# Agentic Learner with Grow-and-Refine Multimodal Semantic Memory

Weihao Bo[1,2]    Shan Zhang[3]    Yanpeng Sun[4*]    Jingjing Wu[2]    Qunyi Xie[2]    Xiao Tan[2]
Kunbin Chen[2]    Wei He[2]    Xiaofan Li[2]    Na Zhao[4]    Jingdong Wang[2]    Zechao Li[1†]

[1]Nanjing University of Science and Technology,    [2]Baidu Inc,
[3]Adelaide AIML,    [4]Singapore University of Technology and Design

## Abstract

*MLLMs exhibit strong reasoning on isolated queries, yet they operate* de novo—*solving each problem independently and often repeating the same mistakes. Existing memory-augmented agents mainly store past trajectories for reuse. However, trajectory-based memory suffers from brevity bias, gradually losing essential domain knowledge. More critically, even in truly multimodal problem-solving settings, it records only a **single-modality** trace of past behavior, failing to preserve how visual attention and logical reasoning jointly contributed to the solution. This is fundamentally misaligned with human cognition: semantic memory is both **multimodal and integrated**, preserving visual and abstract knowledge through coordinated but distinct representational streams. We thus introduce **ViLoMem**, a dual-stream memory framework that constructs compact, schema-based memory. It separately encodes visual distraction patterns and logical reasoning errors, enabling MLLMs to learn from their successful and failed experiences. Following a grow-and-refine principle, the system incrementally accumulates and updates multimodal semantic knowledge—preserving stable, generalizable strategies while avoiding catastrophic forgetting. Across six multimodal benchmarks, **ViLoMem** consistently improves pass@1 accuracy and substantially reduces repeated visual and logical errors. Ablations confirm the necessity of dual-stream memory with explicit distraction–hallucination separation, demonstrating the value of error-aware multimodal memory for lifelong and cross-domain agentic learning. Our project page will be available at* [https://weihao-bo.github.io/ViLoMeo-page/](https://weihao-bo.github.io/ViLoMeo-page/).

## 1. Introduction

Multimodal Large Language Models (MLLMs) have achieved impressive progress in scene understanding, visual question answering, and complex scientific problem solving [5, 38, 50, 52]. Yet despite their growing capability, current MLLMs approach each problem *de novo*—solving every query in isolation, repeatedly re-deriving the same insights and re-committing familiar errors[13, 14, 28, 51]. Although recent memory-augmented models attempt to mitigate this by storing past interactions [27, 49], these memories capture only high-level logical summaries while discarding the visual grounding and perceptual cues essential for multimodal reasoning.

Recent research has demonstrated that MLLMs' visual perception ability remains fundamentally weaker than their linguistic reasoning, with low-level perceptual failures identified as a primary bottleneck for high-level multimodal reasoning tasks [20, 26, 30]. In mathematical multimodal problem-solving in particular, diagram-perception errors exceed logical reasoning errors, and visual mistakes frequently persist in intermediate reasoning steps even when the final answer is correct [45]. This indicates visual attention errors directly cause downstream logical hallucinations that creates a cascading failure pattern [36, 53]. Our ablation studies further confirm this phenomenon: across six multimodal problem-solving benchmarks, the proportion of visual error summaries consistently exceeds that of logical memory errors (Fig. 4). Therefore, when solving problems paired with images, it is essential for models to maintain accurate visual attention to task-relevant regions, avoiding perceptual distractions that propagate into flawed logical inferences.

Logic-only memory is insufficient for multimodal problem solving. While logical theorems and rules are general (e.g., applying the base–height formula for area computation), effective reasoning also requires aligning these abstract rules with their correct visual counterparts (e.g., the shape of triangles). As illustrated in Fig. 1, triangles exhibit diverse visual configurations, and early attempts may contain both logical and visual errors. Through feedback, the model refines its logical memory for question-appropriate theorem application and its visual memory to avoid percep-
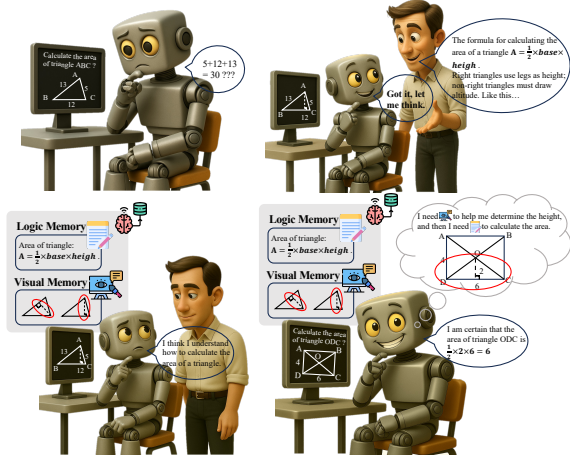
---

Figure 1. Multimodal Semantic Memory Enables Progressive Learning. When solving multimodal problems, early attempts may contain both logical and visual errors; through feedback, the model refines its logical memory for question-appropriate theorem application and its visual memory to avoid perceptual traps—improving by integrating the *where to look* with the *how to reason.*

tual traps, attending to task-relevant regions. This progressive learning mirrors the human cognitive system, where semantic memory maintains *multimodal representations* that integrate visual experience with abstract reasoning [34].

We thus introduce **ViLoMem**, a dual-stream memory framework that separately models visual distraction patterns and logical hallucination errors as structured schemas, coordinating them through unified retrieval. Following a grow-and-refine principle, **ViLoMem** avoids the detail erosion caused by iterative rewriting by filtering similar error patterns and using tailored add/skip and retrieval strategies to incrementally accumulate multimodal semantic knowledge. Specifically, we design custom retrieval strategies for visual and logical streams. For the visual stream, direct image-similarity search is insufficient; the key requirement is helping the model identify question-specific "visually trapped regions". To achieve question-aware attention, we generate cross-modal attention maps guided by keywords (previously observed visual mistakes), enabling the model to highlight regions associated with known error patterns relevant to the current question. For the logical stream, instead of directly retrieving query semantically similar logics, the model first analyzes the problem to identify its underlying subject and reasoning requirements—supporting precise positioning of the task type and precise selection of the relevant logical schema.

Overall, **ViLoMem** automatically attributes successes or failures to the visual or logical stream and updates the corresponding schemas without human supervision. It enables progressive mistake reduction and cross-domain knowledge transfer in multimodal tasks. Our contributions are summarized as follows:

1. We propose **ViLoMem**, the first framework that explicitly separates *visual distraction errors* from *logical hallucination errors*, constructing a dual-stream memory inspired by multimodal semantic memory in the human cognitive system.

2. **ViLoMem** employs a question-aware attention mask for visual images and a *precise-positioning–precise-selection* regime for logical retrieval, together with filtering-based memory update strategies to avoid detail erosion, enabling coordinated retrieval between visual cues and logical constraints. This directs attention to task-relevant regions while suppressing invalid logical inferences.

3. Extensive experiments on six multimodal benchmarks demonstrate that **ViLoMem** consistently improves pass@1 accuracy across diverse model scales, achieving substantial gains on mathematical reasoning tasks (e.g., +6.48 on MathVision for GPT-4.1, +4.38 on MMMU for Qwen3-VL-8B). Ablation studies confirm that both memory streams are essential and complementary, exhibiting heterogeneous effects across benchmarks–task-aligned domains benefit from shared memory, whereas mismatched domains can lead to interference.

## 2. Related Work

### 2.1. Context Engineering

Recent advancements in agent self-improvement have prominently featured *context engineering*, a paradigm that refines model behavior by strategically modifying input prompts rather than altering the model's underlying weights[1, 8, 24, 37]. These methods primarily leverage natural language feedback, enabling a model to analyze its own performance based on execution traces, reasoning steps, or validation signals and then iteratively revise its operational context [2, 25, 32, 44]. This approach has given rise to several influential frameworks. For instance, Re-Act [42] pioneered the integration of reasoning and acting within a synergistic loop. Building on this, Reflexion [25] introduced a mechanism for agents to reflect on past failures, using verbal reinforcement to enhance subsequent planning and decision-making. Other works have focused on optimizing the prompts themselves; TextGrad [44] proposed a novel method to generate gradient-like textual feedback for prompt refinement, while GEPA [2] demonstrated that an evolutionary approach to prompt optimization based on execution traces can achieve performance surpassing that of traditional reinforcement learning in certain scenarios. However, these approaches are limited by their ephemeral nature; the context is constructed for single interactions, preventing long-term knowledge accumulation. Furthermore, they often suffer from a *brevity bias* [15], where iterative refinement strips away crucial details, hindering performance on complex, knowledge-intensive tasks.

