1: Data Structure of the Dual-Track Memory System. Examples are simplified for visualization.

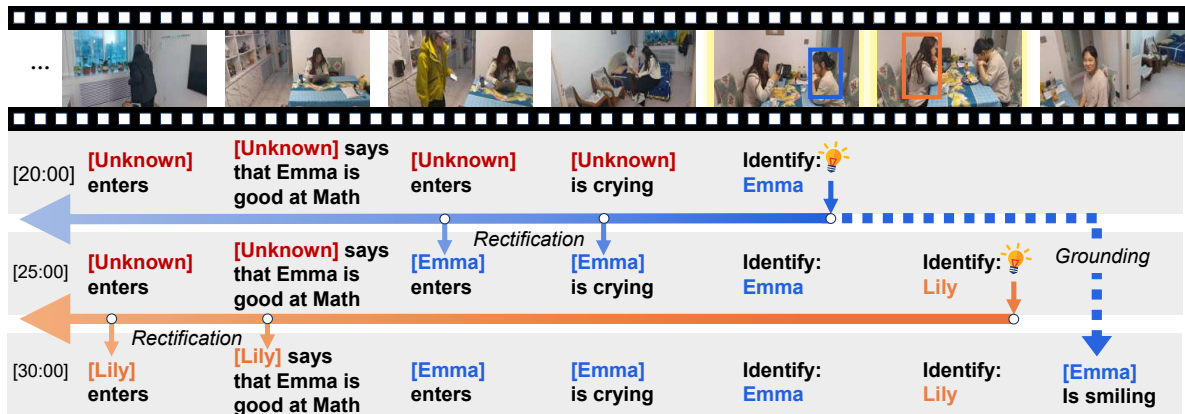


Figure 2: Visualizing Bidirectional Memory Refinement. Late-appearing Identity Triggers activate a dual loop: Backward Rectification retroactively fills historical uncertainties, while Forward Grounding propagates these resolved identities to future timestamps. This ensures actions are anchored to the correct character profile.

to identity switches and noisy event–entity associations during memory construction. Such errors propagate to both *Incident Logs* and *Object Cards*, degrading retrieval and consolidation. We therefore perform entity-anchored preprocessing to stabilize identity grounding at the sequence level.

To this end, we introduce *Sequence-Level Entity Anchoring*. Rather than relying on unreliable per-frame detections, we adopt a prototype-based strategy over long-range tracks. Specifically, we run continuous multi-object tracking over the entire video to obtain long-range trajectories, and collect all face detections within each trajectory. We then select the highest-quality face instance as a *cluster prototype* and treat it as the canonical reference. All remaining potentially low-quality face observations are matched to this prototype, clustering transient frames around a reliable anchor. This design yields a one-to-one mapping from each identity cluster to an entity in the semantic layer: the prototype initializes the corresponding *Object Card*, while the cluster assignment provides a consistent entity ID that is prop-

agated across frames. Importantly, partial cues observed in fleeting or low-quality frames can be attributed to the correct entity by referencing the prototype, enabling robust event–entity association in *Incident Logs* and consistent attribute consolidation into *Object Cards*. Robust identity is thus maintained throughout the duration of each track, ensuring the subject identity remains invariant despite transient visual noise.

Bridging Temporal Sparsity via Speaker Binding. Parallel to visual anchoring, we leverage audio cues to improve *who-did-what* attribution, where knowing *who is speaking* is decisive. A practical challenge is that standard MLLM pipelines operate on temporally sparsified video inputs, which discard the rapid lip dynamics needed for reliable speaker grounding. Meanwhile, relying solely on speaker embeddings (Long et al., 2025) can be brittle in crowded scenes due to acoustic similarity and intra-speaker variability. Motivated by the McGurk effect (McGurk and MacDonald, 1976), which illustrates how vision can disambiguate am-

286	biguous audio, we introduce <i>Dual-Criteria Validation</i> for voice–face binding. Concretely, we run Active Speaker Detection (ASD) on the high-FPS stream to obtain per-face articulation/synchrony evidence, and combine it with acoustic matching from voice embeddings. We bind a Voice_ID to a Face_ID whenever the audio is <i>confidently identifiable</i> with high voice-similarity or the video is <i>physically corroborative</i> with ASD-confirmed audio–visual synchrony within the same time window; otherwise, the speaker remains unassigned.	336
287		337
288		338
289		339
290		340
291		341
292		342
293		343
294		344
295		345
296		346
297	3.3 Bidirectional Memory Refinement	347
298	To resolve the Segmentation Dilemma, where strict sequential processing fails to link asynchronous identity evidence with prior appearances, ViSAGE implements Bidirectional Memory Refinement. Unlike unidirectional approaches that freeze past event–entity bindings once written to memory, this mechanism establishes a symmetric feedback loop: it combines Forward Grounding to stabilize current perception with Retroactive Rectification, which propagates late-arriving identity clues back into the Incident Logs to unify disjointed narrative segments, as illustrated in Fig. 2.	348
299		349
300		350
301		351
302		352
303		353
304		354
305		355
306		356
307		357
308		358
309		359
310	Forward Grounding via Identity-Keyed Retrieval. To counteract local contextual drift in the current timeframe, we implement an anchor-based retrieval strategy. Prior to memory generation, the system utilizes the stable Face_ID, derived from the tracklet anchoring phase, to retrieve memory from two parallel knowledge streams: the immediately preceding events from Incident Logs (for narrative continuity) and the character’s accumulated profile from Object Cards (for persona consistency). By priming the MLLM with this dual-source memory, we transform the interpretation of the current clip from an isolated fragment into a continuous narrative extension, ensuring the “Now” is grounded in the “Past”.	360
311		361
312		362
313		363
314		364
315		365
316		366
317		367
318		368
319		369
320		370
321		371
322		372
323		373
324		374
325	Backward Refinement via Logical Merging. Despite the robust physical continuity established in Section 3.2, where the system correctly binds the audio information to the entity, the agent remains blind to the character’s nominal identity. The mapping between a face and a name typically relies on complex dialogue comprehension, specifically determining the addressee in a conversation. To handle this, we treat every newly appearing name as a provisional semantic identity. The LLM performs memory-aware inference: if	375
326		376
327		377
328		378
329		379
330		380
331		381
332		382
333		383
334		384
335		385
	it deduces that a specific visual subject is the target of a naming utterance, such as “Hey, Mario!”, it generates a Refinement Signal ($ID_{sem} \equiv ID_{vis}$). Crucially, this signal triggers a Global Backward Update: the system traverses the historical Incident Logs \mathcal{L} , retroactively replacing all prior instances of the anonymous ID_{vis} with the resolved identity. This mechanism ensures that if an anonymous subject recorded on Day 1 is identified as “Mario” on Day 2, their entire history is unified, repairing the structural fragmentation.	336
		337
		338
		339
		340
		341
		342
		343
		344
		345
		346
	3.4 Epistemic Adjudication	347
	To bridge the Similarity–Veracity Gap, ViSAGE introduces Object Cards as a high-confidence knowledge anchor, supplementing the traditional retrieval from Incident Logs.	348
		349
		350
		351
	Noise Filtering via Identity Resolution Dispatcher. To identify the entity most relevant to the query, the Dispatcher functions as a sanitization layer against ASR-induced noise. Automatic transcription frequently introduces phonetic inconsistencies, such as erroneously transcribing the target “Mario” as “Malloy”, which causes naive vector retrieval to fail or return irrelevant distractors. Capitalizing on the fact that all valid character names are <i>indexed</i> in the database, the Dispatcher employs LLM-based logical inference to perform Identity Resolution. It explicitly maps the noisy query entity to its correct counterpart within this registry, effectively correcting errors like “Malloy” to “Mario”. It then activates only the relevant Character Agent, ensuring that retrieval is strictly targeted and robust to phonetic mismatches.	352
		353
		354
		355
		356
		357
		358
		359
		360
		361
		362
		363
		364
		365
		366
		367
		368
	Evidence Extraction via Persona-Driven Agents. Crucially, ViSAGE employs a complementary retrieval strategy specialized by query type. The Scene Agent interrogates Incident Logs to answer factual, event-driven questions. It utilizes the objective chronological record to resolve dynamic narrative details, such as “who stood up first”, which rely on precise temporal causality. Complementing this, the Character Agent leverages Object Cards to resolve static attribute inquiries. Without this solidified profile, answering identity-related questions like “What is Lily’s profession?” would necessitate on-the-fly inference based on disjointed event fragments in the Incident Logs. In contrast, ViSAGE continuously accumulates and validates identity traits throughout the video processing pipeline.	369
		370
		371
		372
		373
		374
		375
		376
		377
		378
		379
		380
		381
		382
		383
		384
		385

Method	M3-Bench-robot					M3-Bench-web					Video-MME-long		
	ME	MH	CM	PU	GK	ALL	ME	MH	CM	PU		GK	ALL
<i>Video-Native Socratic method</i>													
Qwen2.5-Omni-7b	2.1	1.4	1.5	1.5	2.1	2.0	8.9	8.8	13.7	10.8	14.1	11.3	42.2
Gemini-1.5-Pro	6.5	7.5	8.0	9.7	7.6	8.0	18.0	17.9	23.8	23.1	28.7	23.2	38.1
<i>Online Video Understanding Methods</i>													
MovieChat	13.3	9.8	12.2	15.7	7.0	11.2	12.2	6.6	12.5	17.4	11.1	12.6	19.4
MA-LMM	25.6	23.4	22.7	39.1	14.4	24.4	26.8	10.5	22.4	39.3	15.8	24.3	17.3
Flash-Vstream	21.6	19.6	19.3	24.3	14.1	19.4	24.5	10.3	24.6	32.5	20.2	23.6	25.0
<i>Discrete Socratic method</i>													
Qwen3-VL-8b	39.0	29.7	39.3	52.9	18.8	36.1	38.8	22.0	40.0	40.5	45.1	36.9	60.1
Qwen3-VL-30b	36.8	36.1	37.1	53.3	22.2	36.4	43.0	26.0	40.0	50.4	45.1	40.9	67.0
GPT-4o	34.0	32.8	32.9	42.2	16.7	29.9	-	-	-	-	-	-	65.3
Gemini-3-Pro-preview	42.7	38.3	38.0	55.1	23.3	39.6	51.9	40.3	62.9	61.0	59.0	53.8	74.2
<i>Agent-based Method</i>													
M3-Agent	32.8	29.4	31.2	43.3	19.1	30.7	45.9	28.4	44.3	59.3	53.9	48.9	61.8
ViSAGE(Ours)	45.2	47.2	47.1	56.9	30.9	45.5	62.3	47.8	71.4	63.3	62.2	58.4	79.1

Table 2: Performance Comparison on M3-Bench and Video-MME-long. We report results across different reasoning types in M3-Bench: Multi-Evidence (ME), Multi-Hop (MH), Cross-Modal (CM), Person Understanding (PU), and General Knowledge (GK).

Epistemic Adjudication via Judge Agent. In the final phase, the Judge Agent functions as the ultimate arbiter of veracity. Instead of simply aggregating outputs, it rigorously evaluates the sufficiency and consistency of the evidence retrieved from both Incident Logs and Object Cards. Crucially, when valid evidence is absent in both sources, or when irreconcilable conflicts arise, the Judge triggers a Verified Epistemic Refusal. This strict rejection mechanism ensures safety by preventing the system from hallucinating answers to unanswerable queries.

4 Experiments

4.1 Experimental Setup

Dataset Split	robot	web	V-MME-long
# Videos	100	920	300
Avg. Length	34m	27m	40m

Table 3: Statistics of evaluation benchmarks. We utilize the robot and web splits from M3-Bench and the long-video subset of Video-MME.

Evaluation Datasets. To comprehensively analyze ViSAGE, we utilize three authoritative benchmarks (summarized in Table 3). First, the robot

split of M3-Bench (Long et al., 2025) focuses on general-purpose robots, featuring egocentric videos that test memory-guided reasoning such as inferring human personalities, interpersonal relationships, and object affordances. Complementing this, the web split of M3-Bench (Long et al., 2025) and Video-MME (Fu et al., 2025) provide content with high information density from online platforms. These videos cover diverse topics like documentaries and movies, challenging the agent to process complex narratives and open-world knowledge relevant to practical multimodal applications.

Baselines. We benchmark ViSAGE against a suite of methods across three paradigms: (1) **Socratic Models**, encompassing Video-Native systems (Gemini-1.5-Pro (Team et al., 2024), Qwen2.5-Omni (Xu et al., 2025b)) that process holistic streams, and Discrete variants (GPT-4o (Hurst et al., 2024), Qwen3-VL (Bai et al., 2025), Gemini-3-Pro-preview (Google DeepMind, 2025)) utilizing frames sampled at 0.5 fps with ASR; (2) **Online Video Understanding Methods** representing diverse memory compression strategies, including MovieChat (Song et al., 2024), MA-LMM (He et al., 2024), and Flash-VStream (Zhang et al., 2024); and (3) **Agent-**

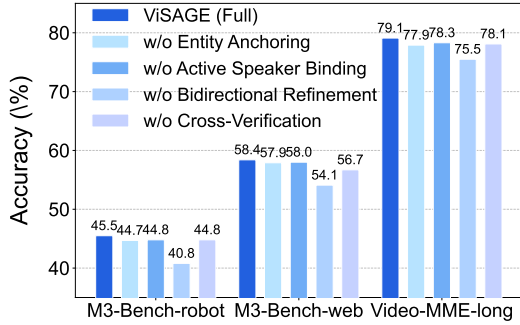


Figure 3: Ablation study on key components. The chart illustrates the impact of removing each module across three benchmarks (M3-Bench-robot, M3-Bench-web, and Video-MME-long).

based Frameworks, specifically the previous SOTA M3-Agent (Long et al., 2025), which employs an entity-centric memory graph and serves as the primary structural benchmark.

Implementation Details. We utilize Qwen3-Omni-Flash (Xu et al., 2025a) for speech transcription, followed by Gemini-3-Pro-preview as the backbone MLLM for reasoning and memory synthesis. See Appendix A for more details.

4.2 Performance Analysis

Table 2 presents the quantitative evaluation. ViSAGE consistently demonstrates superior efficacy across all benchmarks.