## 2.2. Long-term Memory

To address the limitations of transient context, a parallel line of research has focused on equipping agents with *long-term memory*, enabling them to learn from experience and retain knowledge persistently[3, 11, 23, 40, 51]. This vision is rooted in the cognitive science principle that true, deep learning extends beyond formal training and arises from the continuous accumulation of experience [6, 12, 19, 31, 35]. Research in this area explores various architectures for building durable memory systems. For example, Dynamic Cheatsheet [27] constructs an external memory that explicitly stores successful and unsuccessful strategies from past inferences, allowing the agent to consult its history. Similarly, ACE [49] develops an incremental "context playbook" through a generate-reflect-curate cycle, which is designed to avoid the simplification and catastrophic forgetting associated with simple iterative rewriting. The mechanisms for populating these memories are also diverse, ranging from learning through early, formative experiences [48] and reinforcement learning-based exploration [46] to interactive learning from noisy, real-time human feedback [4, 41].

However, these frameworks exhibit a critical blind spot: they are overwhelmingly logic-centric, capturing reasoning patterns while neglecting the visual dimension of multimodal tasks. In contrast, the human brain adopts a hub-and-spoke semantic memory architecture. Visual–semantic associations and error patterns are encoded in the inferotemporal and perirhinal cortex (visual spoke), while abstract reasoning rules and logical error patterns are maintained in the temporal–parietal cortex (logic spoke)[9, 17, 18]. The anterior temporal lobe (ATL) serves as the central hub that integrates these modality-specific representations into unified conceptual knowledge. Inspired by this architecture, our AI system implements an *error-aware multimodal semantic memory*, where visual and logical error patterns are stored in separate modality-specific modules, integrated through a semantic hub, and monitored by an executive verifier that detects redundant visual–logical information and modulates attention to prevent recurring mistakes in multimodal scientific reasoning tasks.

## 3. Method

We propose **ViLoMem**, a plug-in dual-stream memory framework for multimodal reasoning in large language models, featuring a closed-loop *Memory Cycle* that enables the agent to continuously learn from its reasoning and perception errors—facilitating progressive, lifelong learning.

**Problem Formulation.** Consider a sequence of multimodal inputs $(x_1, x_2, \ldots, x_n)$, where each input $x_i = (I_i, q_i)$ consists of an image $I_i$ and a question text $q_i$. The system maintains two memory banks: a logic memory $\mathcal{M}_i^L = \{m_1^L, m_2^L, \ldots, m_{|L|}^L\}$ storing textual reasoning guidelines, and a visual memory $\mathcal{M}_i^V = \{(m_1^V, I_1^V), (m_2^V, I_2^V), \ldots, (m_{|V|}^V, I_{|V|}^V)\}$ storing visual guidelines paired with source images.

As illustrated in Figure 2(a), the cycle operates as follows: given problem $x_i$, the system performs parallel *Retrieval* from both memory banks to obtain relevant memories $R_i^L$ and $R_i^V$. These retrieved memories are then fed to the **Solver** for *Utilization*, which generates a candidate answer $\tilde{y}_i$. The **Verifier** evaluates this answer against the ground truth $y_i$. Upon detecting an error ($\tilde{y}_i \neq y_i$), the system activates the *Generation* process to update both memory banks in parallel, yielding $\mathcal{M}_{i+1}^L$ and $\mathcal{M}_{i+1}^V$. This mechanism enables the agent to progressively refine its capabilities through iterative self-correction.

**Core Operations.** We define several key operations used throughout the framework. Let $\phi^T(\cdot)$ and $\phi^M(\cdot)$ denote text and multimodal embedding functions, respectively. The cosine similarity between two embeddings is computed as:

$$\text{Sim}(u, v) = \frac{u \cdot v}{\|u\|\|v\|} \tag{1}$$

For problem analysis during retrieval, we employ an LLM to extract structured information from the question and reasoning trace:

$$a_i = \text{Analyze}^L(q_i, \tilde{y}_i) \tag{2}$$

The process identifies the problem's subject domain and key concepts. An enriched query is then constructed by combining the original question with this analysis:

$$\tilde{q}_i = [q_i; a_i] \tag{3}$$

### 3.1. Memory Generation

When errors are detected, the system activates a parallel memory-generation framework, as illustrated in Figure 2(b). This framework conducts detailed error attribution and constructs structured memory units corresponding to two distinct error types.

### 3.1.1. Visual Memory Generation

The visual analysis module, powered by an MLLM, simultaneously identifies the error type and generates corrective guidance. Given the original image $I_i$, question $q_i$, erroneous reasoning trace $\tilde{y}_i$, and ground truth $y_i$, the module produces both an error indicator and a corresponding guideline within a single model invocation, formally expressed as:

$$(e_i^V, g_i^V) = \text{AnalyzeGenerate}^V(I_i, q_i, \tilde{y}_i, y_i), \tag{4}$$

where $e_i^V \in \text{True}, \text{False}$ indicates whether the error originates from visual misinterpretation (e.g., object confusion, overlooked visual symbols, or spatial relationship misunderstandings), and $g_i^V$ denotes the generated *Visual*
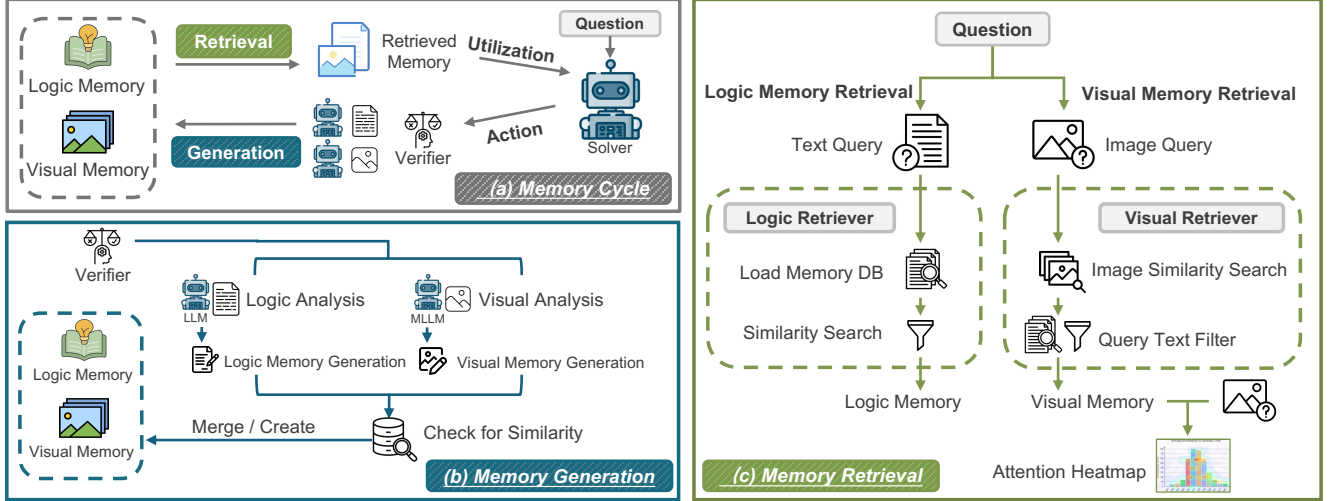
3

Figure 2. Overview of the **ViLoMem** framework. (a) *Memory Cycle*: A closed-loop learning mechanism where both logical and visual memories are retrieved and utilized by the solver. Retrieval is conditioned on the textual question and its paired image. The solver then performs reasoning steps (actions), which are evaluated by the verifier to filter redundant or invalid trajectories. The remaining trajectories are used to update both memory streams according to their respective types. (b) *Memory Generation*: An error-attribution framework that employs an LLM for logical analysis and an MLLM for visual analysis, producing structured memory schemas through similarity-based **merge** and **create** operations. (c) *Memory Retrieval*: Specialized dual-stream retrieval mechanism. Visual memories undergo a two-stage process involving image-embedding retrieval followed by question-specific retrieval, since visual information must be conditioned on both image content and the textual query. Logical memories are retrieved through problem analysis and text-embedding similarity.

*Guideline*—an instruction prescribing the correct observation strategy. All information is stored in a structured JSON dictionary for persistent memory updating. For example, when addressing shape and attribution-related errors in 3D solid objects, the guideline may state:

> "When an object has a uniform, reflective, or metallic-looking surface—even if it appears matte under diffuse lighting—treat it as metallic if it matches the visual style of other known metallic objects in the scene."

Before storage, a *similarity check* is performed against existing memories in $\mathcal{M}_i^V$ using text embeddings. The system computes similarity scores $s_j^V = \mathrm{Sim}(\phi^T(g_i^V), \phi^T(m_j^V))$ for all $m_j^V \in \mathcal{M}_i^V$. If $\max_j s_j^V > \tau^V$ (where $\tau^V$ is a similarity threshold), a *merge* operation consolidates the knowledge:

$$\mathcal{M}_{i+1}^V = \mathcal{M}_i^V \setminus \{(m_{j^*}^V, I_{j^*}^V)\} \cup \{(\mathrm{Merge}^V(m_{j^*}^V, g_i^V), I_{j^*}^V)\}, \tag{5}$$

where $j^* = \arg\max_j s_j^V$. Otherwise, a new memory entry is created: $\mathcal{M}_{i+1}^V = \mathcal{M}_i^V \cup \{(g_i^V, I_i)\}$.