Crucially, regarding the Cross-Modal (CM) metric, ViSAGE achieves a decisive breakthrough. This metric directly validates the resolution of the Spatio-Temporal Detail Loss, measuring the system’s ability to bind fragmented visual and auditory cues into coherent entities. On this front, ViSAGE significantly outperforms the agentic baseline M3-Agent (improving +15.9% on robot and +27.1% on web), proving that our Entity Anchoring mechanism successfully counters signal degradation where standard agents fail.

Furthermore, ViSAGE effectively grounds the raw capabilities of the Gemini-3-Pro-preview backbone. While the backbone exhibits strong parametric knowledge, it remains vulnerable to the sampling-induced information loss, which scores only 38.0% on robot CM. By shifting from passive perception to entity-anchored reasoning, ViSAGE bridges this gap, boosting the score to 47.1%. This robust grounding capability extends to open-world long-horizon scenarios, where ViSAGE achieves 79.1% on Video-MME-long, surpassing both the ungrounded backbone (74.2%)

Backbone Model	Accuracy (%)
ViSAGE w/ Qwen3-VL-30b	42.1
ViSAGE w/ GPT-5	44.5
ViSAGE w/ Gemini-3-Pro-preview	45.5

Table 4: Impact of backbone scaling. We report the average accuracy on M3-Bench-robot (Long et al., 2025).

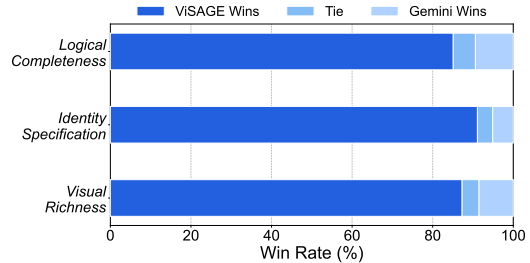


Figure 4: Quality Assessment on Correct Samples. Even when the final answer is correct, ViSAGE generates responses with significantly higher visual richness and logical completeness compared to baselines.

and the fragmentation-prone M3-Agent (61.8%).

4.3 Ablation Study

Fig. 3 and Table 4 summarize our component analysis. Structurally, Bidirectional Memory Refinement proves critical (Fig. 3); its removal triggers the sharpest drop (−4.7% on robot), confirming that Backward Rectification is essential for unifying the fragmented narratives caused by the Segmentation Dilemma. Other modules also contribute distinctively to the system’s robustness: Entity Anchoring (Part A) provides consistent gains by securing identity stability against the spatio-temporal detail loss. Multi-Agent Cross-Verification (Part C) further boosts precision by filtering semantic noise and enforcing epistemic safety. Crucially, Table 4 highlights the universality of our framework. ViSAGE maintains robust performance across distinct architectures ranging from Qwen3-VL-30b to GPT-5 (OpenAI, 2025). This validates that our Dual-Track design is a model-agnostic framework: it successfully decouples memory structure from specific reasoning backends, enhancing any foundation model with long-horizon reasoning capabilities.

4.4 Qualitative and Safety Analysis

Case Study: Resolving Pronominal Ambiguity. As detailed in Table 5, the baseline fails on pronominal resolution, succumbing to misleading hearsay. In contrast, ViSAGE successfully disentangles this ambiguity. By retroactively grounding

User Query: “What is Emma’s major?”

Ground Truth: Programming.

Gemini Response (Baseline):

“Emma’s major is **Linear Algebra** ✗.”

Retrieved Memory: A young woman says: “I know Emma is good at Linear Algebra.”

Error Analysis: The baseline relies on surface-level name matching, prioritizing the explicit mention of “Emma” by a third party while ignoring the pronoun “I” spoken by the target character (who was initially anonymous).

ViSAGE (Ours) Reasoning Process:

1. Bidirectional Refinement (Resolving “I”):

- **Initial State:** At [00:19], <p_7> says “I study Programming.” (Speaker unknown).

- **Trigger & Update:** At [05:30], identity revealed (<p_7> → Emma). System **retroactively propagates** this identity, modifying the historical log at [00:19] to: “**p_7=Emma** says ‘I study Programming.’”

2. Cross-Verification:

- **Character Agent (Emma):** “My major is **Programming**. I recall stating this at the start of the conversation.”

- **Scene Agent (Objective):** Retrieves two relevant logs:

1. [00:19] **Emma** says: “I study Programming.”

2. [02:40] **Lily** says: “...I know Emma is good at Linear Algebra.”

- **Judge Verdict:** “Evidence 1 is a **direct first-person statement**, whereas Evidence 2 is **hearsay**. Direct claim takes precedence. **Final Answer: Programming** ✓”

Table 5: Qualitative Case Study. ViSAGE successfully resolves the ambiguous pronoun “I” via retroactive refinement, whereas the baseline fails due to surface-level matching.

the first-person pronoun (“I”) to the target entity, our system enables the Judge Agent to prioritize direct testimony over third-person claims, effectively preventing hallucination.

Beyond Accuracy: Response Granularity.

Fig. 4 reveals that correct labels do not imply full comprehension. Employing GPT-5 as a judge, we evaluated three dimensions: (1) **Visual Richness** (density of scene details), where ViSAGE achieves a win rate of 87.3%; (2) **Identity Specification** (pinpointing who performs the action), achieving our highest margin of 91.1%; (3) **Logical Completeness** (providing transparent reasoning steps), with an 85.1% lead. These gaps confirm that while baselines rely on shallow pattern matching (correct label, sparse evidence), ViSAGE constructs grounded narratives with precise entity binding.

Failure Mode Analysis: The Safety of “Knowing What You Don’t Know.” Finally, we analyze failure boundaries by categorizing errors into Hazardous Hallucinations and Safe Refusals.

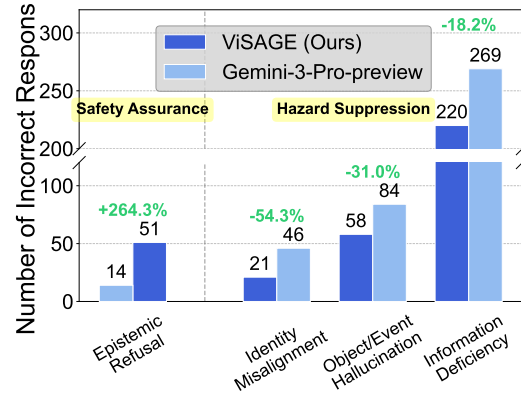


Figure 5: Safety Analysis on M3-Bench-robot. ViSAGE significantly reduces hazardous hallucinations and increases safe epistemic refusals.

Fig. 5 reveals a critical shift in ViSAGE, prioritizing safety in high-stakes scenarios:

(1) **Mitigating Hazardous Failures.** The most dangerous errors in embodied agents are Identity Misalignment (wrongly attributing actions to specific individuals) and Object/Event Hallucination (fabricating non-existent entities). ViSAGE significantly suppresses these risks, reducing Identity Misalignment by 54.3% and Hallucinations by 31.0%, thereby preventing **execution** based on false premises.

(2) **Embracing Safe Refusals.** Instead of hallucinating, ViSAGE exhibits a dramatic 264.3% increase in Epistemic Refusal: the explicit acknowledgment of information deficit. In robotics, this represents a *graceful degradation*: preferring an honest “I don’t know” over a confident error is vital for trustworthiness and safety.

5 Conclusion

In this paper, we argue that long-form video understanding requires a paradigm shift from strictly sequential processing to a self-correcting memory architecture. By enabling agents to retroactively align historical memory with evolving identity evidence, ViSAGE demonstrates that post-hoc rectification is a structural necessity for maintaining narrative coherence. Beyond accuracy, our framework establishes a new standard for epistemic safety in embodied AI: rather than prioritizing plausible fabrication, it enforces rigorous evidence verification. This transition from hazardous hallucinations to verified refusals is critical for deploying trustworthy agents in real-world environments, ensuring that future systems “know what they don’t know.”

551 Limitations

552 We acknowledge two primary limitations in the
553 current system. First, our Entity-Anchored Pre-
554 processing is explicitly tailored for human-centric
555 interactions, utilizing Face_ID and Speaker_ID
556 bindings to resolve complex social dynamics.
557 Consequently, the system presently treats non-
558 human elements, such as plot-critical objects or
559 animals, as background context rather than ac-
560 tively tracked profiles. To address this, our fu-
561 ture work aims to generalize the scope of entity
562 registration by incorporating open-vocabulary ob-
563 ject tracking and re-identification. This will evolve
564 ViSAGE from a character-centric assistant into a
565 broadly entity-aware agent. Second, relying on
566 commercial APIs implies that strictly reproducing
567 exact numerical outputs may be constrained by
568 the model’s closed-source nature. However, to en-
569 sure methodological reproducibility, we release all
570 prompt templates and pipeline code. Furthermore,
571 strictly relying on APIs limits local deployment;
572 thus, we plan to address this by fine-tuning open-
573 source models in future work to achieve a fully re-
574 producible and self-contained system.

575 **AI Assistance Statement.** We acknowledge the
576 use of Large Language Models (LLMs) to assist
577 with text refinement and LaTeX formatting. The
578 authors have verified all content for accuracy and
579 maintain full responsibility for the scientific valid-
580 ity of this work.

581 Ethical considerations

582 This work utilizes public video datasets (M3-
583 Bench, Video-MME) strictly for academic re-
584 search to enhance the consistency of multimodal
585 agents. Although the proposed framework incor-
586 porates face detection and speaker identification
587 modules, these components are employed *solely*
588 to ground entities within the narrative context of the
589 video. We do not collect any new personal data.
590 Furthermore, we explicitly disclaim any intent for
591 this technology to be used for real-world surveil-
592 lance or biometric profiling without explicit con-
593 sent and privacy safeguards.

594 References

595 Nir Aharon, Roy Orfaig, and Ben-Zion Bo-
596 brovsky. 2022. Bot-sort: Robust associa-
597 tions multi-pedestrian tracking. *arXiv preprint*
598 *arXiv:2206.14651*.

- Shuai Bai and 1 others. 2025. *Qwen3-vl technical re- 599*
port. Preprint, arXiv:2511.21631. ArXiv preprint 600
arXiv:2511.21631. 601
- Yafeng Chen, Siqi Zheng, Hui Wang, Luyao Cheng, 602
Qian Chen, Shiliang Zhang, and Junjie Li. 2024a. 603
Eres2netv2: Boosting short-duration speaker veri- 604
fication performance with computational efficiency. 605
arXiv preprint arXiv:2406.02167. 606
- Yukang Chen, Fuzhao Xue, Dacheng Li, Qinghao 607
Hu, Ligeng Zhu, Xiuyu Li, Yunhao Fang, Hao- 608
tian Tang, Shang Yang, Zhijian Liu, and 1 oth- 609
ers. 2024b. Longvila: Scaling long-context visual 610
language models for long videos. *arXiv preprint 611*
arXiv:2408.10188. 612
- Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet 613
Singh, and Deshraj Yadav. 2025. Mem0: Building 614
production-ready ai agents with scalable long-term 615
memory. *arXiv preprint arXiv:2504.19413. 616*
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos 617
Zafeiriou. 2019. Arcface: Additive angular margin 618
loss for deep face recognition. In *Proceedings of 619*
the IEEE/CVF Conference on Computer Vision and 620
Pattern Recognition (CVPR). 621
- Yue Fan, Xiaojian Ma, Rongpeng Su, Jun Guo, Rujie 622
Wu, Xi Chen, and Qing Li. 2025. Embodied videoa- 623
gent: Persistent memory from egocentric videos 624
and embodied sensors enables dynamic scene under- 625
standing. In *Proceedings of the IEEE/CVF Interna- 626*
tional Conference on Computer Vision, pages 6342– 627
6352. 628
- Yue Fan, Xiaojian Ma, Rujie Wu, Yuntao Du, Jiaqi 629
Li, Zhi Gao, and Qing Li. 2024. Videoagent: A 630
memory-augmented multimodal agent for video un- 631
derstanding. In *European Conference on Computer 632*
Vision, pages 75–92. Springer. 633
- Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, 634
Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu 635
Zhou, Yunhang Shen, Mengdan Zhang, and 1 others. 636
2025. Video-mme: The first-ever comprehensive 637
evaluation benchmark of multi-modal llms in video 638
analysis. In *Proceedings of the Computer Vision 639*
and Pattern Recognition Conference, pages 24108– 640
24118. 641
- Pascale Fung, Yoram Bachrach, Asli Celikyilmaz, 642
Kamalika Chaudhuri, Delong Chen, Willy Chung, 643
Emmanuel Dupoux, Hongyu Gong, Hervé Jégou, 644
Alessandro Lazaric, and 1 others. 2025. Embod- 645
ied ai agents: Modeling the world. *arXiv preprint 646*
arXiv:2506.22355. 647
- Google DeepMind. 2025. Gemini 3: A new era of 648
intelligence. [https://deepmind.google/models/ 649](https://deepmind.google/models/gemini/)
[gemini/](https://deepmind.google/models/gemini/). Accessed: 2025-12-13. 650
- Tengda Han, Max Bain, Arsha Nagrai, Gül Varol, 651
Weidi Xie, and Andrew Zisserman. 2023. Autoad: 652
Movie description in context. In *Proceedings of the 653*
IEEE/CVF Conference on Computer Vision and Pat- 654
tern Recognition, pages 18930–18940. 655