### 3.1.2. Logical Memory Generation

In parallel, the logic analysis module, powered by an LLM, examines the reasoning chain for non-visual errors such as computational mistakes, formula misapplications, or logical fallacies. This module focuses solely on textual reasoning without accessing visual information. As formalized in Equation (6), the module produces both error classification and guideline in a single model invocation:

$$(e_i^L, g_i^L) = \mathrm{AnalyzeGenerate}^L(q_i, \tilde{y}_i, y_i), \tag{6}$$

where $e_i^L \in \{\mathrm{Logical}, \mathrm{Non\text{-}Logical}\}$ classifies whether the error involves reasoning failures, and $g_i^L$ represents the abstracted *Logic Guideline*. The model outputs a structured text response containing error type, analysis summary, and guideline fields. For example, when encountering a geometry error arising from incorrect assumptions (i.e., textual biases), the generated guideline may state:

> "In geometry problems involving perpendicular bisectors, remember that only points lying on the perpendicular bisector segment are guaranteed to be equidistant from the endpoints of the segment. Do not assume a point lies on the bisector unless this is explicitly stated or can be proven from the given construction. Always verify the position of intersection points relative to the bisector before applying the equidistance property."

This guideline then undergoes the same *similarity check* and *merge/create* process as visual memory. Similarity scores $s_j^L = \mathrm{Sim}(\phi^T(g_i^L), \phi^T(m_j^L))$ are computed for all $m_j^L \in \mathcal{M}_i^L$, and the memory bank is updated accordingly:

$$\mathcal{M}_{i+1}^L = \begin{cases} \mathcal{M}_i^L \setminus \{m_{j^*}^L\} \cup \{m_{\mathrm{new}}^L\} & s_{j^*}^L > \tau^L \\ \mathcal{M}_i^L \cup \{g_i^L\} & s_{j^*}^L \leq \tau^L \\ \mathcal{M}_i^L & \text{otherwise}, \end{cases} \tag{7}$$

where $j^* = \arg\max_j s_j^L$, $m_{\text{new}}^L = \text{Merge}^L(m_{j^*}^L, g_i^L)$, and the update is triggered when $e_i^L = \text{Logical}$ and $g_i^L \neq \emptyset$.

## 3.2. Memory Retrieval and Utilization

When addressing a new problem $x_i = (I_i, q_i)$, the solver initiates parallel retrieval procedures from both memory banks, as illustrated in Figure 2(c). Unlike conventional single-stage retrieval, our framework employs specialized strategies for each memory type: visual memory uses a two-stage multimodal-to-text pipeline, while logical memory leverages problem analysis to construct enriched queries.

### 3.2.1. Visual Memory Retrieval

Visual memory retrieval employs a two-stage pipeline that progressively refines candidates from visual similarity to semantic relevance.

**Stage 1: Image Embedding Similarity.** The system first employs multimodal embeddings to compute visual similarity between the query image $I_i$ and all stored memory images. For each memory $(m_j^V, I_j^V) \in \mathcal{M}_i^V$, the similarity is computed as $s_j^M = \text{Sim}(\phi^M(I_i), \phi^M(I_j^V))$. This rapidly recalls a set of top-$k^M$ candidate memories:

$$\mathcal{C}_i^V = \{(m_j^V, I_j^V) \mid j \in \text{TopK}(\{s_j^M\}, k^M)\} \qquad (8)$$

**Stage 2: Text Embedding Filtering .** Visual similarity alone is insufficient for semantic matching. The system subsequently performs text-based reranking using the enriched query $\tilde{q}_i$ from Equation (3). For each candidate guideline $m_j^V \in \mathcal{C}_i^V$, text similarity is computed as $s_j^T = \text{Sim}(\phi^T(\tilde{q}_i), \phi^T(m_j^V))$. The final retrieved visual memories are obtained by filtering with threshold $\tau^V$ and selecting top-$k^V$ by similarity score:

$$R_i^V = \{m_j^V \mid j \in \text{TopK}(\{s_j^T \mid s_j^T \geq \tau^V\}, k^V)\} \qquad (9)$$

This two-stage process ensures that the retrieved visual memories are both semantically relevant to the current problem and specifically address common visual pitfalls encountered when interpreting similar images.

**Focusing on where to look via visual attention maps.** Beyond textual guidelines, we further introduce an auxiliary visual representation of memory cues. Leveraging the retrieved visual memory and its associated error patterns, the system generates question-aware attention maps that highlight historically error-prone regions in the query image $I_i$. These attention maps serve as supplementary visual inputs alongside the original image, providing explicit spatial guidance that directs the model's focus toward task-relevant areas while avoiding known perceptual traps. Experimental results demonstrate that this visual augmentation yields additional performance improvements (refer to Section 4.4).

### 3.2.2. Logical Memory Retrieval

Logical memory retrieval is a text-based semantic matching process. The system constructs an enriched query $\tilde{q}_i$ using Equations 2-3 to capture both the problem text and structured domain information. For each memory $m_j^L \in \mathcal{M}_i^L$, text embedding similarity is computed as $s_j^L = \text{Sim}(\phi^T(\tilde{q}_i), \phi^T(m_j^L))$. The top-$k^L$ most relevant guidelines are retrieved by applying similarity threshold $\tau^L$ and ranking by similarity score:

$$R_i^L = \{m_j^L \mid j \in \text{TopK}(\{s_j^L \mid s_j^L \geq \tau^L\}, k^L)\} \qquad (10)$$

### 3.2.3. Solution Generation with Dual Memory

Finally, the solver generates the answer by conditioning on both the original inputs and the retrieved memories from the visual and logical streams:

$$\tilde{y}_i = \text{Gen}(I_i, q_i, R_i^L, R_i^V), \qquad (11)$$

where Gen denotes the MLLM solver that integrates visual perception, question understanding, and dual-stream memory guidance. The retrieved logical guidelines $R_i^L$ provide structured and context-relevant reasoning frameworks, while the visual guidelines $R_i^V$ supply explicit perceptual priors. Together, they enable more robust and accurate multimodal reasoning.

## 4. Experiments

### 4.1. Experimental Setup

**Tasks and Datasets.** We evaluate **ViLoMem** on three multimodal reasoning benchmarks that are particularly sensitive to cumulative visual–logical errors: (1) *Hallucination and real-world robustness*, which emphasize language hallucination, visual illusion, and spatial grounding; (2) *Multimodal mathematical reasoning*, which couples logic reasoning with visual grounding; and (3) *Vision-dependent knowledge*, which requires expert-level visual understanding across multiple disciplines.

HallusionBench [16] diagnoses intertwined language hallucination and visual illusion through 1,129 control-paired questions; RealWorldQA [39] assesses spatial reasoning over 765 natural scenes; MathVista (mini) [21] and MathVision (mini) [33] test visual-grounded mathematical reasoning across diverse diagrams and competition-style problems; MMMU (val) [43] covers 1050 college-level questions across six academic domains (Art & Design, Business, Science, Health & Medicine, Humanities & Social Science, Tech & Engineering); and MMStar [7] offers 1,500 high-quality samples evaluating vision-dependent reasoning across 18 fine-grained dimensions.

**Models and Implementation.** To assess the effectiveness and generalizability of **ViLoMem**, we evaluate it across models of varying scale and accessibility: the proprietary

Table 1. Main results across six multimodal reasoning benchmarks. Baseline metrics for Qwen3 series models are sourced from official reports, while GPT-4.1 baselines are from OpenCompass. Metrics marked with * indicate self-evaluated results where official reports are unavailable or show substantial discrepancies. Models with "(step)" and "(+ ViLoMem)" are prompted by step-by-step reasoning.

| Method | MMMU | MathVista | MathVision | HallusionBench | MMStar | RealWorldQA |
|---|---|---|---|---|---|---|
| GPT-4.1 (baseline) | 74.00 | 70.40 | 46.12* | 58.50 | 69.80 | 73.72 |
| GPT-4.1 (step) | 74.16 | 74.27 | 47.47 | 74.44 | 70.43 | 72.03 |
| GPT-4.1 (+ ViLoMem) | **77.26** | **76.88** | **53.95** | **75.29** | **72.43** | **74.38** |
| Qwen3-VL-235B-A22B-Instruct (baseline) | 78.70 | 84.90 | 61.28* | 63.20 | 78.40 | 79.30 |
| Qwen3-VL-235B-A22B-Instruct (step) | 75.97 | 83.66 | 62.17 | 74.58 | 76.16 | 78.66 |
| Qwen3-VL-235B-A22B-Instruct (+ ViLoMem) | **79.40** | **84.98** | **62.83** | **75.21** | <u>78.31</u> | 77.22 |
| Qwen3-VL-8B-Instruct (baseline) | 66.38* | 77.20 | 48.13* | 61.10 | 70.91 | 71.50 |
| Qwen3-VL-8B-Instruct (step) | 65.52 | 77.80 | 48.35 | 73.08 | 70.22 | 70.85 |
| Qwen3-VL-8B-Instruct (+ ViLoMem) | **69.90** | **77.87** | **49.34** | **73.19** | **72.13** | **73.59** |



Figure 3. Visual memory generation and retrieval examples. Each case shows the original error, the extracted visual pattern, and successful retrieval in analogous scenarios.