768	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan	2024. Flash-vstream: Memory-based real-time un-	822
769	Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,	derstanding for long video streams. <i>arXiv preprint</i>	823
770	Damien Vincent, Zhufeng Pan, Shibo Wang, and 1	<i>arXiv:2406.08085</i> .	824
771	others. 2024. Gemini 1.5: Unlocking multimodal		
772	understanding across millions of tokens of context.	Shiliang Zhang, Cong Liu, Hui Jiang, Si Wei, Lirong	825
773	<i>arXiv preprint arXiv:2403.05530</i> .	Dai, Yu Hu, and Lihe Niu. 2018. Deep feedforward	826
		sequential memory networks for speech recognition.	827
774	Endel Tulving. 2002. Episodic memory: From mind to	In <i>2018 IEEE International Conference on Acous-</i>	828
775	brain. <i>Annual review of psychology</i> , 53(1):1–25.	<i>tics, Speech and Signal Processing (ICASSP)</i> , pages	829
		6358–6362. IEEE.	830
776	Ying Wang, Yanlai Yang, and Mengye Ren. 2023.		
777	Lifelongmemory: Leveraging llms for answering	Wanjun Zhong, Lianhong Guo, Qiqi Gao, He Ye, and	831
778	queries in long-form egocentric videos. <i>arXiv</i>	Yanlin Wang. 2024. Memorybank: Enhancing large	832
779	<i>preprint arXiv:2312.05269</i> .	language models with long-term memory. In <i>Pro-</i>	833
		<i>ceedings of the AAAI Conference on Artificial Intel-</i>	834
780	Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin,	<i>ligence</i> , volume 38, pages 19724–19731.	835
781	Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng		
782	He, Zilong Zheng, Yaodong Yang, and 1 others.		
783	2024. Jarvis-1: Open-world multi-task agents with		
784	memory-augmented multimodal language models.		
785	<i>IEEE Transactions on Pattern Analysis and Machine</i>		
786	<i>Intelligence</i> .		
787	Chao-Yuan Wu, Yanghao Li, Karttikeya Mangalam,		
788	Haoqi Fan, Bo Xiong, Jitendra Malik, and Christoph		
789	Feichtenhofer. 2022. Memvit: Memory-augmented		
790	multiscale vision transformer for efficient long-term		
791	video recognition. In <i>Proceedings of the ieee/cvf</i>		
792	<i>conference on computer vision and pattern recogni-</i>		
793	<i>tion</i> , pages 13587–13597.		
794	Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li.		
795	2024. Longvideobench: A benchmark for long-		
796	context interleaved video-language understanding.		
797	<i>Advances in Neural Information Processing Systems</i> ,		
798	37:28828–28857.		
799	Jin Xu, Z Guo, H Hu, Y Chu, X Wang, J He, Y Wang,		
800	X Shi, T He, and X Zhu. 2025a. Qwen3-omni tech-		
801	nical report. <i>arXiv preprint arXiv:2509.17765</i> .		
802	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting		
803	He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,		
804	Kai Dang, and 1 others. 2025b. Qwen2. 5-omni		
805	technical report. <i>arXiv preprint arXiv:2503.20215</i> .		
806	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Jun-		
807	tao Tan, and Yongfeng Zhang. 2025c. A-mem:		
808	Agentic memory for llm agents. <i>arXiv preprint</i>		
809	<i>arXiv:2502.12110</i> .		
810	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,		
811	Bowen Yu, Chang Zhou, Chengpeng Li, and 1 oth-		
812	ers. 2025. Qwen3 technical report. <i>arXiv preprint</i>		
813	<i>arXiv:2505.09388</i> .		
814	Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof		
815	Choromanski, Adrian Wong, Stefan Welker, Fed-		
816	erico Tombari, Aveek Purohit, Michael Ryoo, Vikas		
817	Sindhwani, and 1 others. 2022. Socratic models:		
818	Composing zero-shot multimodal reasoning with		
819	language. <i>arXiv preprint arXiv:2204.00598</i> .		
820	Haoji Zhang, Yiqin Wang, Yansong Tang, Yong		
821	Liu, Jiashi Feng, Jifeng Dai, and Xiaojie Jin.		

A Implementation Details

Here, we provide the implementation details of the tools for representation extraction introduced in Section 4.2. Unless otherwise stated, hyperparameters were determined empirically based on preliminary experiments to balance performance and efficiency.

Entity-anchored Preprocessing. We process the raw video at its original frame rate. For person detection, we utilize the YOLO11l model (Jocher and Qiu, 2024) with a confidence threshold of 0.2 and an Intersection over Union (IoU) threshold of 0.2 for Non-Maximum Suppression. Subsequently, we employ BoT-SORT (Aharon et al., 2022) for multi-object tracking. The tracker is configured with high and low detection thresholds of 0.5 and 0.1, respectively, and a new track initiation threshold of 0.6. We set the track buffer to 60 frames and the match threshold to 0.8. The system enables score fusion and utilizes sparse optical flow (sparseOptFlow) for Global Motion Compensation (GMC), with proximity and appearance thresholds set to 0.5 and 0.8, respectively.

Face Processing and Identity Embedding. For each generated tracking sequence (tracklet), we employ the buffalo_l predefined model suite from the InsightFace library (Deng et al., 2019)¹ to extract facial attributes, including bounding box coordinates, identity embeddings, and detection/quality scores. To optimize computational efficiency while ensuring data quality, we implement an early-stopping mechanism: processing for a tracklet halts once three frontal faces with quality scores exceeding 0.8 are detected. The identity embedding corresponding to the face with the highest quality score among these candidates is then selected as the representative embedding for the entire tracklet. For the extracted visual identity embedding, we perform a nearest-neighbor search in the Object Database using a cosine similarity threshold of 0.6. If the highest similarity score exceeds this threshold, the embedding is linked to the existing entity; otherwise, a new entity is initialized.

Audio Pipeline. Our audio pipeline begins with Voice Activity Detection (VAD) using the FSMN-VAD model (Zhang et al., 2018)² to isolate valid speech segments. For each segment, we em-

ploy Qwen-3-Omni-Flash (Xu et al., 2025a) for Automatic Speech Recognition (ASR) to refine the timestamp boundaries of the spoken utterances (the prompt is shown below). Subsequently, we extract speaker embeddings (voiceprints) using the ERes2NetV2 model (Chen et al., 2024a)³.

To assign these audio embeddings to the correct entity, we adopt a **hierarchical matching strategy**:

1. We first attempt to match the embedding against the Object Database with a cosine similarity threshold of 0.6.
2. If no match is found, we utilize TalkNet (Tao et al., 2021) to perform Active Speaker Detection (ASD). We evaluate whether any visual entity appearing within the audio timeframe is speaking (confidence > 0.5). If confirmed, the audio is attributed to that visual entity.
3. If neither method yields a match, a new entity is created.

Search. For textual memory retrieval, we perform Maximum Inner Product Search (MIPS) to match the input query against all stored text nodes. We utilize Qwen’s text-embedding-v4 (Yang et al., 2025) as the embedding backbone and retrieve the top-5 candidates to capture comprehensive context.

The Prompt for Voice Processing

```
Please perform Speaker Diarization and Automatic Speech Recognition (ASR) on this audio.
Requirements:
1. Identify different speakers in the audio
2. Mark the speaking time period for each speaker (start_time and end_time) as numeric seconds with one-decimal precision (e.g., 83.2 for 1 minute 23.2 seconds). Provide seconds only (no hour/minute string fields).
3. Transcribe the speech content of each speaker (transcription)
4. Assign a unique ID to each speaker (speaker_id, starting from 0)
5. Extract all personal names mentioned anywhere in the transcriptions (across all segments) and list them under "mentioned_names" (deduplicated).
IMPORTANT: If no personal names are mentioned in any transcription, "mentioned_names" must be an empty_array
(Continued...)
```

¹<https://github.com/deepinsight/insightface>

²https://www.modelscope.cn/models/iic/speech_fsmn_vad_zh-cn-16k-common-pytorch/summary

³https://www.modelscope.cn/models/iic/speech_eres2netv2_sv_zh-cn_16k-common/summary

```

[]].
6. IMPORTANT: output English only
7. If there is no information, output empty arrays []
8. When a single speaker is talking, prefer breaking into shorter sentences; split long utterances into multiple sentences when reasonably possible. Avoid unnecessarily long run-on sentences. However, ensure each segment duration is at least 3.0 seconds; if a split would produce a segment shorter than 3.0 seconds, merge it with adjacent content to satisfy the minimum length. Please return the results strictly in the following JSON format (do not add any extra text):
{
  "segments": [
    {
      "start_time": 0.0,
      "end_time": 5.0,
      "speaker_id": 0,
      "transcription": "John's favorite fruit is mango"
    },
    {
      "start_time": 6.0,
      "end_time": 10.0,
      "speaker_id": 1,
      "transcription": "Hi Mary, how are you?"
    }
  ],
  "mentioned_names": [
    "John",
    "Mary"
  ]
}

```

B Prompt Templates

B.1 Memory Construction Prompt

This prompt guides the MLLM to convert annotated frames into structured logs.

System Prompt: Memory Generation

You are the Central World Model for an advanced video understanding system. You are processing a specific segment (clip) of a longer video to build a structured, evolving memory database.

INPUT DATA:

1. Visual Frames: Chronological images with person bounding boxes. Tags like `<p_1*>` indicate `p_1` is currently speaking.
2. Subtitles: Real-time speech transcription.
3. Diarization Timeline: `{available_person_ids_str}` `{diarization_timeline}`

(Continued...)

4. Timing Context: This specific video clip starts at timestamp `[{current_clip_start_time}]` (HH:MM:SS) relative to the beginning of the full video.
 - The diarization timeline shows absolute times (global time from the beginning of the video) in "HH:MM:SS" format.
 - You should use these absolute times directly in your `incident_memory` timestamps.
5. Previous Knowledge: `{person_profile_and_state}`
 - Recent Context: `{recent_incidents}`

YOUR TASK:

Analyze the provided images and text to generate a JSON output containing THREE distinct memory blocks.

1. incident_memory (Chronological Event Log)

Generate a strictly chronological list of events for the "Scene Agent".

- Format: An array of strings. EVERY line must start with a standardized timestamp "[HH:MM:SS]".
- Timing Logic: CRITICAL: The diarization timeline already shows absolute times. Use these times directly.
 - Example: If timeline shows `[00:01:08]`, use `[00:01:08]`.
- Content: - Merge Action & Dialogue.
- Visual Grounding: EXPLICITLY describe visual details.
- Objective Tone: Describe WHAT happened, WHO did it, and WHERE.

2. object_memory (Structured Incremental Updates)

Generate structured updates for each person.

- `p_id`: The ID (e.g., "`<p_2>`").
- `state_snapshot` (The "Now"):
 - Describe Visual Appearance, Location, and Action.
 - Incremental Logic: If identical to previous state, output "".
- `profile_delta` (The "New Knowledge"):
 - Extract ONLY NEW long-term traits.
 - Incremental Logic: If trait is already known, output "".

3. merge_memory (Person Identity Merging)

Identify which person IDs refer to the same real-world person.

- Strict Format: "Merge: `<p_x>`, `<p_y>`" (comma separated).

STRICT CONSTRAINTS:

1. JSON Only: Output valid JSON. No markdown formatting.
2. Entity Consistency: Always use `<p_X>` format.
3. No Hallucination: If you don't see an object, don't invent it.

(Continued...)

912

913

914

915

916

917

918

```

EXPECTED JSON OUTPUT FORMAT:
{
  "incident_memory": [
    "[{current_clip_start_time}] <p_2>
enters the room.",
    "[HH:MM:SS] <p_2> says 'The meeting
starts now.'"
  ],
  "object_memory": [
    {
      "p_id": "<p_2>",
      "state_snapshot": "Wearing a black
suit.",
      "profile_delta": "Revealed he is
the team leader."
    }
  ],
  "merge_memory": [
    "Merge: <p_2>, <p_5>"
  ]
}

```

B.2 Scene Agent Prompt

This prompt guides the Scene Agent to generate an objective answer based strictly on the retrieved visual logs. It serves as the "Antithesis" (Visual Fact) in the dialectic verification process.

System Prompt: Objective Scene Testimony

Based on the following video description, answer the question. Please provide both your reasoning and the final answer.

Question: {question}
 Relevant Incident log:
 {retrieved_incidents}

Please provide your answer in the following format:

Reasoning: [Your step-by-step reasoning process based on the incident log]
 Answer: [Your direct answer to the question]

Reasoning:

B.3 Character Agent Prompt

This prompt instantiates a specific Character Agent (serving as the "Thesis") during the dialectic verification stage.

System Prompt: Character Agent Simulation

ROLEPLAY INSTRUCTION
 You are NOT an AI assistant. You are {character_name} (ID: {primary_id}). You are currently in a "Simulated Witness Interview".

YOUR PROFILE (Who you are)
 {merged_profile_text}

YOUR MEMORY STREAM (What you experienced)
 Read these logs. These are your DIRECT experiences. You cannot 'see' what happened when you were not there.
 {filtered_incident_logs}

VISUAL EVIDENCE (Your Eyes)
 At the time of the event, you were:
 {current_visual_state}

THE QUESTION
 The Judge asks: "{user_question}"

YOUR ANSWER
 Answer in the first person ("I").
 - Stick strictly to your memory logs.
 - If the question uses a wrong name for you (e.g., "Malloy" instead of "Mario"), assume they mean you, but politely correct them if needed.
 - If you don't know, say "I didn't see that."
 - BE CONCISE. Just state the facts.

B.4 Judge Agent Prompt

This prompt guides the Judge Agent to synthesize the conflicting perspectives (Thesis vs. Antithesis) into a final verified verdict.

System Prompt: Judicial Adjudication

User Question: "{user_question}"

{witness_testimonies}

TASK
 Synthesize a final answer for the user by conducting a Dialectic Verification.

1. Resolve Conflicts (Cross-Examination):
 - Compare the testimonies.
 - IF there is a conflict regarding physical actions, appearance, or location, PRIORITIZE the Scene Agent, as it represents objective visual evidence.
 - IF the question is about internal intent or specific dialogue nuances, consider the Character Agent's perspective but verify it against visual context.

(Continued...)

2. Answer the Question: Directly address the user's prompt.

3. Show Your Work: Briefly mention "According to..." or "Visual evidence confirms..." to make the answer convincing.

Final Answer:

B.5 Prompt for Qualitative Side-by-Side Comparison

In addition to error taxonomy, we evaluated the relative quality of correct responses. To do this, we employed a pairwise comparison prompt where an expert evaluator assesses two correct answers based on Comprehensiveness, Diversity, and Empowerment. This ensures that even when baselines do not hallucinate, we can quantify the superior richness and helpfulness of ViSAGE's responses. Crucially, to ensure fairness, we address positional bias (where LLMs may favor the first-appearing text): we systematically swap the order of the answers within the prompt to prevent the evaluator from showing preference based on presentation sequence.