GPT-4.1 as a strong closed-source baseline, the open-source Qwen3-VL-235B-A22B-Instruct as a state-of-the-art large multimodal model, and Qwen3-VL-8B-Instruct as a smaller model to test whether memory benefits extend to resource-constrained settings. For memory generation, we employ Qwen3-235B-A22B-Instruct for logical memory (pure language reasoning analysis) and Qwen3-VL-235B-A22B-Instruct for visual memory (image-grounded error attribution). Memory retrieval uses Qwen3-Embedding for text similarity and Qwen2.5-VL-Embedding for image similarity, enabling efficient semantic matching. Additional implementation details are provided in the Appendix.

**Evaluation Metrics.** We report pass@1 accuracy using VLMEvalKit [22]. When rule-based matching detects potential errors, we apply an LLM-as-a-judge mechanism for

Table 2. Ablation study on the contribution of dual stream memory components. We evaluate GPT-4.1 with different memory configurations on two representative benchmarks.

| Method | MMMU | MathVista |
|---|---|---|
| GPT-4.1 (baseline) | 74.00 | 70.40 |
| GPT-4.1 (step) | 74.16 | 74.27 |
| GPT-4.1 (w/o logic memory) | 76.64 | 75.59 |
| GPT-4.1 (w/o visual memory) | 76.88 | 75.66 |
| GPT-4.1 (+ ViLoMem) | 77.26 | **76.88** |
| GPT-4.1 (+ ViLoMem & attention) | **78.21** | 76.87 |

verification, enhancing scoring accuracy and reducing false negatives from format variations.

## 4.2. Main Results on Multimodal Benchmarks

Table 1 shows evaluations across six multimodal benchmarks covering mathematical reasoning, hallucination robustness, and visual knowledge understanding. We compare three MLLMs under three configurations: *Baseline*: following the official default prompting setup; *Step*: using explicit step-by-step reasoning prompts; and *+ViLoMem*: integrating our dual-stream memory framework. The comparison between *Step* and *+ViLoMem* highlights the effectiveness of memory in mitigating *de novo* reasoning and promoting experience-driven problem solving.

**ViLoMem** achieves consistent improvements across all models, with particularly notable gains on mathematical reasoning benchmarks. This result aligns with our motivation, as mathematical reasoning tasks demand more visually grounded chains of thoughts. Prior studies have shown that visual perception errors significantly degrade reasoning accuracy [20, 45]. By tracking visual errors and integrating them with logical reasoning, **ViLoMem** effectively enhances overall mathematical reasoning performance. Among the three MLLMs, GPT-4.1 shows the largest improvement—particularly on MathVision (+6.48) and MathVista (+2.61)—owing to its stronger contextual learning ability and superior capacity to utilize and inter-
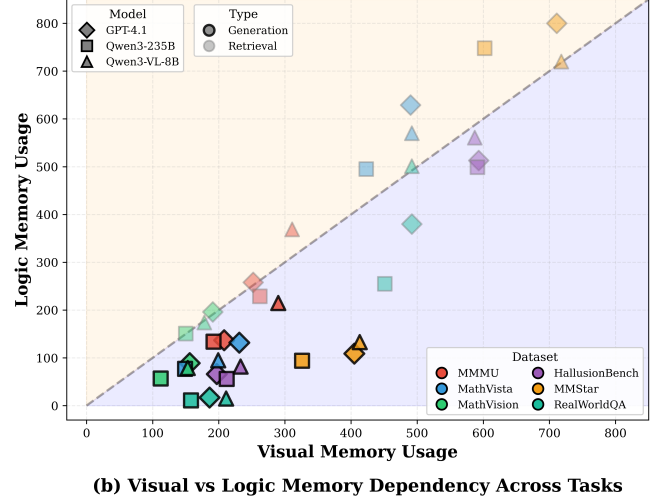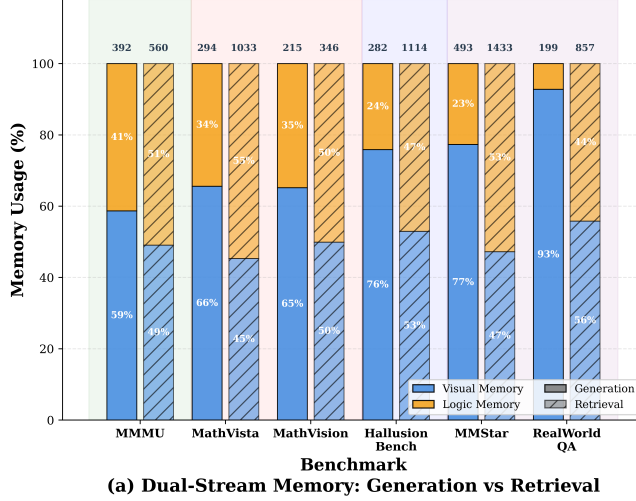
Figure 4. Analysis of dual stream memory usage patterns across six benchmarks. (a) Memory generation and retrieval statistics show that visual errors dominate generation (59% to 93%), while retrieval operations significantly exceed generation events. (b) Cross task dependency analysis reveals balanced utilization of both memory streams during retrieval across diverse tasks and models.

pret past errors for solving similar problems. Smaller models benefit more substantially from memory augmentation: Qwen3-VL-8B-Instruct achieves notable gains on MMMU (+4.38) and RealWorldQA (+2.74), indicating that structured memory provides complementary knowledge beyond the model's limited parametric capacity.

Among the evaluated tasks, improvements on knowledge-intensive benchmarks are moderate, as these tasks primarily rely on factual recall rather than multi-step reasoning. Moreover, manual inspection of the stored memory information from both streams reveals two primary performance bottlenecks. First, when the solver exhibits strong textual bias, over-relying on linguistic reasoning while paying limited attention to visual cues, the resulting reasoning traces contain insufficient visual information for the verifier to generate effective visual memory. Second, when the solver struggles to perceive complex diagrams and generates low-quality visual descriptions, the verifier finds it difficult to identify clear visual errors and tends to attribute all errors to the logical stream, often resulting in mixed memory updates. Therefore, a promising direction for future work is to design more specialized mechanisms to further enhance the decoupling the dual memory streams.

### 4.3. Case Study

Figure 3 illustrates the operation of dual-stream memory in practice. Cases 1, 2, and 4 expose a key limitation of logical memory: it retrieves guidelines irrelevant to the visual context (e.g., recalling perpendicular bisector principles when material discrimination is required). Visual memory effectively addresses this gap by identifying surface reflectivity (Case 1), numerical digits in diagrams (Case 2), and background luminance for color perception (Case 4). The at-

Table 3. Cross model memory transfer analysis. For each solver, we replace its self-generated memory with memories generated by the other two models on the same benchmark.

| Method | MMMU | MathVista |
|---|---|---|
| GPT-4.1 (step) | 74.16 | 74.27 |
| GPT-4.1 (+ **ViLoMem**) | 77.26 | **76.88** |
| GPT-4.1 (+ **ViLoMem** Cross) | **78.21** | 76.58 |
| Qwen3-VL-235B (step) | 75.97 | 83.66 |
| Qwen3-VL-235B (+ **ViLoMem**) | **79.40** | **84.98** |
| Qwen3-VL-235B (+ **ViLoMem** Cross) | 79.26 | 84.21 |
| Qwen3-VL-8B (step) | 65.52 | 77.80 |
| Qwen3-VL-8B (+ **ViLoMem**) | 69.90 | 77.87 |
| Qwen3-VL-8B (+ **ViLoMem** Cross) | **71.26** | **79.20** |

tention maps confirm that retrieved visual cues guide the model toward task-relevant regions (Case 2/ 4). Case 3 highlights the plausibility of our memory generation process: when a problem can be solved without visual cues (the question already providing complete visual descriptions), logical memory alone suffices. Overall, visual memory supports perception-intensive tasks, while logical memory governs reasoning-driven problems.