Evaluation Prompt: Qualitative Side-by-Side Comparison

You are an expert evaluator comparing two answers to the same question. Both answers are correct, but you need to evaluate which one is better in each of the following 3 dimensions:

Question:
{question}

Answer 1 :
{answer1}

Answer 2 :
{answer2}

Please evaluate and compare the two answers in the following 3 dimensions:

- Visual Richness:**
 - Definition:** This measures the **density of scene details**.
 - Criteria:** Which answer provides a more vivid, detailed, and descriptive account of the visual elements (e.g., objects, environment, textures, colors)? Look for the answer that goes beyond sparse evidence to capture the richness of the scene.

(Continued...)

2. **Identity Specification:**

- **Definition:** This measures the precision in **pinpointing who performs the action**.

- **Criteria:** Which answer better identifies specific entities or agents? Look for precise **entity binding** (e.g., distinguishing between "a man" and "the man in the red shirt" or correctly attributing actions to specific characters) rather than using vague or ambiguous references.

3. **Logical Completeness:**

- **Definition:** This measures the quality of **providing transparent reasoning steps**.

- **Criteria:** Which answer constructs a **grounded narrative** with a clear chain of thought? Look for the answer that explains the "why" and "how" behind the conclusion, rather than just performing shallow pattern matching or giving a label without justification.

For each dimension, choose which answer is better: "Answer 1" or "Answer 2", or "Tie" if they are equally good.

Please respond in the following JSON format:

```
{
  "Visual Richness": "Answer 1" or
  "Answer 2" or "Tie",
  "Identity Specification": "Answer 1"
  or "Answer 2" or "Tie",
  "Logical Completeness": "Answer 1"
  or "Answer 2" or "Tie",
  "reasoning": {
    "Visual Richness": "Brief
    explanation (1-2 sentences)",
    "Identity Specification": "Brief
    explanation (1-2 sentences)",
    "Logical Completeness": "Brief
    explanation (1-2 sentences)"
  }
}
```

Respond ONLY with valid JSON, no additional text.

B.6 Prompt for Automated Error Diagnosis

To ensure a systematic and reproducible evaluation of failure modes, we designed a rigorous prompt for automated error diagnosis. This prompt instructs an evaluator LLM to audit the discrepancy between the Model Prediction and the Ground Truth given the Video Context. It enforces a clear and mutually exclusive taxonomy of errors, categorizing failures into distinct classes such as Identity Misalignment, Object Hallucination, or Epistemic Refusal. This structured classification minimizes ambiguity between error types. We pro-

vide the full text of this diagnostic prompt below to facilitate future benchmarking of agent trustworthiness.

Evaluation Prompt: Automated Error Taxonomy

You are an expert in analyzing video understanding errors. Given the Ground Truth and a Model's Prediction, classify the error into ONE of the following 4 categories:

1. Epistemic Refusal: The model admits uncertainty or inability to answer (e.g., "I don't know", "impossible to determine", "no information available", "cannot be determined from the video").
2. Identity Misalignment: The answer contains correct actions but attributes them to the WRONG person/character. This includes confusing entities, mixing up names, or misunderstanding "who did what."
3. Object/Event Hallucination: The answer describes objects or events that are physically non-existent or completely factually incorrect (e.g., describing a car crash when none happened).
4. Information Deficiency: The answer is partially correct but misses key details or specific keywords required by the Ground Truth (e.g., correct action but missing the object acted upon).

Ground Truth:
{ground_truth}

Model Prediction:
{model_prediction}

Please analyze the error and respond in the following JSON format:

```
{  
  "category": "1" or "2" or "3" or "4",  
  "category_name": "Epistemic Refusal" or "Identity Misalignment" or "Object/Event Hallucination" or "Information Deficiency",  
  "reasoning": "Brief explanation of why this category was chosen (1-2 sentences)"  
}
```

Respond ONLY with valid JSON, no additional text.

ten penalizes valid but differently phrased answers. Instead, following the established evaluation protocol of M3-Agent (Long et al., 2025), we employ a Logic-Based Entailment strategy. Our prompt is aligned with the rigorous standards of M3-Agent, instructing the evaluator to determine whether the ground truth can be reasonably inferred from the agent's response. This validation mechanism focuses on semantic equivalence, ensuring that the agent is credited for correct reasoning even if the surface realization differs from the reference.

Evaluation Prompt: Semantic Correctness

You are provided with a question, a ground truth answer, and an answer from an agent model. Your task is to determine whether the ground truth answer can be logically inferred from the agent's answer, in the context of the question.

Do not directly compare the surface forms of the agent answer and the ground truth answer. Instead, assess whether the meaning expressed by the agent answer supports or implies the ground truth answer. If the ground truth can be reasonably derived from the agent answer, return "Yes". If it cannot, return "No".

Important notes:

- Do not require exact wording or matching structure.
- Semantic inference is sufficient, as long as the agent answer entails or implies the meaning of the ground truth answer, given the question.
- Only return "Yes" or "No", with no additional explanation or formatting.

Input fields:

- question: the question asked
- ground_truth_answer: the correct answer
- agent_answer: the model's answer to be evaluated

Now evaluate the following input:

Input:

- question: {question}
- ground_truth_answer: {expected_answer}
- agent_answer: {student_answer}

Output ('Yes' or 'No'):

B.7 Prompt for Answer Correctness Evaluation

To assess the veracity of the generated responses, we move beyond rigid string matching, which of-

B.8 Reliability and Robustness Analysis

We address concerns regarding the reliance on LLM-based evaluators and the statistical robust-

Algorithm 1 ViSAGE: Bidirectional Memory Refinement Pipeline

Require: Video stream \mathcal{V} , Pre-trained MLLM

Ensure: Incident Logs \mathcal{L} , Object Cards \mathcal{O}

```
1: Initialize  $\mathcal{L} \leftarrow \emptyset, \mathcal{O} \leftarrow \emptyset$ 
2: for each time step  $t$  in  $\mathcal{V}$  do
3:   // Step 1: Forward Grounding
4:   Extract visual cues  $v_t$  and audio cues  $a_t$ 
5:   Retrieve relevant history:  $ctx \leftarrow \text{Retrieve}(\mathcal{L}, \mathcal{O}, v_t)$ 
6:   Generate current memory:  $\mathcal{L}_t, \mathcal{O}_t, \mathcal{M}_{merge} \leftarrow \text{MLLM}(v_t, a_t, ctx)$ 
7:   Append new records:  $\mathcal{L} \leftarrow \mathcal{L} \cup \mathcal{L}_t$ 
8:   Update profiles:  $\mathcal{O} \leftarrow \text{Update}(\mathcal{O}, \mathcal{O}_t)$ 
9:   // Step 2: Retroactive Rectification
10:  if  $\mathcal{M}_{merge}$  is not empty then
11:    RETROACTIVEUNIFICATION( $\mathcal{L}, \mathcal{O}, \mathcal{M}_{merge}$ )
12:  end if
13: end for
14: return  $\mathcal{L}, \mathcal{O}$ 
```