### 4.4. Ablation Study

We validate the necessity of dual-stream memory by selectively disabling each component on GPT-4.1. As shown in Table 2, removing either stream consistently degrades performance, confirming that both memory types are essential. Removing logical memory leads to larger drops on *Math-Vista*, where systematic reasoning and formula-related errors frequently recur. In contrast, removing visual memory produces comparable degradation across both benchmarks, indicating that visual distraction errors are pervasive in mul-

Table 4. Cross benchmark memory generalization analysis. For each benchmark, we exclude its task specific memory and merge memories from all other benchmarks as the retrieval source. Results demonstrate that while cross domain memories provide partial benefits, task aligned memories remain essential for optimal performance.

| Method | MMMU | MathVista | MathVision | HallusionBench | MMStar | RealWorldQA |
|---|---|---|---|---|---|---|
| Qwen3-VL-8B (baseline) | 66.38* | 77.20 | 48.13* | 61.10 | 70.91 | 71.50 |
| Qwen3-VL-8B (step) | 65.52 | 77.80 | 48.35 | 73.08 | 70.22 | 70.85 |
| Qwen3-VL-8B (+ **ViLoMem**) | **69.90** | **77.87** | 49.34 | **73.19** | **72.13** | **73.59** |
| Qwen3-VL-8B (+ **ViLoMem** Cross) | 65.14 | 76.10 | **50.00** | 70.66 | 70.93 | 71.63 |

timodal reasoning tasks. The gap between the single-stream variants and the full **ViLoMem** model demonstrates their complementarity: the visual and logical streams capture distinct, rather than redundant, error patterns. Augmenting visual memory with question-aware attention maps yields notable gains on MMMU, but only marginal improvements on MathVista, because diagram-based tasks require more fine-grained visual understanding, e.g., smaller-scale vertex attention and higher spatial precision. More detailed analyses are provided in the Appendix.

### 4.5. Memory Usage Analysis

Figure 4 analyzes memory usage patterns across all benchmarks. Visual memory generation dominates the error collection, accounting for 59%–93% of stored cases in Figure 4(a), demonstrating that visual perception remains the primary bottleneck in multimodal reasoning. Despite this generation asymmetry, both streams contribute comparably during retrieval, indicating effective memory reuse. Figure 4(b) further confirms consistent dual-stream coordination across all three MLLMs, as reflected by the distribution of translucent retrieval points along the diagonal, indicating balanced contributions from both visual and logical streams. Moreover, our memory mechanism is not biased toward any specific model, as all three models exhibit similar patterns of memory utilization.

### 4.6. Cross-Model Memory Transfer

To evaluate the reusability and composability of the dual-stream memory framework, we conduct cross-model memory transfer experiments where each solver retrieves memories generated by other models. As shown in Table 3, the 8B model benefits most from cross-model memories (+1.36 on MMMU, +1.33 on MathVista), surpassing its self-generated performance, indicating that memories distilled from stronger models encode higher-quality error patterns and generalization strategies. In contrast, larger models show comparable or slightly reduced performance, as their reasoning capabilities already yield near-optimal memory formation. These results highlight that dual-stream memory supports effective knowledge distillation from stronger to weaker models, enabling collaborative learning without explicit fine-tuning or ensembling.

### 4.7. Cross-Benchmark Memory Generalization

We assess memory transferability across task domains using Qwen3-VL-8B-Instruct. For each target benchmark in Table 4, we exclude its task-specific memory bank and instead retrieve from memories accumulated across *all other benchmarks*. The results reveal substantial heterogeneity: MathVision and RealWorldQA benefit from cross-domain memories, as both require strong spatial reasoning. In contrast, tasks with large domain gaps, such as MathVista and HallusionBench (diagram-grounded vs. natural image reasoning), exhibit conflicts in memory utilization. Overall, the persistent gap between cross-domain and **ViLoMem** underscores that task-aligned memories are essential for optimal performance, validating our design choice to maintain distinct memory banks for different domains.

## 5. Conclusion

We introduce **ViLoMem**, a dual-stream memory framework that separately models visual distraction patterns and logical hallucination errors for multimodal large language models. Inspired by human semantic memory systems, **ViLoMem** coordinates visual and logical memory streams through specialized retrieval strategies and grow-and-refine update mechanisms. Comprehensive evaluations across six multimodal benchmarks demonstrate consistent improvements, with particularly pronounced gains on mathematical reasoning tasks where visual-logical coupling is most acute. Ablation studies confirm that both memory streams are complementary; joint operation enables synergistic error correction. Further analyses reveal heterogeneous cross-domain transfer behavior—task-aligned domains benefit from shared memory, whereas domain-mismatched tasks exhibit mild interference. Moreover, cross-model transfer experiments highlight that our memory can distill error patterns and reasoning strategies from stronger models to smaller ones, demonstrating its potential as a lightweight knowledge-sharing mechanism without explicit fine-tuning. By enabling progressive error reduction without catastrophic forgetting, **ViLoMem** builds a foundation for continual learning in multimodal reasoning.

# References

[1] Rishabh Agarwal, Avi Singh, Lei Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, et al. Many-shot in-context learning. *Advances in Neural Information Processing Systems*, 37:76930–76966, 2024. 2

[2] Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziems, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, et al. Gepa: Reflective prompt evolution can outperform reinforcement learning. *arXiv preprint arXiv:2507.19457*, 2025. 2

[3] Keivan Alizadeh, Seyed Iman Mirzadeh, Dmitry Belenko, S Khatamifard, Minsik Cho, Carlo C Del Mundo, Mohammad Rastegari, and Mehrdad Farajtabar. Llm in a flash: Efficient large language model inference with limited memory. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12562–12584, 2024. 3

[4] Ali Ayub, Chrystopher L Nehaniv, and Kerstin Dautenhahn. Interactive continual learning architecture for long-term personalization of home service robots. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11289–11296. IEEE, 2024. 3

[5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1

[6] Yuxuan Cai, Yipeng Hao, Jie Zhou, et al. Building self-evolving agents via experience-driven lifelong learning: A framework and benchmark. *arXiv preprint arXiv:2508.19005*, 2025. 3

[7] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. MMStar: Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024. 5

[8] Qizhou Chen, Taolin Zhang, Xiaofeng He, Dongyang Li, Chengyu Wang, Longtao Huang, et al. Lifelong knowledge editing for llms with retrieval-augmented continuous prompt learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13565–13580, 2024. 2

[9] Alex Clarke and Lorraine K Tyler. Object-specific semantic coding in human perirhinal cortex. *Journal of Neuroscience*, 34(14):4766–4775, 2014. 3

[10] 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

[11] 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. Springer, 2024. 3

[12] Jizhan Fang, Xinle Deng, Haoming Xu, Ziyan Jiang, Yuqi Tang, Ziwen Xu, Shumin Deng, Yunzhi Yao, Mengru Wang, Shuofei Qiao, et al. Lightmem: Lightweight and efficient memory-augmented generation. *arXiv preprint arXiv:2510.18866*, 2025. 3

[13] Jinyuan Fang, Yanwen Peng, Xi Zhang, Yingxu Wang, Xinhao Yi, Guibin Zhang, Yi Xu, Bin Wu, Siwei Liu, Zihao Li, et al. A comprehensive survey of self-evolving ai agents: A new paradigm bridging foundation models and lifelong agentic systems. *arXiv preprint arXiv:2508.07407*, 2025. 1

[14] Huan-ang Gao, Jiayi Geng, Wenyue Hua, Mengkang Hu, Xinzhe Juan, Hongzhang Liu, Shilong Liu, Jiahao Qiu, Xuan Qi, Yiran Wu, et al. A survey of self-evolving agents: On path to artificial super intelligence. *arXiv preprint arXiv:2507.21046*, 2025. 1

[15] Shuzheng Gao, Chaozheng Wang, Cuiyun Gao, Xiaoqian Jiao, Chun Yong Chong, Shan Gao, and Michael Lyu. The prompt alchemist: Automated llm-tailored prompt optimization for test case generation. *arXiv preprint arXiv:2501.01329*, 2025. 2

[16] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. HallusionBench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*, 2023. 5

[17] Philipp Kuhnke, Curtiss A Chapman, Vincent KM Cheung, Sabrina Turker, Astrid Graessner, Sandra Martin, Kathleen A Williams, and Gesa Hartwigsen. The role of the angular gyrus in semantic cognition: a synthesis of five functional neuroimaging studies. *Brain Structure and Function*, 228 (1):273–291, 2023. 3

[18] Matthew A Lambon Ralph, Karen Sage, Roy W Jones, and Emily J Mayberry. Coherent concepts are computed in the anterior temporal lobes. *Proceedings of the National Academy of Sciences*, 107(6):2717–2722, 2010. 3

[19] Michael M Lombardo and Robert W Eichinger. *The Career Architect Development Planner*. Lominger, Minneapolis, 1st edition, 1996. 3

[20] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 1, 6

[21] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 5

[22] OpenCompass Contributors. VLMEvalKit: Open-source evaluation toolkit for large vision-language models, 2024. 6, 3

[23] César Santos, Fumio Machida, and Ermeson Andrade. Experimental investigation of memory-related software aging in llm systems. *Journal of Systems and Software*, page 112653, 2025. 3

[24] Rulin Shao, Jacqueline He, Akari Asai, Weijia Shi, Tim Dettmers, Sewon Min, Luke Zettlemoyer, and Pang Wei W

Koh. Scaling retrieval-based language models with a trillion-token datastore. *Advances in Neural Information Processing Systems*, 37:91260–91299, 2024. 2

[25] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023. 2

[26] Hao Sun et al. MathGlance: A benchmark for math at a glance understanding. *arXiv preprint*, 2025. Placeholder - shows visual perception bottleneck in mathematical reasoning; update with full citation when available. 1

[27] Mirac Suzgun, Mert Yuksekgonul, Federico Bianchi, Dan Jurafsky, and James Zou. Dynamic cheatsheet: Test-time learning with adaptive memory. *arXiv preprint arXiv:2504.07952*, 2025. 1, 3, 2

[28] Zhen Tan, Jun Yan, I-Hung Hsu, Rujun Han, Zifeng Wang, Long Le, Yiwen Song, Yanfei Chen, Hamid Palangi, George Lee, et al. In prospect and retrospect: Reflective memory management for long-term personalized dialogue agents. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8416–8439, 2025. 1

[29] V Team, Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, et al. Glm-4.5v and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. *arXiv preprint arXiv:2507.01006*, 2025. 1

[30] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal LLMs. *arXiv preprint arXiv:2401.06209*, 2024. Identifies nine visual patterns where MLLMs systematically fail; shows CLIP encoder limitations cascade to reasoning failures. 1

[31] Boshi Wang, Weijian Xu, Yunsheng Li, Mei Gao, Yujia Xie, Huan Sun, and Dongdong Chen. Improving code localization with repository memory. *arXiv preprint arXiv:2510.01003*, 2025. 3

[32] Fei Wang, Xingchen Wan, Ruoxi Sun, Jiefeng Chen, and Sercan O Arik. Astute rag: Overcoming imperfect retrieval augmentation and knowledge conflicts for large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 30553–30571, 2025. 2

[33] Ke Wang, Junting Ren, Weikang Yuan, Sicong Wang, Zihao Yang, Wentao Ma, and Wanli Ouyang. MATH-Vision: A challenging mathematical reasoning benchmark requiring visual understanding. *arXiv preprint arXiv:2409.13925*, 2024. 5

[34] Xiaohan Wang et al. MEG evidence that modality-independent conceptual representations contain semantic and visual features. *Journal of Neuroscience*, 44(28), 2024. Evidence that ATL acts as hub integrating sensory-motor features into coherent conceptual representations. 2

[35] Rong Wu, Xiaoman Wang, Jianbiao Mei, Pinlong Cai, Daocheng Fu, Cheng Yang, Licheng Wen, Xuemeng Yang, Yufan Shen, Yuxin Wang, et al. Evolver: Self-evolving llm agents through an experience-driven lifecycle. *arXiv preprint arXiv:2510.16079*, 2025. 3

[36] Shengqiong Wu, Hao Fei, Liangming Pan, William Yang Wang, Shuicheng Yan, and Tat-Seng Chua. Combating multimodal LLM hallucination via bottom-up holistic reasoning. *arXiv preprint arXiv:2412.11124*, 2024. Shows insufficient visual comprehension causes hallucinations; identifies object, attribute, and relationship perception errors. 1

[37] Yingsheng Wu, Yuxuan Gu, Xiaocheng Feng, Weihong Zhong, Dongliang Xu, Qing Yang, Hongtao Liu, and Bing Qin. Extending context window of large language models from a distributional perspective. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7288–7301, 2024. 2

[38] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024. 1

[39] xAI Team. RealWorldQA: A benchmark for real-world visual understanding, 2024. Real-world spatial understanding benchmark with 765 images. 5

[40] Wujiang Xu, Kai Mei, Hang Gao, Juntao Tan, Zujie Liang, and Yongfeng Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025. 3

[41] Yutao Yang, Jie Zhou, Junsong Li, Qianjun Pan, Bihao Zhan, Qin Chen, Xipeng Qiu, and Liang He. Reinforced interactive continual learning via real-time noisy human feedback. *arXiv preprint arXiv:2505.09925*, 2025. 3

[42] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations*, 2023. 2

[43] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023. 5

[44] Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. Textgrad: Automatic "differentiation" via text. *arXiv preprint arXiv:2406.07496*, 2024. 2

[45] Author Names Zhang. Primitive visual perception for multimodal reasoning. *arXiv preprint*, 2025. Placeholder - shows 72-78% of math reasoning failures stem from perception errors exceeding logic errors; update with full citation. 1, 6

[46] Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin, Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li, Xiangyuan Xue, Yijiang Li, et al. The landscape of agentic reinforcement learning for llms: A survey. *arXiv preprint arXiv:2509.02547*, 2025. 3

[47] Jiarui Zhang, Mahyar Khayatkhoei, Prateek Chhikara, and Filip Ilievski. Mllms know where to look: Training-free perception of small visual details with multimodal llms. In *The Thirteenth International Conference on Learning Representations*, 2025. 2

[48] Kai Zhang, Xiangchao Chen, Bo Liu, Tianci Xue, Zeyi Liao, Zhihan Liu, Xiyao Wang, Yuting Ning, Zhaorun Chen, Xiaohan Fu, et al. Agent learning via early experience. *arXiv preprint arXiv:2510.08558*, 2025. 3

[49] Qizheng Zhang, Changran Hu, Shubhangi Upasani, Boyuan Ma, Fenglu Hong, Vamsidhar Kamanuru, Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen Li, et al. Agentic context engineering: Evolving contexts for self-improving language models. *arXiv preprint arXiv:2510.04618*, 2025. 1, 3

[50] Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu, Min Zhang, Wen Zhang, and Huajun Chen. Abstractive visual understanding of multi-modal structured knowledge: A new perspective for mllm evaluation. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 12323–12332, 2025. 1

[51] Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model-based agents. *ACM Transactions on Information Systems*, 43 (6):1–47, 2025. 1, 3

[52] Zijia Zhao, Yuqi Huo, Tongtian Yue, Longteng Guo, Haoyu Lu, Bingning Wang, Weipeng Chen, and Jing Liu. Efficient motion-aware video mllm. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24159–24168, 2025. 1

[53] Chenyue Zhou, Mingxuan Wang, Yanbiao Ma, Chenxu Wu, Wanyi Chen, et al. From perception to cognition: A survey of vision-language interactive reasoning in multimodal large language models. *arXiv preprint arXiv:2509.25373*, 2025. Comprehensive survey on perception-cognition disconnect; shows static visual processing causes decoupling between answers and visual facts. 1

[54] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1

# Agentic Learner with Grow-and-Refine Multimodal Semantic Memory

## Supplementary Material

## 6. Additional Results and Ablation Study

### 6.1. Integration with more models

To verify the flexibility of **ViLoMem**, we extend our evaluation beyond the main experiments to recent reasoning-enhanced models, including GLM-4.1v [29], InternVL3-38B [54], and Gemini 2.5 [10]. As shown in Table 5, **ViLoMem** demonstrates robust adaptability across different architecture designs and inference regimes, consistently improving performance over both baseline and step-by-step configurations. This pattern echoes our observations in the main paper that visual perception remains a dominant bottleneck for multimodal reasoning [20, 45] and that decoupling visual distraction from logical hallucination yields complementary gains across tasks. Notably, models equipped with "thinking" or long-chain reasoning capabilities exhibit superior compatibility with the step-by-step format required for memory retrieval: their extended inference process allows for tighter integration of retrieved visual and logical guidelines into the reasoning chain, enabling them to correct potential errors before they propagate. These results suggest that **ViLoMem** is particularly well-suited to models with strong deliberative reasoning, while still offering consistent benefits to smaller or less capable solvers.

### 6.2. Attention Mechanism Ablation

Table 6 presents the ablation study of the attention mechanism. In general, the integration of attention maps yields consistent performance gains across hallucination and general reasoning benchmarks (e.g., HallusionBench, MM-Star), corroborating the critical importance of visual memory in refining perceptual grounding. However, we observe a performance plateau or marginal decline on mathematics-centric datasets (MathVista and MathVision). We attribute this limitation to two primary factors: (1) **Visualization Precision**: Current attention visualization methods struggle to faithfully preserve fine-grained geometric structures and chart details, which are essential for mathematical reasoning. (2) **Contextual Interpretation**: While serving as an auxiliary image to enhance visual context, the attention map imposes higher demands on the model's intrinsic capability to interpret heatmap overlays. The benefit of this enriched context is contingent on the model's ability to align these explicit visual cues with the raw image features without information loss.

### 6.3. Additional Case Study

Figure 5 summarizes representative qualitative cases. For many vision-intensive questions (e.g., traffic-light color,

visible portion of the sun, object localization, and optical-illusion setups), logical memory is either not retrieved or fails to offer useful guidance, while visual memory provides concrete viewing strategies such as checking the actual illuminated region, reading tiny objects and relative positions from the viewer's frame, or isolating targets from distracting backgrounds. In these cases, attention maps concentrate on the queried regions (e.g., the active light, visible solar arc, or relevant segments), so that the retrieved visual guidelines directly steer the solver toward task-relevant evidence.

For geometry and chart-reading tasks, visual and logical memories are complementary: logical memory provides reusable rules for measurement and graph interpretation, while visual memory focuses on concrete inspection behaviors such as aligning with gridlines, following step edges, or checking true line orientation under strong illusions. Together, these cases highlight a clear division of labor: visual memory governs "where to look" and mitigates systematic perceptual traps, whereas logical memory refines "how to reason" once the correct visual evidence has been attended.