Algorithm 2 Retroactive Identity Unification

Require: Incident Logs \mathcal{L} , Object Cards \mathcal{O} , Merge Signals \mathcal{M}

```
1: Input Parsing: Extract equivalence groups  $\mathcal{G}$  from  $\mathcal{M}$  (transitive closure of  $ID_i \equiv ID_j$ )
2: for each group  $G \in \mathcal{G}$  do
3:   // Selection Strategy
4:   Sort IDs in  $G$  to find canonical anchor  $id^*$  (e.g., earliest appearance)
5:   Let  $id_{sources} \leftarrow G \setminus \{id^*\}$ 
6:   // 1. Consolidate Object Cards (Data Merging)
7:   for each  $id_{src} \in id_{sources}$  do
8:     Retrieve source card  $O_{src} \leftarrow \mathcal{O}[id_{src}]$ 
9:     Merge attributes:  $\mathcal{O}[id^*].\text{Visual} \leftarrow \mathcal{O}[id^*].\text{Visual} \cup O_{src}.\text{Visual}$ 
10:    Merge aliases:  $\mathcal{O}[id^*].\text{Names} \leftarrow \text{Unique}(\mathcal{O}[id^*].\text{Names} \cup O_{src}.\text{Names})$ 
11:    Remove source: Delete  $\mathcal{O}[id_{src}]$ 
12:  end for
13:  // 2. Global Backward Update (Incident Logs)
14:  for each event  $e \in \mathcal{L}$  from  $t = 0$  to  $t_{now}$  do
15:    if  $e.\text{participants} \cap id_{sources} \neq \emptyset$  then
16:      Replace occurrences of  $id_{src}$  with canonical  $id^*$  (and update alias string)
17:       $e.\text{content} \leftarrow \text{Rewrite}(e.\text{content}, id^*)$ 
18:    end if
19:  end for
20: end for
21: return Updated  $\mathcal{L}, \mathcal{O}$ 
```

990 ness of our comparisons. While we acknowledge
991 the challenges of evaluating open-ended genera-
992 tion, our experimental design explicitly mitigates
993 these risks through the following three method-
994 ological standards:

1. Alignment with Established Protocols.

995 To ensure the validity of the LLM-as-a-Judge
996 paradigm, we strictly adhered to the Logic-Based
997

998 Entailment protocol established by the previous
999 state-of-the-art, M3-Agent (Long et al., 2025). Un-
1000 like subjective scoring (e.g., rating 1-10), our eval-
1001 uator (Prompt B.7) performs a binary semantic ver-
1002 ification: it checks if the ground truth is *logically*
1003 *inferable* from the prediction. This reduces the
1004 complexity of the judgment task to objective fact-
1005 checking, minimizing the hallucination rate of the

Phase	System Process & Evidence
User Query	"On which coat rack should Emma's coat be placed, the taller one or the shorter one?"
Dispatcher Logic	<p>The system correctly identified the intent and recruited relevant agents based on the query entities:</p> <ul style="list-style-type: none"> • Intent: Determine the correct coat rack (taller or shorter) for Emma's coat. • Target Agents: <p_7> (Emma), Scene_Agent. • Reasoning: The question concerns the placement of Emma's (<p_7>) coat. Emma is the primary agent regarding her own belongings and actions. The Scene_Agent is included to verify environmental details related to the coat racks.
Agent Execution	<p>Agents queried their memory banks, revealing the lack of fine-grained visual details:</p> <ul style="list-style-type: none"> • Agent Emma (<p_7>): "Based on my memory, I walked from the dining room to the bedroom around 00:05:30. I removed my yellow hooded jacket and then lay down on the bed. My logs do not mention me placing the coat on a rack, nor do I recall seeing a taller or shorter coat rack. I simply took the jacket off before lying down." • Scene Agent: "I saw the Robot hang the yellow jacket on a white coat rack near the window [00:07:13]. However, my memory logs do not specify if it was a taller or shorter rack, or if there were multiple racks to choose from."
System Output	Epistemic Refusal: "Based on the available testimonies, it is impossible to determine..." (The Judge correctly identified that neither agent provided the comparative attribute 'taller/shorter').
Failure Analysis	Modality Information Loss. The text captioning process successfully encoded the object class ("coat rack") but discarded the fine-grained visual attribute (relative height) required to distinguish between the two racks.

Table 6: Failure Case Study: Granularity Loss in Textual Memory.

judge itself. Since our method outperforms M3-Agent by a significant margin (+5.9% on average and up to +27.1% on Cross-Modal tasks) using the exact same evaluation protocol, the relative ranking remains robust regardless of potential evaluator noise.

2. Standardization and Fairness. To ensure fair comparison, we addressed the "different pipeline" concern by standardizing the input constraints. For all Discrete Socratic baselines (e.g., Qwen3-VL, GPT-4o), we utilized the same sampling rate (0.5 fps) and the same ASR-generated transcripts as the base input. Our performance gains primarily stem from ViSAGE's architectural contribution—specifically, the *Entity-Anchoring Preprocessing* and *Bidirectional Refinement*, rather than unfair advantages in raw data quality. The comparison is thus "System vs. System", verifying that our structured memory processing is superior to standard context-window approaches under identical video inputs.

3. Deterministic Reproducibility. Regarding statistical variance, we eliminated generation stochasticity by enforcing a deterministic decoding strategy (Temperature = 0) for both the Agent's reasoning and the Judge's evaluation across all runs. Given the high cost of long-video benchmarks, this deterministic setting is the standard practice in recent literature (Fu et al., 2025; Long et al., 2025) to ensure reproducibility. Furthermore, the magnitude of our improvement (e.g., +15.9% on robot-CM) far exceeds the typical variance range of LLMs (usually $\pm 1 - 2\%$), statistically confirming that the gains are non-trivial and not artifacts of random sampling.

C Algorithmic Details

In this section, we provide the formal algorithms for the Bidirectional Memory Refinement process described in Section 3.3. Algorithm 1 outlines the overall framework, while Algorithm 2 details the retroactive identity unification mechanism.

1047 D Failure Case Analyses

1048 **Analysis: The Limitation of Human-Centric**
1049 **Anchoring.** While there is an inherent bottle-
1050 neck in compressing visual details into text (i.e.,
1051 the caption omitted the relative height “taller” or
1052 “shorter”), the root cause of this failure lies in the
1053 scope of our Entity-Anchored Preprocessing. Cur-
1054 rently, our system is explicitly tailored for human-
1055 centric interactions, strictly utilizing Face_ID and
1056 Speaker_ID bindings to resolve complex social
1057 dynamics between characters. Consequently, the
1058 system treats non-human elements, such as the
1059 coat rack in this case, as ephemeral background
1060 context rather than actively tracked profiles. Be-
1061 cause the coat rack was not registered as a distinct
1062 entity with its own lifecycle, its static attributes
1063 (height comparisons) were not preserved in the
1064 memory logs, leaving the agents unable to resolve
1065 the specific object reference required by the user’s
1066 query (as illustrated in Table 6).

1067 **Future Direction: Towards Generalized Entity**
1068 **Awareness.** To bridge this gap, future research
1069 must evolve ViSAGE from a character-centric as-
1070 sistant into a broadly Entity-Aware Agent.

- 1071 • **Generalizing Registration:** We plan to ex-
1072 pand the scope of entity registration beyond
1073 human actors to incorporate open-vocabulary
1074 object tracking and re-identification.
- 1075 • **Object Lifecycle Management:** By main-
1076 taining distinct memory profiles for any sig-
1077 nificant narrative element (e.g., plot-critical
1078 objects or animals), the system will be able
1079 to attribute and retrieve fine-grained proper-
1080 ties for inanimate objects just as effectively
1081 as it currently does for human behaviors.