### 6.4. Comparison with Existing Memory Methods

We benchmark **ViLoMem** against state-of-the-art memory mechanisms [27, 49]. While the original Dynamic-Cheetsheet [27] employs cumulative memory, its unbounded context growth is infeasible for our large-scale setting (approx. 1,000 cases per benchmark), so we adopt the retrieval-based configuration from the open-source Dynamic-Cheetsheet codebase, which follows the similar methodology as ACE [49]. For a fair multimodal comparison, we replicate the official prompt structure and use the same MLLM for both memory generation and inference. In this setup, the retrieval module relies purely on text similarity without image-aware matching.

Experimental results in Table 6 show that this direct adaptation of logical memory methods is suboptimal in multimodal settings and can even underperform the baseline, especially for smaller models. In practice, such text-only retrieval often surfaces visually dissimilar examples with similar questions, resurfacing prior misperceptions as salient "hints" that misdirect attention away from the correct regions of the current problem. Qualitative inspection further reveals that Dynamic-Cheetsheet and ACE are tailored to code- or logic-centric schemas: even when driven by an MLLM, they mainly produce fine-grained corrections of specific visual details (digits, colors, marks) rather than robust guidance on how to inspect diagrams. These detail-level cues lack stable visual grounding and easily conflict

Table 5. Additional evaluation results on GLM4.1v, InternVL3-38B, and Gemini2.5-flash across six multimodal reasoning benchmarks. Models with "(step)" and "(+ **ViLoMem**)" are prompted by step-by-step reasoning. Results demonstrate consistent improvements from **ViLoMem** across diverse model architectures.

| Method | MMMU (dev) | MathVista (mini) | MathVision (mini) | HallusionBench | MMStar | RealWorldQA |
|---|---|---|---|---|---|---|
| GLM4.1v (baseline) | 69.14 | 72.57 | 56.88 | 73.08 | 72.90 | 73.33 |
| GLM4.1v (step) | 70.29 | 73.47 | 58.22 | 72.77 | 73.40 | 72.54 |
| GLM4.1v (+ **ViLoMem**) | **71.52** | **73.97** | **61.51** | **74.02** | **73.47** | <u>72.68</u> |
| InternVL3-38B (baseline) | 62.92 | 70.80 | 35.53 | 67.40 | 69.33 | 71.99 |
| InternVL3-38B (step) | 64.18 | 71.90 | 35.56 | 71.50 | 67.80 | 72.42 |
| InternVL3-38B (+ **ViLoMem**) | **65.97** | **73.80** | **36.84** | **72.34** | **69.73** | **73.20** |
| Gemini2.5-flash (baseline) | 72.18 | 81.10 | 53.21 | 72.67 | 72.07 | 76.99 |
| Gemini2.5-flash (step) | 71.90 | 81.41 | 53.94 | 76.34 | 72.40 | 71.50 |
| Gemini2.5-flash (+ **ViLoMem**) | **72.86** | **83.40** | **58.22** | **78.33** | **73.20** | <u>76.42</u> |

Table 6. Comprehensive ablation study and comparison with existing memory methods across six multimodal reasoning benchmarks. We compare **ViLoMem** with attention mechanism variants and the Dynamic-Cheetsheet [27] baseline adapted for multimodal tasks.

| Method | MMMU (dev) | MathVista (mini) | MathVision (mini) | HallusionBench | MMStar | RealWorldQA |
|---|---|---|---|---|---|---|
| *GPT-4.1* | | | | | | |
| baseline | 74.00 | 70.40 | 46.12* | 58.50 | 69.80 | 73.72 |
| step | 74.16 | 74.27 | 47.47 | 74.44 | 70.43 | 72.03 |
| + dynamic-cheetsheet | 70.95 | 73.87 | 48.68 | 75.30 | 68.68 | 70.13 |
| + **ViLoMem** | 77.26 | **76.88** | **53.95** | 75.29 | **72.43** | **74.38** |
| + **ViLoMem** & attention | **78.21** | 76.87 | 50.66 | **75.73** | 71.76 | **74.38** |
| *Qwen3-VL-235B-A22B-Instruct* | | | | | | |
| baseline | 78.70 | 84.90 | 61.28* | 63.20 | 78.40 | 79.30 |
| step | 75.97 | 83.66 | 62.17 | 74.58 | 76.16 | 78.66 |
| + dynamic-cheetsheet | 72.13 | 83.25 | 60.06 | 70.62 | 75.49 | 77.11 |
| + **ViLoMem** | **79.40** | **84.98** | **62.83** | 75.21 | 78.31 | 77.22 |
| + **ViLoMem** & attention | 78.14 | 83.87 | 60.86 | **75.95** | **78.46** | <u>77.88</u> |
| *Qwen3-VL-8B-Instruct* | | | | | | |
| baseline | 66.38* | 77.20 | 48.13* | 61.10 | 70.91 | 71.50 |
| step | 65.52 | 77.80 | 48.35 | 73.08 | 70.22 | 70.85 |
| + dynamic-cheetsheet | 63.39 | 74.92 | 46.81 | 68.39 | 69.12 | 69.98 |
| + **ViLoMem** | **69.90** | **77.87** | **49.34** | 73.19 | **72.13** | **73.59** |
| + **ViLoMem** & attention | 67.52 | 77.07 | 48.72 | **74.87** | 72.67 | 73.46 |

with the actual image, inducing additional hallucinations that smaller models are particularly vulnerable to. This contrast highlights the need for **ViLoMem**'s decoupled visual stream and question-aware retrieval, which explicitly organize and retrieve perception-oriented error patterns instead of repurposing logic-only memories.

# 7. Additional Experimental Details

This section provides additional implementation details that complement the experimental setup.

**Model Deployment.** For open-source models, we deploy most checkpoints using `vLLM` for efficient batched inference. Due to its scale, *Qwen3-VL-235B-A22B-Instruct* is accessed via its official API instead of local deployment, and all proprietary models (e.g., GPT-4.1, Gemini 2.5 flash) are evaluated through their corresponding APIs. For

API-based evaluations, certain images or prompts may be flagged as unsafe by the provider's safety filters and thus rejected, which introduces a small amount of noise into the reported scores.

**Decoding Hyperparameters.** Unless otherwise specified, we use a temperature of $0.7$ and a maximum generation length of $8,192$ tokens for all models. Within our memory pipeline, the maximum generation length is set to $1,024$ tokens for problem analysis and $2,048$ tokens for memory generation to balance expressiveness and efficiency. Baseline evaluations directly feed benchmark questions to the models without additional prompts, whereas the *Step* configuration prepends a simple step-by-step system prompt; the full template is shown in Figure 6.

**Attention Map Generation.** Attention maps are generated following the training-free small-detail perception framework of Zhang et al. [47], instantiated with *Qwen2.5-VL-*

*3B* as the backbone model. This setup produces token-level saliency over input images, which we overlay as heatmaps to visualize and interpret visual memory retrieval.

**Evaluation Protocol.** We adopt VLMEvalKit [22] as the primary evaluation framework. When automatic matching fails or produces ambiguous results (e.g., due to formatting variations), we further apply *Math-Verify* and an LLM-as-a-judge protocol to reduce sensitivity to output formatting. The judge model is *Qwen3-8B-Instruct*, which assesses whether a model's response is semantically correct with respect to the reference answer.

## 8. Prompt Templates

We provide the full prompt templates used in our framework, including the step-by-step reasoning prompt used in the *Step* configuration (Figure 6), the Problem Analysis Prompt (Figure 7), the Logical Memory Generation Prompt (Figure 8), and the Visual Memory Generation Prompt (Figure 9), together with the LLM-as-a-judge verification prompt (Figure 10).

| Question | Logic Memory | Visual Memory | Attention Map |
|---|---|---|---|
| **CASE 1**<br>Is the traffic light green? | / | When counting illuminated traffic lights, **verify the color of each light individually and confirm it is part of a traffic signal assembly**, not a decorative or non-traffic light fixture. |  |
| **CASE 2**<br>What percent of the sun is showing? | When calculating the range of a data set, always double-check that the minimum and maximum values are correctly identified by reviewing all data points, especially when values repeat or are close in magnitude.... | When **identifying 'tiny thing' or relative positions**, always **verify object size and spatial layout from the viewer's perspective**; matte finish and right-side positioning must be visually confirmed, not assumed.... |  |
| **CASE 3**<br>Move the ruler to measure the length of the pencil to the nearest inch. The pencil is about (_) inches long. | When solving geometry problems involving squares inscribed in circles or angles formed by intersecting chords, verify key relationships such as the diagonal of the square equaling the circle's diameter and the angle measure being half the sum of the intercepted arcs..... | Always **verify the exact starting point of the object being measured on the ruler**; do not assume alignment with 0 cm unless visually confirmed, as the object may begin at a non-zero mark.... |  |
| **CASE 4**<br>Are blue lines in the image parallel? Yes or No | When applying the Corresponding Angles Postulate, always verify that the angles lie on the same side of the transversal and in matching positions at each intersection..... | When assessing parallelism of lines surrounded by diagonal patterns, use a ruler or grid overlay to verify true orientation and spacing, as the **background can induce a perceptual illusion of verticality or parallelism**.<br><br>When comparing line lengths in an image, **measure or visually align them directly rather than relying on assumptions about known optical illusions, as the actual lengths may differ from the illusion's typical setup....** |  |
| **CASE 5**<br>Where is the sheep? | / | When identifying objects relative to others in a scene, **always verify both the object's identity and its spatial position (left/right/front/back) by comparing their actual locations in the image**, not assumptions based on size or context. |  |
| **CASE 6**<br>What is the position of the sink relative to the refrigerator? | When a question asks for the direction of object A relative to object B, ensure the reference frame is correctly oriented: "**A in relation to B**" means you start from B and **determine where A lies, not the reverse**. Always double-check the subject and reference point in directional questions to avoid reversing the relationship. | When identifying objects relative to others in a scene, **always verify both the object's identity and its spatial position (left/right/front/back) by comparing their actual locations in the image**, not assumptions based on size or context. |  |
| **CASE 7**<br>What is the value of the smallest individual bar in the whole chart? | When interpreting values on a logarithmic scale, remember that tick marks represent powers of 10, and values between $10^2$ and $10^3$ range from 100 to 1000; always verify whether a bar exceeds a linear threshold like 100 by estimating its actual value, not just its position relative to $10^2$... | On a logarithmic scale, a bar ending exactly at a labeled tick (e.g., $10^2$) represents that exact value, not a value greater than it; **always check if the bar exceeds the tick mark to determine if it's strictly larger**.... |  |
| **CASE 8**<br>When does the function value first reach 2? | When matching a correctly identified element (e.g., a region, value, or object) to a multiple-choice option, always verify that the letter of the choice corresponds to the correct labeled item in the diagram or question, not just the reasoning outcome. Double-check the mapping between your conclusion and the answer options to avoid selection errors.... | When interpreting a step function graph, **always verify the exact y-value of each horizontal segment by aligning it with the y-axis gridlines**, rather than assuming values based on adjacent steps or integer patterns.<br><br>When identifying the y-intercept on a graph, **always trace the curve to where it crosses the y-axis (x=0) and read the exact y-value at that point**, not the vertex or any nearby grid line.... |  |
| **CASE 9**<br>Does this figure mainly depict a hen and eggs, no potatoes? | When a question asks whether two colors are different in the context of an optical illusion, always consider whether it is asking about objective color values (e.g., RGB or physical properties) rather than perceived appearance; remember that illusions affect perception, not necessarily reality. | When comparing the color of objects on a gradient background, **verify that perceived differences are not caused by the background's luminance shift by isolating or masking the background to assess the objects' true color**... |  |

Figure 5. Showcase of representative cases demonstrating **ViLoMem**'s memory generation and retrieval process across different types of multimodal reasoning tasks.

---

**Prompt: Step-by-Step Reasoning**

**Objective:** Solve the given problem using a step-by-step process.

**Expected Output Structure:**

```
Step 1:
Step 2:
...
Step n: Final Answer: \boxed{answer}

Question:
```

---

Figure 6. The step-by-step reasoning system prompt used in the *Step* configuration.

---

**Prompt: Problem Analysis**

**Objective:**
Analyze the following problem to identify its subject area and the key concepts, principles, formulas, or laws required for its solution. This analysis will be used to retrieve relevant guiding principles from a knowledge base.
**Instructions:**
- Do not solve the problem.
- First, identify the primary subject (e.g., Physics, Chemistry, Biology, Mathematics).
- Then, list the core concepts or principles involved (e.g., Newton's Second Law, Conservation of Energy, Stoichiometry, Pythagorean theorem).
- Keep the analysis concise and focused.
**Problem:**
{question}
**Output Format:**
Subject: <The primary subject>
Key Concepts: <A brief list of key concepts>

---

Figure 7. The prompt template for analyzing the problem to identify its subject and key concepts.

Figure 8. The prompt template for generating logical memories.

**Prompt: Visual Memory Generation**

**Objective:**
You are an expert in visual reasoning and error analysis. Your task is to first describe the provided image objectively, then analyze an incorrect reasoning process to determine if the error stems from misinterpreting that image. If a visual error is found, you must generate a concise, actionable guideline (a "visual memory") to prevent this mistake in the future.

**Context:**
- Problem: {question}
- Incorrect Reasoning Steps: {reasoning_steps}
- Correct Answer (for reference): {gold_answer}

**Attached Image:** <image>

**Thinking Process and Final Output:**
Your response must follow a strict two-stage process. The first stage is your internal "thought process" which you will write out. The second stage is the final JSON output.

**Stage 1: Internal Thought Process (Write this out first)**
1. **Describe the Image:** Begin by providing an objective, detailed description of the attached image. List all key elements, labels, values, geometric shapes, and their relationships. This description will serve as the "ground truth" for your analysis.
2. **Analyze for Discrepancies:** Compare your image description and the image itself against the text in `Incorrect Reasoning Steps`. Identify any contradictions, misinterpretations, or omissions.

**Stage 2: Final JSON Output (Provide ONLY this JSON block as the final answer)**
After completing your thought process, generate a JSON object based on your analysis. The JSON should adhere to the following structure and guidelines.

**Guidelines for `guideline` Generation:**
- The guideline MUST be about how to correctly interpret a specific visual pattern or element.
- It must be a rule that can be applied to other, similar-looking problems.
- It should be concise (one to two sentences).

**Guideline Examples (Good, Specific Visual Memories):**
- **(Physics/Diagrams):** "In a free-body diagram, always verify that all forces, including friction and normal force, are accounted for before applying Newton's laws."
- **(Geometry):** "When an angle appears to be a right angle in a diagram, do not assume it is 90 degrees unless it is explicitly marked with a square symbol."
- **(Chemistry/Molecules):** "For complex organic molecules, double-check the placement of double bonds and functional groups as they dictate the molecule's reactivity."
- **(Biology/Graphs):** "When reading a bar chart, pay close attention to the Y-axis scale and units to avoid misinterpreting the magnitude of the results."

**Avoid these types of guidelines (Bad, Non-Visual or Too Vague):**
- "The model made a calculation error." (This is a logical error, not visual)
- "You need to look at the image more carefully." (Not actionable)
- "The reasoning about the physics was wrong." (Too general)

**Final Output Format (use this exact JSON structure):**

```
{
    "is_visual_error": true/false,
    "analysis_summary": "A brief, one-sentence summary of the visual
                         misinterpretation.",
    "guideline": "Your 1-2 sentence visual guideline. Provide this
                  only if is_visual_error is true, otherwise it
                  should be null."
}
```

Figure 9. The prompt template for generating visual memories.

**Prompt: LLM-as-a-Judge Verification**

**Objective:**
You are an expert answer verification judge. Your task is to determine whether a model prediction matches the gold answer.

**Core Principle (Critical Rule):**
- All decisions are based *only* on the gold answer; ignore the quality of the reasoning.
- If the extracted final answer from the prediction exactly matches the gold answer, set `verified=true`; otherwise, set `verified=false`.
- Do not consider whether the prediction's reasoning is correct or sensible.
- Do not give partial credit for "close" answers (e.g., $2 \neq 9$, $C \neq A$).

**Verification Steps:**
1. **Identify the gold answer format.** Determine whether the gold answer is:
   - a single letter (`A/B/C/D/E`) for multiple-choice questions (compare letters only);
   - a number for numerical questions (compare numeric values, ignoring formatting such as 7 vs. 7.0);
   - a text span for open-ended questions (compare semantic meaning, allowing minor wording differences).
2. **Extract the final answer from the prediction.**
   - For multiple-choice questions, locate the final chosen letter (often after "Final Answer:" or "Answer:") and compare it with the gold letter.
   - For numerical questions, locate the final numeric value, ignoring units and extra text, then compare it to the gold number.
   - For text answers, extract the final answer phrase and compare its semantic meaning with the gold text (e.g., "Yes, the baby is crawling to the right." matches "Yes").
3. **Apply the strict matching rule.**
   - Only compare the final extracted answers.
   - Do not use external knowledge to judge whether an answer is reasonable.
   - If the extracted answer and the gold answer match under the appropriate format, output `verified=true`; otherwise, output `verified=false`.

**Input Fields:**
- Question: {question}
- Gold Answer: {gold_answer}
- Choices (optional, for multiple choice): {choices_text}
- Prediction: {prediction}

**Output Format (JSON, exact structure):**

```
{
  "reasoning": "Step 1: Extract answer from prediction: [extracted_value].
Step 2: Compare with gold: [gold_value].
Step 3: Match result: [yes/no].",
  "verified": true or false
}
```

Figure 10. The LLM-as-a-judge prompt template used to verify whether a model prediction matches the gold answer, independent of reasoning quality.